

# Maths for Intelligent Systems

Marc Toussaint

January, 2017

## Contents

<b>1</b>	<b>Speaking Maths</b>	<b>4</b>
1.1	Describing systems . . . . .	4
1.2	Should maths be taught with many application examples? Or abstractly? . . .	4
1.3	Notation: Some seeming trivialities . . . . .	5
<b>2</b>	<b>Linear Algebra</b>	<b>6</b>
2.1	Vector Spaces . . . . .	6
	Why should we care for vector spaces in intelligent systems research? (6) What is a vector? (7) What is a vector space? (7)	
2.2	Vectors, dual vectors, coordinates, matrices, tensors . . . . .	8
	A taxonomy of linear functions (8) Bases and coordinates (9) The dual vector space – and its coordinates (9) Coordinates for every linear thing: tensors (10) Finally: Matrices (10)	
2.3	Scalar product and orthonormal basis . . . . .	12
	Properties of orthonormal bases (13)	
2.4	The Structure of Transforms & Singular Value Decomposition . . . . .	14
	The Singular Value Decomposition Theorem (14)	
2.5	Point of departure from the coordinate-free notation . . . . .	16
2.6	Filling SVD with life . . . . .	17
	Understand $vv^T$ as projection (17) SVD for symmetric matrices (18) SVD for general matrices (19)	
2.7	Eigendecomposition . . . . .	20
	Power Method (20) Power Method including the smallest eigenvalue (20) Why should I care about Eigenvalues and Eigenvectors? (21)	

2.8	Point of Departure: Numerics to compute these things . . . . .	21
2.9	Examples and Exercises . . . . .	22
	Projection (22) Matrix equations (22) PCA (22) Exercise: The Linear Robot (Nathan's) (22)	
<b>3</b>	<b>Derivatives</b>	<b>22</b>
3.1	The coordinate-free view: The meaning of life for a derivative is to eat a vector and output a scalar . . . . .	22
3.2	Don't confuse partial and total derivative . . . . .	23
	Partial derivative (23) Chain Rule, Function Networks and the total derivative (23)	
3.3	"Gradient vectors", Co- and contra-variance, Jacobian, Hessian . . . . .	26
	Gradient, Jacobian, and Hessian (26) Why the partial derivative coordinates are co-variant (27)	
3.4	Derivative Rules and Examples . . . . .	29
	GP regression (30) Logistic regression (31)	
3.5	Taylor expansion . . . . .	31
3.6	Steepest descent and the covariant gradient vector . . . . .	32
3.7	Check your gradients numerically! . . . . .	32
3.8	Examples and Exercises . . . . .	33
	Vector derivatives (33) More derivatives (33) Minima (34) Finite Difference Gradient Checking (34) "Backprop" in a Neural Net (34) Logistic Regression Gradient & Hessian (35)	
<b>4</b>	<b>Optimization</b>	<b>35</b>
4.1	The general optimization problem – a mathematical program . . . . .	35
4.2	The KKT conditions . . . . .	36
4.3	Unconstrained problems to tackle a constrained problem . . . . .	37
	Log Barrier (39) Augmented Lagrangian* (39)	
4.4	The Lagrangian . . . . .	40
	How the Lagrangian relates to the KKT conditions (40) Solving mathematical programs analytically, on paper. (41) Solving the dual problem, instead of the primal. (42) Finding the "saddle point" directly with a primal-dual Newton method. (42) Properties of the dual problem* (43) Log barrier again (43)	
4.5	Downhill algorithms for unconstrained optimization . . . . .	44
	Why you shouldn't trust the magnitude of the gradient (44) Ensuring monotone and sufficient decrease: Backtracking line search, Wolfe conditions, & convergence (45) The Newton direction (46) Gauss-Newton: a super important special case (49) Quasi-Newton & BFGS: approximating the hessian from gradient observations (50) Conjugate Gradient* (51) Rprop* (52)	

4.6	Convex Problems . . . . .	53
	Convex sets, functions, problems (54)    Linear and quadratic programs (54)    The Simplex Algorithm (55)    Sequential Quadratic Programming (57)	
4.7	Blackbox & Global Optimization: It's all about learning . . . . .	57
	A sequential decision problem formulation (58)    Acquisition Functions for Bayesian Global Optimization* (59)    Classical model-based blackbox optimization (non-global)* (60)    Evolutionary Algorithms* (62)	
4.8	Examples and Exercises . . . . .	62
	Optimize a constrained problem (62)    Lagrangian and dual function (63)    Convergence proof (63)    Robust unconstrained optimization (64)    Network flow problem (65) Minimum fuel optimal control (65)    Reformulating Norms (65)    Solve one of these only, the others are optional (65)	
<b>5</b>	<b>Probabilities &amp; Information</b>	<b>65</b>
5.1	Basics . . . . .	66
	Axioms, definitions, Bayes rule (66)    Standard discrete distributions (68)    Conjugate distributions (68)    Distributions over continuous domain (68)    Gaussian (69)    "Particle distribution" (70)	
5.2	Between probabilities and optimization: neg-log-probabilities, exp-neg-energies, exponential family, Gibbs and Boltzmann . . . . .	72
5.3	Information, Entropie & Kullback-Leibler . . . . .	75
5.4	The Laplace approximation: A 2nd-order Taylor of $\log p$ . . . . .	76
5.5	Variational Inference . . . . .	77
5.6	The Fisher information metric: 2nd-order Taylor of the KLD . . . . .	77
5.7	Examples and Exercises . . . . .	77
	Maximum Entropy and Maximum Likelihood (77)    Maximum likelihood and KL- divergence (78)    Laplace Approximation (79)    Learning = Compression (79)    A gzip experiment (79)    Maximum Entropy and ML (80)	
<b>A</b>	<b>Gaussian identities</b>	<b>81</b>
<b>B</b>	<b>3D geometry basics (for robotics)</b>	<b>85</b>
B.1	Rotations . . . . .	85
B.2	Transformations . . . . .	89
	Static transformations (89)    Dynamic transformations (89)    A note on affine coor- dinate frames (91)	
B.3	Kinematic chains . . . . .	92
	Rigid and actuated transforms (92)    Jacobian & Hessian (93)	

# 1 Speaking Maths

## 1.1 Describing systems

Systems can be described in many ways. Biologists describe their systems often using text, and lots and lots of data. Architects describe buildings using drawings. Physics describe nature using differential equations, or optimality principles, or differential geometry and group theory. The whole point of science is to find descriptions of systems—in the natural science descriptions that allow prediction, in the synthetic/engineering sciences descriptions that enable the design of good systems, problem-solving systems.

And how should we describe intelligent systems? Robots, perception systems, machine learning systems? I guess there are two main categories: the imperative way in terms of literal algorithms (code), or the declarative way in terms of formulating the problem. I prefer the latter.

The point of this lecture is to teach you to *speak maths*, to use maths to describe systems or problems. I feel that most maths courses rather teach to consume maths, or solve mathematical problems, or prove things. Clearly, this is also important. But for the purpose of intelligent systems research, it is essential to be skilled in *expressing* problems mathematically, before even thinking about resulting algorithms.

If you happen to attend a Machine Learning or Robotics course you'll see that *every* problem is addressed the same way: You have an "intuitively formulated" problem; the first step is to find a mathematical formulation; the second step to solve it. The second step is often technical. The first step is really the interesting and creative part. This is where you have to nail down the problem, i.e., nail down what it means to be intelligent/good/well-performing, and thereby describe "intelligence"—or at least a tiny aspect of it.

The "Maths for Intelligent Systems" course will recap essentials of linear algebra, optimization, probabilities, and statistics. These fields are essential to formulate problems in intelligent systems research and hopefully will equip you with the basics of speaking maths.

## 1.2 Should maths be taught with many application examples? Or abstractly?

Maybe this is the wrong question and implies a view on maths I don't agree with. I think (but this is arguable) maths is nothing but abstractions of real-world things. At least I aim to teach maths as abstractions of real-world things. It is misleading to think that there is "pure maths" and then "applications". Instead mathematical concepts, such as a vector, are abstractions of real-world things, such as cars, faces, scenes, images, documents; and theorems, methods and algorithms that apply on

vectors of course also apply to all the real-world things—subject to the limitations of this abstraction. So, the goal is not to teach you a lookup table of which method can be used in which application, but rather to teach which concepts maths offers to abstract real-world things—so that you find such abstractions yourself once you'll have to solve a real-world problem.

### 1.3 Notation: Some seeming trivialities

Equations and mathematical expressions have a syntax. This is hardly ever made explicit<sup>1</sup> and might seem trivial. But it is surprising how buggy mathematical statements can be in scientific papers (and oral exams). I don't want to write much text about this, just some bullet points:

- Always declare objects.
- Be aware of variable/index scoping. For instance, if you have an equation, and one side includes a variable  $i$ , the other doesn't, this often is a notational bug. (Unless this equation incidently makes a statement about independence on  $i$ .)
- Type checking. Within an equation, be sure to know exactly of what type each term is: vector? matrix? scalar? tensor?
- Decorations are ok, but really not necessary. It is much more important to declare all things. E.g., there are all kinds of decorations used for vectors,  $\mathbf{v}$ ,  $\underline{v}$ ,  $\vec{v}$ ,  $|v\rangle$  and matrices. But these are not necessary. Properly declaring all symbols is much more important.
- sets:  $\{x_i\}_{i=1}^n$ , tuples:  $(x_i)_{i=1}^n, (x_1, \dots, x_n), x_{1:n}$
- set notations  $\{f(x) : x \in \mathbb{R}\}, \{n \in \mathbb{N} : \exists \{v_i\}_{i=1}^n \text{ linearly independent}, v_i \in V\}$
- indicator functions  $[a = b]$  or  $\mathbb{I}(a = b)$  and Kronecker  $\delta_{ab}$
- $(a, b) \in A \times B$ ,
- outer product  $x \otimes y \in X \times Y$  is almost the same as  $(x, y) \in X \times Y$  (slightly different in standard vector and tensor context); don't confuse with the cross product in 3D!
- direct sum  $A \oplus B$  is same as  $A \times B$  in finite-dimensional spaces
- write *argmin* instead of *argmin*
- Never use multiple letters for one thing. E.g.  $length = 3$  means  $l$  times  $e$  times  $n$  times  $g$  times  $t$  times  $h$  equals 3.

<sup>1</sup>Except perhaps by Gödel's incompleteness theorems and areas like automated theorem proving.

- Difference between  $\rightarrow$  and  $\mapsto$

$$f : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto \cos(x)$$

- The dot thing:

$$f : A \times B \rightarrow C, \quad f(a, \cdot) : B \rightarrow C, b \mapsto f(a, b)$$

The dot is in particular used for the inner product (a 2-form):  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$

- The  $p$ -norm  $\|x\|_p = [\sum_i x_i^p]^{1/p}$ . By default  $p = 2$ , that is  $\|x\| = \|x\|_2 = |x|$ , which is the length of a vector. When saying “2-norm”, I often actually mean  $\|x\|^2 = \sum_i x_i^2$ .

- $\text{diag}(a_1, \dots, a_n) = \begin{pmatrix} a_1 & & 0 \\ & \ddots & \\ 0 & & a_n \end{pmatrix}$

- A typical convention is

$$\mathbf{0}_n = (0, \dots, 0) \in \mathbb{R}^n, \quad \mathbf{1}_n = (1, \dots, 1) \in \mathbb{R}^n, \quad \mathbf{I}_n = \text{diag}(\mathbf{1}_n)$$

Also,  $e_i = (0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^n$  often denotes the  $i$ th row of the identity matrix, which of course are the coordinates of a basis vector  $e_i \in \mathcal{V}$  in a basis  $\{e_i\}_{i=1}^n$ .

- The element-wise product of two matrices  $A$  and  $B$  is also called Hadamard product and notated  $A \circ B$  (which has nothing to do with the concatenation of two operations!). If there is need, perhaps also use this to notate the element-wise product of two vectors.

## 2 Linear Algebra

### 2.1 Vector Spaces

#### 2.1.1 Why should we care for vector spaces in intelligent systems research?

We want to describe intelligent systems. For this we describe systems, or aspects of systems, as elements of a space:

- The input space  $X$ , output space  $Y$  in ML
- The space of functions (or classifiers)  $f : X \rightarrow Y$  in ML
- The space of world states  $S$  and actions  $A$  in Reinforcement Learning
- The space of policies  $\pi : S \rightarrow A$  in RL
- The space of feedback controllers  $\pi : x \mapsto u$  in robot control

- The configuration space  $Q$  of a robot
- The space of paths  $x : [0, 1] \rightarrow Q$  in robotics
- The space of image segmentations  $s : I \rightarrow \{0, 1\}$  in computer vision

Actually, some of these spaces are not vector spaces at all. E.g. the configuration space of a robot might have ‘holes’, be a manifold with complex topology, or not even that (switch dimensionality at some places). But to do computations in these spaces one always either introduces (local) parameterizations that make them a vector space,<sup>2</sup> or one focusses on local tangent spaces (local linearizations) of these spaces, which are vector spaces.

Perhaps the most important computation we want to do in these spaces is taking derivatives—to set them equal to zero, or do gradient descent, or Newton steps for optimization. But taking derivatives essentially requires the input space to (locally) be a vector space.<sup>3</sup> So, we also we need vector spaces because we need derivatives—and Linear Algebra to deal with the resulting equations.

### 2.1.2 What is a vector?

A vector is nothing but an element of a vector space.

It is in general not a column, array, or tuple of numbers. (But tuples of numbers are a special case of a vector space.)

### 2.1.3 What is a vector space?

**Definition 2.1** (vector space). A vector space<sup>a</sup>  $V$  is a space (=set) on which two operations, addition and multiplication, are defined as follows

- addition  $+$  :  $V \times V \rightarrow V$  is an abelian group, i.e.,
  - $a, b \in V \Rightarrow a + b \in V$  (closed under  $+$ )
  - $a + b = b + a$  (commutative)
  - $\exists$  unique  $0 \in V$  s.t.  $\forall v \in V : 0 + v = v$  (unique identity)
  - $\forall v \in V : \exists$  unique inverse  $-v$  s.t.  $v + (-v) = 0$  (unique inverse)
- multiplication  $\cdot$  :  $\mathbb{R} \times V \rightarrow V$  fulfils
 
$$\alpha(\beta v) = (\alpha\beta)v, \quad 1v = v, \quad \alpha(v + w) = \alpha v + \alpha w$$

<sup>a</sup>We only consider vector spaces over  $\mathbb{R}$ .

Roughly, this definition says that a vector space is “closed under linear transformations”, meaning that we can add and scale vectors and they remain vectors.

<sup>2</sup>E.g. by definition and  $n$ -dimensional manifold  $X$  is locally isomorphic to  $\mathbb{R}^n$ .

<sup>3</sup>Also when the space is actually a manifold; the differential is defined as a 1-form on the local tangent.

## 2.2 Vectors, dual vectors, coordinates, matrices, tensors

### 2.2.1 A taxonomy of linear functions

For simplicity we consider only functions involving a single vector space  $V$ . But all that is said transfers to the case when multiple vector spaces  $V, W, \dots$  were involved.

**Definition 2.2.**  $f : V \rightarrow X$  **linear**  $\iff f(\alpha v + \beta w) = \alpha f(v) + \beta f(w)$ , where  $X$  is any other vector space (e.g.  $X = \mathbb{R}$ , or  $X = V \times V$ ).

**Definition 2.3.**  $f : V \times V \times \dots \times V \rightarrow X$  **multi-linear**  $\iff f$  is linear in each input.

Many names are used for special linear functions—let’s make some explicit:

- $f : V \rightarrow \mathbb{R}$ , called linear functional<sup>4</sup>, or 1-form, or dual vector.
- $f : V \rightarrow V$ , called linear map, or linear transform, or vector-valued 1-form
- $f : V \times V \rightarrow \mathbb{R}$ , called bilinear functional, or 2-form
- $f : V \times V \times V \rightarrow \mathbb{R}$ , called 3-form (or unspecifically ‘multi-linear functional’)
- $f : V \times V \rightarrow V$ , called vector-valued 2-form (or unspecifically ‘multi-linear map’)
- $f : V \times V \times V \rightarrow V \times V$ , called bivector-valued 3-form
- $f : V^k \rightarrow V^m$ , called  $m$ -vector-valued  $k$ -form

This gives us a simple taxonomy of linear functions based on how many vectors a function *eats*, and how many it *outputs*. To give examples, consider some space  $X$  of systems (examples above), which might itself not be a vector space. But locally, around a specific  $x \in X$ , its tangent  $V$  is a vector space. Then

- $f : X \rightarrow \mathbb{R}$  could be a cost function over the system space.
- The differential  $df|_x : V \rightarrow \mathbb{R}$  is a 1-form, telling us how  $f$  changes when ‘making a tangent step’  $v \in V$ .
- The 2nd derivative  $d^2f|_x : V \times V \rightarrow \mathbb{R}$  is a 2-form, telling us how  $df|_x(v)$  changes when ‘making a tangent step’  $w \in V$ .
- The inner product  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$  is a 2-form.

This is simply to show that 1- and 2-forms are very common things. Another example:

- $f : \mathbb{R}^i \rightarrow \mathbb{R}^o$  is a neural network that maps  $i$  input signals to  $o$  output signals.
- Its derivative  $df|_x : \mathbb{R}^i \rightarrow \mathbb{R}^o$  is a vector-valued 1-form, telling us how each output changes with a step  $v \in \mathbb{R}^i$  in the input.
- Its 2nd derivative  $d^2f|_x : \mathbb{R}^i \times \mathbb{R}^i \rightarrow \mathbb{R}^o$  is a vector-valued 2-form.

This is simply to show that vector-valued functions, 1-forms, and 2-forms are common. Instead of being a neural network,  $f$  could also be a mapping from one param-

<sup>4</sup>The word ‘functional’ instead of ‘function’ is especially used when  $V$  is a space of functions.



eterization of a system to another, or the mapping from the joint angles of a robot to its hand position.

## 2.2.2 Bases and coordinates

We need to define some notions. I'm not commenting these definitions—train yourself in reading maths...

**Definition 2.4.**  $\text{span}(\{v_i\}_{i=1}^k) = \{\sum_i \alpha_i v_i : \alpha_i \in \mathbb{R}\}$

**Definition 2.5.**  $\{v_i\}_{i=1}^n$  **linearly independent**  $\iff \left[ \sum_i \alpha_i v_i = 0 \Rightarrow \forall_i \alpha_i = 0 \right]$

**Definition 2.6.**  $\dim(V) = \max_n \{n \in \mathbb{N} : \exists \{v_i\}_{i=1}^n \text{ lin.indep.}, v_i \in V\}$

**Definition 2.7.**  $B = \{e_i\}_{i=1}^n$  is a **basis** of  $V \iff \text{span}(B) = V$  and  $B$  lin.indep.

**Definition 2.8.** The tuple  $(v_1, v_2, \dots, v_n) \in \mathbb{R}^n$  is called **coordinates** of  $v \in V$  in the basis  $\{e_i\}_{i=1}^n$  iff  $v = \sum_i v_i e_i$

Given a basis  $\{e_i\}_{i=1}^n$ , we can describe every vector  $v$  as a linear combination  $v = \sum_i v_i e_i$  of *basic elements*—the basis vectors  $e_i$ . This general idea, that “linear things” can be described as linear combinations of “basic elements” carries over also to functions. In fact, to all the types of functions we described above: 1-forms, 2-forms, bi-vector-valued  $k$ -forms, whatever. And if we describe all these als linear combinations of basic elements we automatically also introduce coordinates for these things. To get there, we first have to introduce a second type of “basic elements”: dual vectors.

## 2.2.3 The dual vector space – and its coordinates

**Definition 2.9.** Given  $V$ , its **dual space** is  $V^* = \{f : V \rightarrow \mathbb{R} \text{ linear}\}$  (the space of linear functionals, or 1-forms). Every  $v^* \in V^*$  is called 1-form or **dual vector** (sometimes also *covector*).

First, it is easy to see that  $V^*$  is also a vector space: We can add two linear functionals,  $f = f_1 + f_2$ , and scale them, and it remains a linear functional.

Second, given a basis  $\{e_i\}_{i=1}^n$  of  $V$ , we define a corresponding dual basis  $\{\bar{e}_i\}_{i=1}^n$  of  $V^*$  simply by

$$\forall_{i,j} : \bar{e}_i(e_j) = \delta_{ij}$$

where  $\delta_{ij} = [i = j]$  is the **Kronecker delta**. Note that

$$\forall v \in V : \bar{e}_i(v) = v_i$$

That is,  $\bar{e}_i$  is the 1-form that simply maps a vector to its  $i$ th coordinate. It can be shown that  $\{\bar{e}_i\}_{i=1}^n$  is in fact a basis of  $V^*$ . (Omitted.) That tells us a lot!

$\dim(V^*) = \dim(V)$ . That is, the space of 1-forms has the same dimension as  $V$ . At this place, geometric intuition should kick in: indeed, every linear function over  $V$  could be envisioned as a “plane” over  $V$ . Such a plane can be illustrated by its iso-lines and these can be uniquely determined by their orientation and distance (same dimensionality as  $V$  itself). Also, (assuming we’d know already what a transpose or scalar product is) every 1-form must be of the form  $f(v) = c^\top v$  for some  $c \in V$ —so every  $f$  is uniquely described by a  $c \in V$ . Showing that the vector space  $V$  and its dual  $V^*$  are really twins.

The dual basis  $\{\bar{e}_i\}_{i=1}^n$  introduces coordinates in the dual space: Every 1-form  $f$  can be described as a linear combination of basis 1-forms,

$$f = \sum_i f_i \bar{e}_i$$

where the tuple  $(f_1, f_2, \dots, f_n)$  are the coordinates of  $f$ . Or

$$\text{span}(\{\bar{e}_i\}_{i=1}^n) = V^* .$$

## 2.2.4 Coordinates for every linear thing: tensors

We now have the *basic elements*: the basis vectors  $\{e_i\}_{i=1}^n$  of  $V$ , and basis 1-forms  $\{\bar{e}_i\}_{i=1}^n$  of  $V^*$ . From these, we can describe, for instance, any bivector-valued 3-form as a linear combination as follows:

$$f : V \times V \times V \rightarrow V \times V$$

$$f = \sum_{ijklm} f_{ijklm} e_i \otimes e_j \otimes \bar{e}_k \otimes \bar{e}_l \otimes \bar{e}_m$$

The  $\otimes$  is called **outer product** (or tensor product), and  $v \otimes w \in V \times W$  if  $V$  and  $W$  are finite vector spaces. For our purposes, we may think of  $v \otimes w = (v, w)$  simply as the tuple of both. Therefore  $e_i \otimes e_j \otimes \bar{e}_k \otimes \bar{e}_l \otimes \bar{e}_m$  is a 5-tuple and we have in total  $n^5$  such basis objects—and  $f_{ijklm}$  denotes the corresponding  $n^5$  coordinates. The first two indices are contra-variant, the last three covariant—these notions are explained in detail later.

## 2.2.5 Finally: Matrices

As a special case of the above, every  $f : V \rightarrow U$  can be described as a linear combination

$$f = \sum_{ij} f_{ij} e_i \otimes \bar{e}_j ,$$

where  $\{\bar{e}_j\}$  is a basis of  $V^*$  and  $\{e_i\}$  a basis of  $U$ .

Let's see how this fits with some easier view, without all this fuss about 1-forms. We already understood that the operator  $\bar{e}_j(v) = v_j$  simply picks the  $j$ th coordinate of a vector. Therefore

$$f(v) = \left[ \sum_{ij} f_{ij} e_i \otimes \bar{e}_j \right](v) = \sum_{ij} f_{ij} e_i v_j .$$

In case it helps, we can 'derive' this more slowly as

$$f(v) = f\left(\sum_k v_k e_k\right) = \sum_k v_k f(e_k) = \sum_k v_k \left[ \sum_{ij} f_{ij} e_i \otimes \bar{e}_j \right] e_k \tag{1}$$

$$= \sum_{ijk} f_{ij} v_k e_i \bar{e}_j(e_k) = \sum_{ijk} f_{ij} v_k e_i \delta_{jk} = \sum_i \left[ \sum_j f_{ij} v_j \right] e_i . \tag{2}$$

As a result, this tells us that the vector  $u = f(v) \in V$  has the coordinates  $u_i = \sum_j f_{ij} v_j$ . And the vector  $f(e_j) \in V$  has the coordinates  $f_{ij}$ , that is,  $f(e_j) = \sum_i f_{ij} e_i$ .

So there are  $n^2$  coordinates  $f_{ij}$  for a linear map  $f$ . The first index is contra-variant, the second covariant (explained later). As it so happens, the whole world has agreed on a convention on how to write such coordinate numbers on sheets of 2-dimensional paper: as a **matrix** !

$$\begin{pmatrix} f_{11} & f_{12} & \cdots & f_{1n} \\ f_{21} & f_{22} & \cdots & f_{2n} \\ \vdots & & & \vdots \\ f_{n1} & f_{n2} & \cdots & f_{nn} \end{pmatrix}$$

The first (contra-variant) index spans columns; the second (covariant) spans rows. We call this and the respective definition of a matrix multiplication as the **matrix convention** .

Note that the identity map  $\mathbf{I} : V \rightarrow V$  can be written as

$$\mathbf{I} = \sum_i e_i \otimes \bar{e}_i , \quad \mathbf{I}_{ij} = \delta_{ij} . \tag{3}$$

Equally, the (contra-variant) coordinates of a vector are written as columns

$$\begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix}$$

and the (covariant) coordinates of a 1-form  $h : V \rightarrow R$  as a row

$$(h_1 \quad h_2 \quad \cdots \quad h_n)$$

$u = f(v)$  is itself a vector, and its coordinates written as a column are

$$\begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix} = \begin{pmatrix} f_{11} & f_{12} & \cdots & f_{1n} \\ f_{21} & f_{22} & \cdots & f_{2n} \\ \vdots & & & \vdots \\ f_{n1} & f_{n2} & \cdots & f_{nn} \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix}$$

where this **matrix multiplication** is defined by  $u_i = \sum_j f_{ij}v_j$ , consistent to the above.

## 2.3 Scalar product and orthonormal basis

Please note that so far we have not in any way referred to a scalar product. All the concepts above, including the dual vector space, bases, coordinates, matrix-vector multiplication, are fully independent of the notion of a scalar product. Only now we need it.

**Definition 2.10.** A **scalar product** (also called inner product) of  $V$  is a symmetric positive definite 2-form

$$\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R} .$$

with  $\langle v, w \rangle = \langle w, v \rangle$  and  $\langle v, v \rangle > 0$  for all  $v \neq 0 \in V$ .

**Definition 2.11.** Given a scalar product, we define for every  $v \in V$  its dual  $v^* \in V^*$  as

$$v^* = \langle v, \cdot \rangle = \sum_i v_i \langle e_i, \cdot \rangle = \sum_i v_i e_i^* .$$

Note that  $\bar{e}_i$  and  $e_i^*$  are different 1-forms! The canonical dual basis  $\{\bar{e}_i\}$  is independent of an introduction of a scalar product, they were the basis to introduce coordinates for linear functions, esp. matrices. And while such coordinates do depend on a choice of basis  $\{e_i\}$ , they do *not* depend on a choice of scalar product.

The vectors  $\{e_i^*\}$  also form a basis for  $V^*$ , but a different one to the canonical basis, and one that depends on the notion of a scalar product. You can see this: the coefficients  $v_i$  of  $v^*$  in the basis  $\{e_i^*\}$  are identical to the coefficients  $v_i$  of  $v$  in the basis  $\{e_i\}$ , but different to the coefficients  $(v^*)_i$  of  $v^*$  in the basis  $\{\bar{e}_i\}$ .

**Definition 2.12.** Given a scalar product, a set of vectors  $\{v_i\}_{i=1}^n$  is called **orthonormal** iff  $\langle v_i, v_j \rangle = \delta_{ij}$ .

**Definition 2.13.** Given a scalar product and basis  $\{e_i\}$ , we define the **metric**

**tensor**  $g_{ij} = \langle e_i, e_j \rangle$ , which are the coordinates of the 2-form  $\langle \cdot, \cdot \rangle$ , that is

$$\langle \cdot, \cdot \rangle = \sum_{ij} g_{ij} \bar{e}_i \otimes \bar{e}_j .$$

This also implies that

$$\langle v, w \rangle = \sum_{ij} v_i w_j \langle e_i, e_j \rangle = \sum_{ij} v_i w_j g_{ij} = v^\top G w .$$

Although related, do not confuse  $g_{ij}$  with the usual definition of a metric  $d(\cdot, \cdot)$  in a metric space.

### 2.3.1 Properties of orthonormal bases

If we have an orthonormal basis  $\{e_i\}$ , many things simplify a lot. Throughout this subsection, we assume  $\{e_i\}$  orthonormal.

- The metric tensor  $g_{ij} = \langle e_i, e_j \rangle = \delta_{ij}$  is the identity matrix.<sup>5</sup> Such a metric is also called **Euclidean**. The norm  $\|e_i\| = 1$ . The canonical dual basis  $\{\bar{e}_i\}$  and the one defined via the scalar product  $\{e_i^*\}$  become identical,  $\bar{e}_i = e_i^* = \langle e_i, \cdot \rangle$ . Consequently,  $v$  and  $v^*$  have the same coordinates  $v_i = (v^*)_i$  w.r.t.  $\{e_i\}$  respectively  $\{\bar{e}_i\}$ .
- The coordinates of vectors can now easily be computed:

$$v = \sum_i v_i e_i \quad \Rightarrow \quad \langle e_i, v \rangle = \langle e_i, \sum_j v_j e_j \rangle = \sum_j \langle e_i, e_j \rangle v_j = v_i \quad (4)$$

- The coordinates of a linear transform can equally easily be computed: Given a linear transform  $f : V \rightarrow U$  an *arbitrary* (e.g. non-orthonormal) input basis  $\{v_i\}$  of  $V$ , but an orthonormal basis  $\{u_i\}$ , then

$$f = \sum_{ij} f_{ij} u_i \otimes \bar{v}_j \quad \Rightarrow \quad (5)$$

$$\langle u_i, f v_j \rangle = \langle u_j, \sum_{kl} f_{kl} u_k \otimes \bar{v}_l(v_j) \rangle = \sum_{kl} f_{kl} \langle u_j, u_k \rangle \bar{v}_l(v_j) = \sum_{kl} f_{kl} \delta_{jk} \delta_{lj} = f_{ij} \quad (6)$$

- The projection onto a basis vector is given by  $e_i \langle e_i, \cdot \rangle$ .
- The projection onto the span of several basis vectors  $\{e_1, \dots, e_k\}$  is given by  $\sum_{i=1}^k e_i \langle e_i, \cdot \rangle$ .

<sup>5</sup>Being picky, a metric is not a matrix but a twice covariant tensor (a row of rows). That's why it is correctly called metric tensor.

- The identity mapping  $\mathbf{I} : V \rightarrow V$  is given by  $\mathbf{I} = \sum_{i=1}^{\dim(V)} e_i \langle e_i, \cdot \rangle$ .
- The scalar product with an orthonormal basis is

$$\langle v, w \rangle = \sum_{ij} v_i w_j \delta_{ij} = \sum_i v_i w_i$$

which, using matrix convention, can also be written as

$$\langle v, w \rangle = (v_1 \ v_2 \ \dots \ v_n) \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix} = v^\top w = \sum_i v_i w_i,$$

where for the first time we introduced the **transpose** which, in the matrix convention, swaps columns to rows and rows to columns.

As a general note, a transposed vector “eats a vector and outputs a scalar”. That is  $v^\top : V \rightarrow \mathbb{R}$  should be thought of as a 1-form! Due to the matrix conventions, it generally is the case that “rows eat columns”, that is, every row index should always be thought of as relating to a 1-form (dual vector), and every column index as relating to a vector. That is totally consistent to our definition of coordinates.

For an orthonormal basis we also have

$$v^*(w) = \langle v, w \rangle = v^\top w.$$

That is,  $v^\top$  is the coordinate representation of the 1-form  $v^*$ . (Which also says, that the coordinates of the 1-form  $v^*$  in the special basis  $\{e_i^*\}_{i=1}^n \subset V^*$  coincide with the coordinates of the vector  $v$ .)

## 2.4 The Structure of Transforms & Singular Value Decomposition

We focus here on linear transforms (or “linear maps”)  $f : V \rightarrow U$  from one vector space to another (or the same). It turns out that such transforms have a very specific and intuitive structure, which is captured by the singular value decomposition.

### 2.4.1 The Singular Value Decomposition Theorem

We state the following theorem:

**Theorem 2.1.** Given two vector spaces  $\mathcal{V}$  and  $\mathcal{U}$  with scalar products,  $\dim(\mathcal{V}) = n$  and  $\dim(\mathcal{U}) = m$ , for every linear transform  $f : \mathcal{V} \rightarrow \mathcal{U}$  there exist a  $k \leq \min(n, m)$  and orthonormal vectors  $\{v_i\}_{i=1}^k \subset \mathcal{V}$ , orthonormal vectors  $\{u_i\}_{i=1}^k \subset \mathcal{U}$ , and

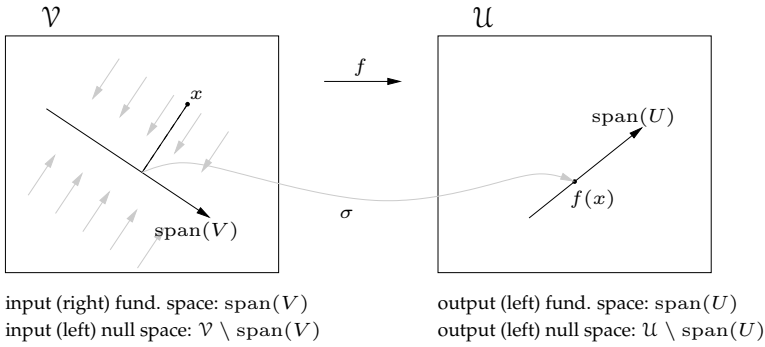


Figure 1: A linear transformation  $f = \sum_{i=1}^k \sigma_i u_i v_i^*$  can be described as: take the input  $x$ , project it onto the first input fundamental vector  $v_1$  to yield a scalar, stretch/squeeze it by  $\sigma_1$ , and project this into the first output fundamental vector  $u_1$ ; repeat this for all  $i = 1, \dots, k$ , and add the results.

positive scalars  $\sigma_i > 0, i = 1, \dots, k$ , such that

$$f = \sum_{i=1}^k \sigma_i u_i v_i^*$$

As above,  $v_i^* = \langle v_i, \cdot \rangle$  is the basis 1-form that picks the  $i$ th coordinate of a vector in the basis  $\{v_i\}_{i=1}^k \subset \mathcal{V}$ .<sup>a</sup>

<sup>a</sup>Note that  $\{v_i\}_{i=1}^k$  may not be a full basis of  $\mathcal{V}$  if  $k < n$ . But because  $\{v_i\}$  is orthonormal,  $\langle v_i, \cdot \rangle$  uniquely picks the  $i$ th coordinate no matter how  $\{v_i\}_{i=1}^k$  is completed with further  $n - k$  vectors to become a full basis.

We first restate this theorem equivalently in coordinates.

**Theorem 2.2.** For every matrix  $A \in \mathbb{R}^{m \times n}$  there exists a  $k \leq n, m$  and orthonormal vectors  $\{v_i\}_{i=1}^k \subset \mathbb{R}^n$ , orthonormal vectors  $\{u_i\} \subset \mathbb{R}^m$ , and positive scalars  $\sigma_i > 0, i = 1, \dots, k$ , such that

$$A = \sum_{i=1}^k \sigma_i u_i v_i^\top = USV^\top$$

where  $V = (v_1, \dots, v_k) \in \mathbb{R}^{n \times k}$ ,  $U = (u_1, \dots, u_k) \in \mathbb{R}^{m \times k}$  contain the orthonormal bases vectors as columns and  $S = \text{diag}(\sigma_1, \dots, \sigma_k)$ .

Just to explicitly show the transition from coordinate-free to the coordinate-based theorem, consider arbitrary orthonormal bases  $\{e_i\}_{i=1}^n \subset \mathcal{V}$  and  $\{\hat{e}_i\}_{i=1}^m \subset \mathcal{U}$ . For

$x \in \mathcal{V}$  we have

$$\begin{aligned} f(x) &= \sum_{i=1}^k \sigma_i u_i \langle v_i, x \rangle = \sum_{i=1}^k \sigma_i \left( \sum_l u_{li} \hat{e}_l \right) \left\langle \sum_j v_{ji} e_j, \sum_k x_k e_k \right\rangle & (7) \\ &= \sum_{i=1}^k \sigma_i \left( \sum_l u_{li} \hat{e}_l \right) \sum_{jk} v_{ji} x_k \delta_{jk} = \sum_l \left[ \sum_{i=1}^k u_{li} \sigma_i \sum_j v_{ji} x_j \right] \hat{e}_l = \sum_l \left[ U S V^T x \right]_l \hat{e}_l & (8) \end{aligned}$$

where  $v_{ji}$  are the coordinates of  $v_i$ ,  $u_{li}$  the coordinates of  $u_i$ ,  $U = (u_1, \dots, u_k)$  is the matrix containing  $\{u_i\}$  as columns,  $V = (v_1, \dots, v_k)$  the matrix containing  $\{v_i\}$  as columns, and  $S = \text{diag}(\sigma_1, \dots, \sigma_k)$  the diagonal matrix with elements  $\sigma_i$ .

We add some definitions based on this:

**Definition 2.14.** The **rank**  $\text{rank}(f) = \text{rank}(A)$  of a transform  $f$  or its matrix  $A$  is the unique minimal  $k$ .

**Definition 2.15.** The **determinant** of a transform  $f$  or its matrix  $A$  is

$$\det(f) = \det(A) = \det(S) = \pm \begin{cases} \chi \prod_{i=1}^n \sigma_i & \text{for } \text{rank}(f) = n = m \\ 0 & \text{otherwise} \end{cases},$$

where  $\chi \in \{-1, +1\}$  depends on whether the transform is a reflection or not.

The last definition is a bit flaky, as  $\chi$  is not properly defined. If, alternatively, in the above theorems we would require  $V$  and  $U$  to be rotations, that is, elements of  $SO(n)$  (of the *special* orthogonal group); then negative  $\sigma$ 's would indicate such a reflection and  $\det(A) = \prod_{i=1}^n \sigma_i$ . But above we requires  $\sigma$ 's to be strictly positive and  $V$  and  $U$  only orthogonal. So fundamental space vectors  $v_i$  and  $u_i$  could flip sign. The  $\phi$  above indicates how many flip sign.

**Definition 2.16.** a) The **row space** (also called right (=input) fundamental space) of a transform  $f$  is  $\text{span}\{v_i\}_{i=1}^{\text{rank}(f)}$ . The **input null space** (or right null space)  $\mathcal{V}_\perp$  is the subspace orthogonal to the row space, such that  $v \in \mathcal{V}_\perp \Rightarrow f(v) = 0$ .

b) The **column space** or (also called left (=output) fundamental space) of a transform  $f$  is  $\text{span}\{u_i\}_{i=1}^{\text{rank}(f)}$ . The **output null space** (or left null space)  $\mathcal{U}_\perp$  is the subspace orthogonal to the column space, such that  $u \in \mathcal{U}_\perp \Rightarrow \langle f(\cdot), u \rangle = 0$ .

## 2.5 Point of departure from the coordinate-free notation

The coordinate-free introduction of vectors and transforms helps a lot to understand what these fundamentally are. That coordinate vectors and matrices are 'just' coordinates and rely on a choice of basis; what a metric  $g_{ij}$  really is; that only for a Euclidean metric the inner product satisfies  $\langle v, w \rangle = v^T w$ . Further, the coordinate-free



view is essential to understand that vector coordinates behave differently to 1-form coordinates (e.g., “gradients”!) under a transformation of the basis. We postpone this discussion of contra- versus covariance.

However, we now understand that columns correspond vectors, rows to 1-forms, and in the Euclidean case the 1-form  $v^*$  directly corresponds to  $v^\top$ , in the non-Euclidean to  $v^\top G$ . In applications we typically represent things from start in orthonormal bases (including perhaps non-Euclidean metrics), there is not much gain sticking to the coordinate-free notation in most cases. Only when the matrix notation gets confusing (and this happens, e.g. when trying to compute something like the “Jacobian of a Jacobian”, or applying the chain and product rule for a matrix expression  $\partial_x f(x)^\top A(x) b(x)$ ) it is always a safe harbour to remind what we actually talk about.

Therefore, in the rest of the notes we rely on the normal coordinate-based view. Only in some explanations we remind on the coordinate-free view when helpful.

## 2.6 Filling SVD with life

In the following we list some statements—all of them relate to the SVD theorem and together they’re meant to give a more intuitive understanding of the equation  $A = \sum_{i=1}^k \sigma_i u_i v_i^\top = USV^\top$ .

### 2.6.1 Understand $vv^\top$ as projection

- The **projection** of a vector  $x \in \mathcal{V}$  onto a vector  $v \in \mathcal{V}$ , is given by

$$x_{\parallel} = \frac{1}{v^2} v \langle v, x \rangle \quad \text{or} \quad \frac{1}{v^2} v v^\top x .$$

Here, the  $\frac{1}{v^2}$  is normalizing in case  $v$  does not have length  $|v| = 1$ .

- The projection-on- $v$ -matrix  $vv^\top$  is symmetric, semi-pos-def, and has  $\text{rank}(vv^\top) = 1$ .
- The projection of a vector  $x \in \mathcal{V}$  onto a subvector space  $\text{span}\{v_i\}_{i=1}^k$  for *orthonormal*  $\{v_i\}_{i=1}^k$  is given by

$$x_{\parallel} = \sum_i v_i v_i^\top x = VV^\top x$$

where  $V = (v_1, \dots, v_k) \in \mathbb{R}^{n \times k}$ .

- The projection matrix  $VV^\top$  for orthonormal  $V$  is symmetric, semi-pos-def, and has  $\text{rank}(VV^\top) = k$ .

### 2.6.2 SVD for symmetric matrices

- Thm 2  $\Rightarrow$  Every symmetric matrix  $A$  is of the form

$$A = \sum_i \lambda_i v_i v_i^\top = V \Lambda V^\top$$

for orthonormal  $V = (v_1, \dots, v_k)$ . Here  $\lambda_i = \pm \sigma_i$  and  $\Lambda = \text{diag}(\lambda)$  is the diagonal matrix of  $\lambda$ 's. This describes nothing but a stretching/squeezing along orthogonal projections.

- The  $\lambda_i$  and  $v_i$  are also the eigenvalues and eigenvectors of  $A$ , that is, for all  $i = 1, \dots, k$ :

$$A v_i = \lambda_i v_i.$$

The SVD  $A = V S V^\top = V S V^{-1}$  is therefore also the eigendecomposition of  $A$ .

- The pseudo-inverse of a symmetric matrix is

$$A^\dagger = \sum_i \lambda_i^{-1} v_i v_i^\top = V S^{-1} V^\top$$

which simply does the reverse stretching/squeezing along the same orthogonal projections. Note that

$$A A^\dagger = A^\dagger A = V V^\top$$

is the projection on  $\{v_i\}_{i=1}^k$ . For  $\text{rank}(A) = n$  we have  $V V^\top = \mathbf{I}$  and  $A^\dagger = A^{-1}$ . For  $\text{rank}(A) < n$ , we have that  $A^\dagger y$  minimizes  $\min_x \|A x - y\|^2$ , but there are infinitely many  $x$ 's that minimize this, spanned by the null space of  $A$ .  $A^\dagger y$  is the minimizer closest to zero (with smallest norm).

- Consider a data set  $D = \{x_i\}_{i=1}^m$ ,  $x_i \in \mathbb{R}^n$ . For simplicity assume it has zero mean,  $\sum_{i=1}^m x_i = 0$ . The **covariance matrix** is defined as

$$C = \frac{1}{n} \sum_i x_i x_i^\top = \frac{1}{n} X^\top X$$

where (consistent to ML lecture convention) the data matrix  $X$  contains  $x_i^\top$  in the  $i$ th row. Each  $x_i x_i^\top$  is a projection.  $C$  is symmetric and semi-pos-dev. Using SVD we can write

$$C = \sum_i \lambda_i v_i v_i^\top$$

and  $\lambda_i$  is the data variance along the eigenvector  $v_i$ ;  $\sqrt{\lambda_i}$  the standard deviation along  $v_i$ ; and  $\sqrt{\lambda_i} v_i$  the principle axis vectors that make the ellipsoid we typically illustrate covariances with.

### 2.6.3 SVD for general matrices

- For  $\text{rank}(A) = n$ , the determinant of a matrix is  $\det(A) = \prod_i \sigma_i$ . We may define the **volume** spanned by any  $\{b_i\}_{i=1}^n$  as

$$\text{vol}(\{b_i\}_{i=1}^n) = \det(B), \quad B = (b_1, \dots, b_n) \in \mathbb{R}^{n \times n}.$$

It follows that

$$\text{vol}(\{Ab_i\}_{i=1}^n) = \det(A) \det(B)$$

that is, the volume is being multiplied with  $\det(A) = \prod_i \sigma_i$ , which is consistent with our understanding of transforms as stretchings/squeezings along projections.

- The pseudo-inverse of a general matrix is

$$A^\dagger = \sum_i \sigma_i^{-1} v_i u_i^\top = V S^{-1} U^\top.$$

If  $k = n$  (full input rank),  $\text{rank}(A^\top A) = n$  and

$$(A^\top A)^{-1} A^\top = (V S U^\top U S V^\top)^{-1} V S U^\top = V^{-\top} S^{-2} V^{-1} V S U^\top = V S^{-1} U^\top = A^\dagger$$

and  $A^\dagger$  is also called left pseudoinverse because  $A^\dagger A = \mathbf{I}_n$ .

If  $k = m$  (full output rank),  $\text{rank}(A A^\top) = m$  and

$$A^\top (A A^\top)^{-1} = V S U^\top (U S V^\top V S U^\top)^{-1} = V S U^\top U^{-\top} S^{-2} U^{-1} = V S^{-1} U^\top = A^\dagger$$

and  $A^\dagger$  is also called right pseudoinverse because  $A A^\dagger = \mathbf{I}_m$ .

- Assume  $m = n$  (same input/output dimension, or  $\mathcal{V} = \mathcal{U}$ ), but  $k < n$ . Then there exist orthogonal  $V, U \in \mathbb{R}^{n \times n}$  such that

$$A = U D V^\top, \quad D = \text{diag}(\sigma_1, \dots, \sigma_k, 0, \dots, 0) = \begin{pmatrix} S & 0 \\ 0 & 0 \end{pmatrix}.$$

Here,  $V$  and  $U$  contain a full orthonormal basis instead of only  $k$  orthonormal vectors. But the diagonal matrix  $D$  projects all but  $k$  of those to zero. Every square matrix  $A \in \mathbb{R}^{n \times n}$  can be written like this.

**Definition 2.17 (Rotation).** Given a scalar-product  $\langle \cdot, \cdot \rangle$  on  $V$ , a linear transform  $f : V \rightarrow V$  is called *rotation* iff it preserves the scalar product, that is,

$$\forall v, w \in V : \langle f(v), f(w) \rangle = \langle v, w \rangle. \tag{9}$$

- Every rotation matrix is orthogonal, i.e., composed of columns of orthonormal vectors.
- Every rotation has rank  $n$  and  $\sigma_{1, \dots, n} = 1$ . (No stretching/squeezing.)
- Every square matrix can be written as

$$\text{rotation}_U \cdot \text{scaling}_D \cdot \text{rotation}_{V^{-1}}$$

## 2.7 Eigendecomposition

**Definition 2.18.** The **diagonalization** or **eigendecomposition** of a square matrix  $A \in \mathbb{R}^{n \times n}$  is (if it exists!)

$$A = Q\Lambda Q^{-1}$$

where  $L = \text{diag}(\lambda_1, \dots, \lambda_n)$  is a diagonal matrix of eigenvalues. Each column  $q_i$  of  $Q$  is an **eigenvector**.

- The set of eigenvalues is the set of roots of the **characteristic polynomial**  $p_A(\lambda) = \det(\lambda I - A)$ . Why? Because then  $A - \lambda I$  has ‘volume’ zero (or rank  $< n$ ), showing that one vector is mapped to zero, that is  $0 = (A - \lambda I)x = Ax - \lambda x$ .
- Is every square matrix diagonalizable? No, but only an  $n^2 - 1$ -dimensional subset of matrices are not diagonalizable; most matrices are. A matrix is *not* diagonalizable if an eigenvalue  $\lambda$  has multiplicity  $k$  (more precisely,  $\lambda$  is a root of  $p_A(\lambda)$  with multiplicity  $k$ ), but  $n - \text{rank}(A - \lambda I)$  (the dimensionality of the span of the eigenvectors of  $\lambda$ !) is less than  $k$ . Therefore, the eigenvectors of  $\lambda$  are not linearly independent; they do not span the necessary  $k$  dimensions. So, only very “special” matrices are not diagonalizable. Random matrices are (with prob 1).
- Symmetric matrices?  $\rightarrow$  SVD
- Rotations? Not real. But complex! Think of oscillating projection onto eigenvector. If  $\phi$  is the rotation angle,  $e^{\pm i\phi}$  are eigenvalues.

### 2.7.1 Power Method

To find the largest eigenvector of  $A$ , initialize  $x$  randomly and iterate

$$x \leftarrow Ax, \quad x \leftarrow \frac{1}{\|x\|} x$$

- If this converges,  $x$  must be an eigenvector and  $\lambda = x^T Ax$  the eigenvalue.
- If  $A$  is diagonalizable, and  $x$  is initially a non-zero linear combination of all eigen vectors, then it is obvious that  $x$  will converge to the “largest” (in *absolute* terms  $|\lambda_i|$ ) eigenvector (=eigenvector with largest eigenvalue). Actually, if the largest (by norm) eigenvector is negative, then it doesn’t really converge by flip sign at every iteration.

### 2.7.2 Power Method including the smallest eigenvalue

A trick, hard to find in the literature, to also compute the smallest eigenvalue and -vector is the following. We assume all eigenvalues to be positive. Initialize  $x$  and  $y$

randomly, iterate

$$x \leftarrow Ax, \quad \lambda \leftarrow \|x\|, \quad x \leftarrow x/\lambda, \quad y \leftarrow (\lambda I - A)y, \quad y \leftarrow y/\|y\|$$

Then  $y$  will converge to the smallest eigenvector, and  $\lambda - \|y\|$  will be its eigenvalue. Note that (in the limit)  $A - \lambda I$  only has negative eigenvalues, therefore  $\|y\|$  should be positive. Finding smallest eigenvalues is a common problem in model fitting.

### 2.7.3 Why should I care about Eigenvalues and Eigenvectors?

- Least squares problems (finding smallest eigenvalue/-vector); (e.g. Camera Calibration, Bundle Adjustment, Plane Fitting)
- PCA
- stationary distributions of Markov Processes
- Page Rank
- Spectral Clustering
- Spectral Learning (rather novel approach to training HMMs, for instance)

## 2.8 Point of Departure: Numerics to compute these things

We will not go into details of numerics. Nathan's script gives a really nice explanation of the QR-method. I just mention two things:

(i) The most important forms of matrices for numerics are diagonal matrices, orthogonal matrices, and upper triangular matrices. One reason is that all three types can very easily be inverted. A lot of numerics is about finding **decompositions** of general matrices into products of these special-form matrices, e.g.:

- QR-decomposition:  $A = QR$  with  $Q$  orthogonal and  $R$  upper triangular.
- LU-decomposition:  $A = LU$  with  $U$  and  $L^T$  upper triangular.
- Cholesky decomposition: (symmetric)  $A = C^T C$  with  $C$  upper triangular
- Eigen- & singular value decompositions

Often, these decompositions are intermediate steps to compute eigenvalue or singular value decompositions.

(ii) Use linear algebra packages. At the origin of all is LAPACK; browse through <http://www.netlib.org/lapack/lug/> to get an impression of what really has been one of the most important algorithms in all technical areas of the last half century. Modern wrappers are: Matlab (Octave), which originated as just a console interface to LAPACK; the C++-library Eigen; or the Python NumPy.

## 2.9 Examples and Exercises

### 2.9.1 Projection

Let's define the projection of a vector  $x$  onto a vector  $v$  via an optimality principle: We search for

$$\alpha^* = \operatorname{argmin}_{\alpha} \|x - \alpha v\|^2$$

such that  $\alpha^* v \in \operatorname{span}\{v\}$  is the vector closest to  $x$ . Compute  $\alpha^*$ . (Consider coordinate vectors only, assuming an arbitrary orthonormal basis.) Hopefully this verifies  $x_{\parallel v} = (1/v^2)vv^T x$ .

### 2.9.2 Matrix equations

a) Let  $X, A$  be arbitrary matrices,  $A$  invertible. Solve for  $X$ :

$$XA + A^T = \mathbf{I}$$

b) Let  $X, A, B$  be arbitrary matrices,  $(C - 2A^T)$  invertible. Solve for  $X$ :

$$X^T C = [2A(X + B)]^T$$

c) Let  $x \in \mathbb{R}^n, y \in \mathbb{R}^d, A \in \mathbb{R}^{d \times n}$ .  $A$  obviously *not* invertible, but let  $A^T A$  be invertible. Solve for  $x$ :

$$(Ax - y)^T A = \mathbf{0}_n^T$$

d) As above, additionally  $B \in \mathbb{R}^{n \times n}, B$  positive-definite. Solve for  $x$ :

$$(Ax - y)^T A + x^T B = \mathbf{0}_n^T$$

### 2.9.3 PCA

### 2.9.4 Exercise: The Linear Robot (Nathan's)

## 3 Derivatives

### 3.1 The coordinate-free view: The meaning of life for a derivative is to eat a vector and output a scalar

We briefly state the general concept of a derivative in coordinate-free form. Practitioners may skip this.

**Definition 3.1.** Given a function  $f : V \rightarrow G$  we define the differential

$$df|_x : V \rightarrow G, v \mapsto \lim_{h \rightarrow 0} \frac{f(x + hv) - f(x)}{h}$$

This definition is very general. It holds whenever  $G$  is a continuous space that allows the definition of this limit and the limit exists ( $f$  is differentiable). The notation  $df|_x$  reads “the differential at location  $x$ ”, i.e., evaluating this derivative at location  $x$ .

Note that  $df|_x$  is a mapping from a “tangent vector”  $v$  to output-change. Further, by this definition  $df|_x$  is linear. Therefore  $df|_x$  is a  $G$ -valued 1-form. As discussed earlier, we can introduce coordinates for 1-forms; these coordinates are what typically called is the “gradient” or “Jacobian”. But here we explicitly see that we speak of coordinates of a 1-form.

## 3.2 Don’t confuse partial and total derivative

### 3.2.1 Partial derivative

In coordinates, we may think of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  simply as a function with  $n$  arguments,  $f(x_1, \dots, x_n)$ . In the following, we think of  $f(x_1, \dots, x_n)$  as an explicit algebraic expression of these  $n$  arguments.

**Definition 3.2.** The *partial* derivative of a function of multiple arguments  $f(x_1, \dots, x_n)$  is the standard derivative w.r.t. only one of its arguments,

$$\frac{\partial}{\partial x_i} f(x_1, \dots, x_n) = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_i + h, \dots, x_n) - f(x)}{h}.$$

### 3.2.2 Chain Rule, Function Networks and the total derivative

An example: We have three real-valued quantities  $x$ ,  $g$  and  $f$  which depend on each other: let’s say

$$f(x, g) = 3x + 2g \quad \text{and} \quad g(x) = 2x.$$

Question: What is the “derivative of  $f$  w.r.t.  $x$ ”?

The correct answer: Which one do you mean? The partial or total?

The partial derivative defined above really thinks of  $f(x, g)$  as a function of two arguments; and does not at all care about whether there might be dependencies of these arguments. It only looks at  $f(x, g)$  alone and takes the partial derivative

(=derivative w.r.t. one function argument):

$$\frac{\partial}{\partial x} f(x, g) = 3$$

HOWEVER, if you suddenly talk about  $f(x, g(x))$  that's a totally different story, because  $f(x, g(x))$  is a function (algebraic expression) of  $x$  only, and

$$\frac{\partial}{\partial x} f(x, g(x)) = \frac{\partial}{\partial x} 3x + 2(2x) = 7$$

Bottom line, the definition of the partial derivative really depends on what you explicitly defined as the arguments of the function. Let's make this more general:

**Definition 3.3.** A **function network** is a DAG of  $n$  quantities  $x_i$  where each quantity is a deterministic function of a set of parents  $\pi(i) \subset \{1, \dots, n\}$ , that is

$$x_i = f_i(x_{\pi(i)})$$

where  $x_{\pi(i)} = (x_j : j \in \pi(i))$  is the tuple of parent values. This could also be called *deterministic Bayes net*.

In a function network all values can be computed deterministically if the input values (which do have no parents) are given. We ask: Assume we have a variation  $dx$  of some input value, how do all other values vary? The answer is the general chain rule, of which there exist two versions. Here is version 1:

**Identities 3.1** (General chain rule (Forward-Version)).

$$\frac{df}{dx} = \frac{\partial f}{\partial x} + \sum_{\substack{g \in \pi(f) \\ g \neq x}} \frac{\partial f}{\partial g} \frac{dg}{dx}$$

Read this as follows: "the change of  $f$  is given as the change from its direct dependence on  $x$  plus the change from its direct dependence on  $g$  times the change of  $g$  plus...".

This rule also defines the **total derivative** of  $\frac{df}{dx}$  w.r.t.  $x$ . Note how different these two notions of derivatives are by definition: a partial derivative only looks at a function itself and takes a limit of differences w.r.t. one argument—no notion of further dependencies. The total derivative asks how, in a function network, one value changes with a change of another.

Here is another version of the general chain rule:



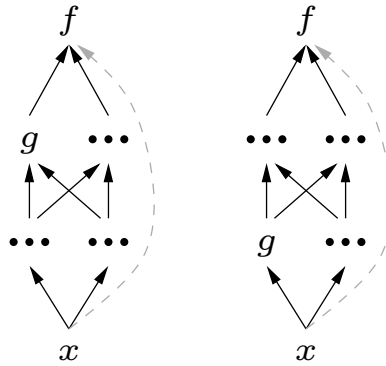


Figure 2: General Chain Rule. Left: Forward-Version, Right: Backward-Version. They gray arc denotes the direct dependence  $\frac{\partial f}{\partial x}$ , which is excluded from the summations by the ' $\neq$ ' constraints.

**Identities 3.2** (General chain rule (Backward-Version)).

$$\frac{df}{dx} = \frac{\partial f}{\partial x} + \sum_{\substack{g: x \in \pi(g) \\ g \neq f}} \frac{df}{dg} \frac{\partial g}{\partial x}$$

Figure 2 illustrates the fwd and bwd versions of the chain rule. The bwd version allows you to propagate back, given gradients  $\frac{df}{dg}$  from top to  $g$ , one step further down, from top to  $x$ . The fwd version allows you to propagate forward, given gradients  $\frac{dg}{dx}$  from  $g$  to bottom, one step further up, from  $f$  to bottom. Both versions are recursive equations. If you would resursively plug the definition for a given functional network, both of them would yield the same expression of  $\frac{df}{dx}$  in terms of partial derivatives only.

Let's compare to the chain rule as it is usually given (written more precisely),

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f(g)}{\partial g} \Big|_{g=g(x)} \frac{\partial g(x)}{\partial x}$$

Note that we here very explicitly notated that  $\frac{\partial f(g)}{\partial g}$  considers  $f$  to be a function of the argument  $g$ , which is evaluated at  $g = g(x)$ . Written like this, the rule is fine. But the above discussion and explicitly distinguishing between partial and total derivative is, when things get complicated, less prone to confusion.

### 3.3 “Gradient vectors”, Co- and contra-variance, Jacobian, Hessian

#### 3.3.1 Gradient, Jacobian, and Hessian

Lets focus again on functions  $f : \mathbb{R}^n \rightarrow \mathbb{R}^d$

**Definition 3.4.** Given  $f : \mathbb{R}^n \rightarrow \mathbb{R}^d$ , we define the derivative (also called **Jacobian matrix**) as

$$\frac{\partial}{\partial x} f(x) = \begin{pmatrix} \frac{\partial}{\partial x_1} f_1(x) & \frac{\partial}{\partial x_2} f_1(x) & \dots & \frac{\partial}{\partial x_n} f_1(x) \\ \frac{\partial}{\partial x_1} f_2(x) & \frac{\partial}{\partial x_2} f_2(x) & \dots & \frac{\partial}{\partial x_n} f_2(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial}{\partial x_1} f_n(x) & \frac{\partial}{\partial x_2} f_n(x) & \dots & \frac{\partial}{\partial x_n} f_n(x) \end{pmatrix}$$

Consistent with this, we also define the derivative of  $f : \mathbb{R}^n \rightarrow \mathbb{R}^1$  as the row vector of partial derivatives

$$\frac{\partial}{\partial x} f(x) = \left( \frac{\partial}{\partial x_1} f, \dots, \frac{\partial}{\partial x_n} f \right).$$

Further, we define the **gradient** as the corresponding column “vector”

$$\nabla f(x) = \left[ \frac{\partial}{\partial x} f(x) \right]^\top.$$

The “purpose” of a derivative is to output a change of function value when being multiplied to a change of input  $\delta$ . That is, in first order approximation, we have

$$f(x + \delta) - f(x) \doteq \frac{\partial}{\partial x} f(x) \delta.$$

This equation holds, no matter if output space is  $\mathbb{R}^d$  or  $\mathbb{R}$ . In the gradient notation we have

$$f(x + \delta) - f(x) \doteq \nabla f(x)^\top \delta.$$

Note that this latter equation is truly independent of the choice of vector space basis and independent of an optional metric or scalar product in  $V$ : This transpose should not be understood as a scalar product between two vectors, but rather as undoing the transpose in the definition of  $\nabla f$ . As mentioned earlier, all this is consistent to understanding these derivatives as coordinates of a 1-form.

In the abstract, coordinate-free notation, the second derivative would be defined as the 2-form

$$d^2 f|_x : V \times V \rightarrow G, \tag{10}$$

$$(v, w) \mapsto \lim_{h \rightarrow 0} \frac{df|_{x+hw}(v) - f|_x(v)}{h} \tag{11}$$

$$= \lim_{h,l \rightarrow 0} \frac{f(x + hw + lv) - f(x + hw) - f(x + lv) + f(x)}{hl} \quad (12)$$

**Definition 3.5.** In coordinates, we define the **Hessian** of a scalar function  $f : \mathbb{R} \rightarrow \mathbb{R}$  as the symmetric matrix

$$\nabla^2 f(x) = \frac{\partial}{\partial x} \nabla f(x) = \begin{pmatrix} \frac{\partial^2}{\partial x_1 \partial x_1} f & \frac{\partial^2}{\partial x_1 \partial x_2} f & \cdots & \frac{\partial^2}{\partial x_1 \partial x_n} f \\ \frac{\partial^2}{\partial x_2 \partial x_1} f & \frac{\partial^2}{\partial x_2 \partial x_2} f & \cdots & \frac{\partial^2}{\partial x_2 \partial x_n} f \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial x_n \partial x_1} f & \frac{\partial^2}{\partial x_n \partial x_2} f & \cdots & \frac{\partial^2}{\partial x_n \partial x_n} f \end{pmatrix}$$

The Hessian can be thought of as the coordinates of the 2-form defined above (which would actually be a row vector of row vectors), or as the Jacobian of  $\nabla f$ . Using the Hessian, we can express the 2nd order approximation of  $f$  as:

$$f(x + \delta) \doteq f(x) + \partial f(x) \delta + \delta^\top \nabla^2 f(x) \delta.$$

As a Hessian  $H$  is symmetric we can decompose it as  $H = \sum_i \lambda_i h_i h_i^\top$  with eigenvalues  $\lambda_i$  and eigenvectors  $h_i$ .

$\lambda_i$  gives the curvature (as in the normal 1D 2nd derivative) along an eigenvector  $h_i$ . If  $\lambda_i > 0$ ,  $f$  is convex along  $h_i$ ; if  $\lambda_i < 0$ ,  $f$  is concave along  $h_i$ . If some  $\lambda_i$ 's are positive, some negative, and we are at a zero-gradient point  $\partial_x f = 0$ , this is called a **saddle point**.

### 3.3.2 Why the partial derivative coordinates are co-variant

Consider the following example: We have the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f(x) = x_1 + x_2$ . The function's partial derivative is of course  $\frac{\partial f}{\partial x} = (1 \ 1)$ . Now let's transform the coordinates of the space: we introduce new coordinates  $(z_1, z_2) = (2x_1, x_2)$  or  $z = Bx$  with  $B = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}$ . The same function, written in the new coordinates, is  $f(z) = \frac{1}{2}z_1 + z_2$ . The partial derivatives of that same function, written in these new coordinates, is  $\frac{\partial f}{\partial z} = (\frac{1}{2} \ 1)$ .

Considering a metric tensor  $A_x$  in  $x$ -coordinates, if we require that  $\langle z, z \rangle = \langle x, x \rangle = x^\top A_x x$  it follows that the metric tensor in  $z$ -coordinates is  $A_z = B^{-\top} A_x B^{-1}$ . We summarize:

- Let  $B$  be a matrix that describes a linear transformation in coordinates
- A coordinate vector  $x$  transforms as  $z = Bx$
- The gradient vector  $\nabla_x f(x)$  transforms as  $\nabla_z f(z) = B^{-\top} \nabla_x f(x)$

- The metric  $A_x$  transforms as  $A_z = B^\top A_x B^{-1}$
- The steepest descent direction  $A_x^{-1} \nabla_x f(x)$  transforms as  $A_z^{-1} \nabla_z f(z) = B A_x^{-1} \nabla_x f(x)$ , like a normal coordinate vector.

Generally, if we have two kinds of mathematical objects and when we multiply them together this gives us a scalar. The scalar shouldn't depend on any choice of coordinate system and is therefore invariant against coordinate transforms. Then, if one of the objects transforms in a **covariant** ("transforming *with* the transformation") manner, the other object must transform in a **contra-variant** ("transforming *contrary* to the transformation") manner to ensure that the resulting scalar is invariant. This is a general principle: whenever two things multiply together to give an invariant scalar, one should transform co- the other contra-variant. Above we introduced  $m$ -vector-valued  $k$ -forms, and tensors as their coordinates. Actually, co- and contra-variant is a statement about tensor indices. The metric tensor (coordinates) of a 2-form  $A$  is, for instance, twice co-variant. The inverse metric tensor is twice covariant.

Now, the above avoids extelling you exactly which one is call co- and which one called contra-variant. Here are some confusions possible. First let's stick to the standard explanation that explicitly talks about basis transformations. Citing Wikipedia:

- "For a vector to be basis-independent, the components of the vector must contra-vary with a change of basis to compensate. That is, the matrix that transforms the vector of components must be the inverse of the matrix that transforms the basis vectors. The components of vectors (as opposed to those of dual vectors) are said to be contravariant.
- For a dual vector (also called a covector) to be basis-independent, the components of the dual vector must co-vary with a change of basis to remain representing the same covector. That is, the components must be transformed by the same matrix as the change of basis matrix. The components of dual vectors (as opposed to those of vectors) are said to be covariant."

So far so good. This gives a clear explanation for why vector coordinates are called contra- and 1-form coordinates co-variant.

Some people though, like me, even use the words *co-/contra-variant* also as adjective for a vector or 1-form itself, not their coordinates. That must be confusing. One origin<sup>6</sup> is in the physicist's convention to use the word co-variance almost synonymously to "invariant"—and in the above we definitely want the "vector [itself] to be basis-independent", so invariant under change of basis. Such invariant vectors are sometimes called co-variant vectors; especially in the context of the gradient vector: **co-variant gradient** ("basis-independent gradient"). That's confusing because

---

<sup>6</sup>A second origin is to use an entirely different convention to define co-variance without reference to coordinates at all: Assume vectors (including, e.g., basis vectors) transform *with* a linear transformation  $v \mapsto f v$ . Postulate additionally that scalars must be invariant. Then it follows that 1-forms must transform contra-variantly,  $g \in V^* \mapsto g \circ f^{-1} \in V^*$ , such that  $g(x) \in \mathbb{R} \mapsto g \circ f^{-1}(f x) = g(x)$  invariantly. If we now go down to coordinates, if  $w^\top$  are the coordinates of the 1-form  $g$ , and  $A$  the coordinates of the transform  $f$ , then  $w \mapsto A^\top w$  such that  $w^\top x \mapsto (A^\top w)^\top A x = w^\top x$  invariantly.

the co-variant gradient (e.g.  $A^{-1}\nabla f$ ) has actually contra-variant coordinates, just as (invariant) vectors.

**Ordinary gradient** descent of the form  $x \leftarrow x + \alpha\nabla f$  adds objects of different types: a basis-invariant vector  $x$  with a non-basis-invariant gradient  $\nabla f$ . Clearly, adding two such different types leads to an object who's transformation under coordinate transforms is strange—and indeed the ordinary gradient descent is not invariant under transformations.

### 3.4 Derivative Rules and Examples

Learn this:

**Identities 3.3.**

$$\frac{\partial}{\partial x} f(x)^\top g(x) = f(x)^\top \frac{\partial}{\partial x} g(x) + g(x)^\top \frac{\partial}{\partial x} f(x) \tag{13}$$

Note that using the 'gradient column' convention this reads

$$\nabla_x f(x)^\top g(x) = [\frac{\partial}{\partial x} g(x)]^\top f(x) + [\frac{\partial}{\partial x} f(x)]^\top g(x)$$

which I find impossible to remember, and mixes gradients-in-columns ( $\nabla$ ) with gradients-in-rows (the Jacobian).

We have:

$$\frac{\partial}{\partial x} [whatever]x = [whatever], \quad \text{if } whatever \text{ is indep. of } x \tag{14}$$

$$\frac{\partial}{\partial x} a^\top x = a^\top \tag{15}$$

$$\frac{\partial}{\partial x} Ax = A \tag{16}$$

$$\frac{\partial}{\partial x} (Ax - b)^\top (Cx - d) = (Ax - b)^\top C + (Cx - d)^\top A \tag{17}$$

$$\frac{\partial}{\partial x} x^\top Ax = x^\top A + x^\top A = 2x^\top A \tag{18}$$

$$\frac{\partial}{\partial x} \|x\| = \frac{\partial}{\partial x} (x^\top x)^{\frac{1}{2}} = \frac{1}{2} (x^\top x)^{-\frac{1}{2}} 2x^\top = \frac{1}{\|x\|} x^\top \tag{19}$$

$$\frac{\partial^2}{\partial x^2} (Ax + a)^\top C (Bx + b) = A^\top CB + B^\top C^\top A \tag{20}$$

Further useful identities:

**Identities 3.4** (Derivative Rules).

$$\frac{\partial}{\partial \theta} |A| = |A| \operatorname{tr}(A^{-1} \frac{\partial}{\partial \theta} A) \quad (21)$$

$$\frac{\partial}{\partial \theta} A^{-1} = -A^{-1} \left( \frac{\partial}{\partial \theta} A \right) A^{-1} \quad (22)$$

$$\frac{\partial}{\partial \theta} \operatorname{tr}(A) = \sum_i \frac{\partial}{\partial \theta} A_{ii} \quad (23)$$

VERY AWKWARD: Taking a derivative of a scalar value w.r.t. a matrix. Some write this again as a matrix:

$$\frac{\partial}{\partial X} a^\top X b = ab^\top \quad (24)$$

$$\frac{\partial}{\partial X} d/dX(a^\top X^\top C X b) = C^\top X ab^\top + C X ba^\top \quad (25)$$

$$\partial_X \operatorname{tr}(X) = \mathbf{I} \quad (26)$$

But I would recommend (in standard contexts) to never take a derivative w.r.t. a matrix. Instead, perhaps take the derivative w.r.t. a matrix element:  $\frac{\partial}{\partial X_{ij}} a^\top X b = a_i b_j$ .

**Identities 3.5** (Matrix Identities).

$$(A^{-1} + B^{-1})^{-1} = A(A+B)^{-1}B = B(A+B)^{-1}A \quad (27)$$

$$(A^{-1} - B^{-1})^{-1} = A(B-A)^{-1}B \quad (28)$$

$$(A + UBV)^{-1} = A^{-1} - A^{-1}U(B^{-1} + VA^{-1}U)^{-1}VA^{-1} \quad (29)$$

$$(A^{-1} + B^{-1})^{-1} = A - A(B+A)^{-1}A \quad (30)$$

$$(A + J^\top B J)^{-1} J^\top B = A^{-1} J^\top (B^{-1} + J A^{-1} J^\top)^{-1} \quad (31)$$

$$(A + J^\top B J)^{-1} A = \mathbf{I} - (A + J^\top B J)^{-1} J^\top B J \quad (32)$$

(100)=Woodbury; (102,103) holds for pos def  $A$  and  $B$ . See also the [matrix cookbook](#).

**3.4.1 GP regression**

An example from GP regression: The log-likelihood gradient w.r.t. a kernel hyperparameter:

$$\log P(y|X, b) = -\frac{1}{2} y^\top K^{-1} y - \frac{1}{2} \log |K| - \frac{n}{2} \log 2\pi \quad (33)$$

$$\text{where } K_{ij} = e^{-b(x_i - x_j)^2} + \sigma^2 \delta_{ij} \quad (34)$$

$$\frac{\partial}{\partial b} y^\top K^{-1} y = y^\top (-K^{-1} (\frac{\partial}{\partial b} K) K^{-1}) = y^\top (-K^{-1} A K^{-1}), \quad \text{with } A_{ij} = -(x_i - x_j)^2 e^{-b(x_i - x_j)} \quad (35)$$

$$\frac{\partial}{\partial b} \log |K| = \frac{1}{|K|} \frac{\partial}{\partial b} |K| = \frac{1}{|K|} |K| \operatorname{tr}(K^{-1} \frac{\partial}{\partial b} K) = \operatorname{tr}(K^{-1} A) \quad (36)$$

### 3.4.2 Logistic regression

An example from logistic regression: We have the loss gradient and want the Hessian:

$$\nabla_\beta L = X^\top (p - x) + 2\lambda \mathbf{I} \beta \quad (37)$$

$$\text{where } p_i = \sigma(x_i^\top \beta), \quad \sigma(z) = \frac{e^z}{1 + e^z}, \quad \sigma'(z) = \sigma(z) (1 - \sigma(z)) \quad (38)$$

$$\nabla_\beta^2 L = \frac{\partial}{\partial \beta} \nabla_\beta L = X^\top \frac{\partial}{\partial \beta} p + 2\lambda \quad (39)$$

$$\frac{\partial}{\partial \beta} p_i = p_i (1 - p_i) x_i^\top \quad (40)$$

$$\frac{\partial}{\partial \beta} p = \operatorname{diag}([p_i(1 - p_i)]_{i=1}^n) X = \operatorname{diag}(p \circ (1 - p)) X, \quad \text{where } \circ \text{ is the element-wise} \quad (41)$$

$$\nabla_\beta^2 L = X^\top \operatorname{diag}(p \circ (1 - p)) X + 2\lambda \mathbf{I} \quad (42)$$

## 3.5 Taylor expansion

In 1D, we have

$$f(x + v) \approx f(x) + f'(x)v + \frac{1}{2} f''(x)v^2 + \dots + \frac{1}{k!} f^{(k)}(x)v^k \quad (43)$$

For  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , we have

$$f(x + v) \approx f(x) + \nabla f(x)^\top v + \frac{1}{2} v^\top \nabla^2 f(x) v + \dots \quad (44)$$

For  $f : \mathbb{R}^n \rightarrow \mathbb{R}^d$  is more convenient to go back to index notation:

$$f_i(x + v) \approx f_i(x) + \sum_j \frac{\partial}{\partial x_j} f_i(x) v_j + \frac{1}{2} \sum_{jk} \frac{\partial^2}{\partial x_j \partial x_k} f_i(x) v_j v_k + \frac{1}{6} \sum_{jkl} \frac{\partial^3}{\partial x_j \partial x_k \partial x_l} f_i(x) v_j v_k v_l \quad (45)$$

### 3.6 Steepest descent and the covariant gradient vector

Let's define the steepest descent direction is the one where, when you make a step of length 1, you get the largest decrease of  $f$  in its linear (=1st order Taylor) approximation.

**Definition 3.6.** Given  $f : V \rightarrow \mathbb{R}$  and a norm  $\|x\|^2 = \langle x, x \rangle$  (or scalar product) defined on  $V$ , we define the **steepest descent** vector  $\delta^* \in V$  as the vector:

$$\delta^* = \underset{\delta}{\operatorname{argmin}} df|_x(\delta) \quad \text{s.t.} \quad \|\delta\|^2 = 1$$

Note that, for any choice of coordinates (and independent from a definition of a scalar product)  $df|_x(\delta) = \partial_x f(x)\delta = \nabla f(x)^\top \delta$ .

Clearly, if we have coordinates in which the norm is Euclidean then

$$\|\delta\|^2 = \delta^\top \delta \quad \Rightarrow \quad \delta^* \propto -\nabla f(x)$$

However, if we have coordinates in which the metric is non-Euclidean, we have:

**Theorem 3.1** (Steepest Descent Direction (Covariant gradient)). For a general scalar product  $\langle v, w \rangle = v^\top G w$  (with metric tensor  $G$ ), the steepest descent direction is

$$\delta^* \propto -G^{-1} \nabla f(x) \tag{46}$$

Proof: Let  $G = B^\top B$  (Cholesky decomposition) and  $z = B\delta$

$$\begin{aligned} \delta^* &= \underset{\delta}{\operatorname{argmin}} \nabla f^\top \delta \quad \text{s.t.} \quad \delta^\top A \delta = 1 \\ &= B^{-1} \underset{z}{\operatorname{argmin}} \nabla f^\top B^{-1} z \quad \text{s.t.} \quad z^\top z = 1 \\ &\propto B^{-1} [-B^{-\top} \nabla f] = -A^{-1} \nabla f \end{aligned}$$

There is an important special case of this: Let  $p \in \Lambda^X$ , that is,  $p$  is a probability distribution over the space  $X$ . Further, let  $\theta \in \mathbb{R}^n$  and  $\theta \mapsto p(\theta)$  is some parameterization of the probability distribution. Then the derivative  $d_\theta p(\theta) \in T_p \Lambda^X$  is a vector in the tangent space of  $\Lambda^X$ . Now, for such vectors, for tangent vectors of the space of probability distributions, there is a generic metric, the **Fisher metric**: [TODO: move to 'probabilities' section]

### 3.7 Check your gradients numerically!

This is your typical work procedure when implementing a Machine Learning or AI-ish or Optimization kind of methods:





### 3.8.3 Minima

a) A core problem in Machine Learning: For  $\beta \in \mathbb{R}^d$ ,  $y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times d}$ , compute

$$\operatorname{argmin}_{\beta} \|y - X\beta\|^2 + \lambda \|\beta\|^2$$

b) A core problem in Robotics: For  $q, q_0 \in \mathbb{R}^n$ ,  $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^d$ ,  $y^* \in \mathbb{R}^d$  non-linear but smooth, compute

$$\operatorname{argmin}_q \|\phi(q) - y^*\|_C^2 + \|q - q_0\|_W^2$$

### 3.8.4 Finite Difference Gradient Checking

a) Implement the following pseudo code for empirical gradient checking:

---

**Input:**  $x \in \mathbb{R}^n$ , function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , function  $df : \mathbb{R}^n \rightarrow \mathbb{R}^{d \times n}$   
 1: initialize  $\hat{J} \in \mathbb{R}^{d \times n}$ , and  $\epsilon = 10^{-6}$   
 2: **for**  $i = 1 : n$  **do**  
 3:      $\hat{J}_{:,i} = [f(x + \epsilon e_i) - f(x - \epsilon e_i)]/2\epsilon$                      // assigns the  $i$ th column of  $\hat{J}$   
 4: **end for**  
 5: if  $\|\hat{J} - df(x)\|_{\infty} < 10^{-4}$  return true; else false

---

Here  $e_i$  is the  $i$ th standard basis vector in  $\mathbb{R}^n$ .

b) Test this for

- $f : x \mapsto Ax$ ,  $df : x \mapsto A$ ,  $x \sim \text{@randn}(n,1)$  @  $(N(0, 1))$  in Matlab and  $A \sim \text{@randn}(n,n)$  @
- $f : x \mapsto x^{\top}x$ ,  $df : x \mapsto 2x^{\top}$ ,  $x \sim \text{@randn}(n,1)$  @

### 3.8.5 “Backprop” in a Neural Net

Consider the function

$$f : \mathbb{R}^{h_0} \rightarrow \mathbb{R}^{h_3}, \quad f(x_0) = W_2 \sigma(W_1 \sigma(W_0 x_0))$$

where  $W_i \in \mathbb{R}^{h_{i+1} \times h_i}$  and  $\sigma(z) = 1/(e^{-z} + 1)$  is the sigmoid function, which here is applied element-wise. Choose  $(h_0, h_1, h_2, h_3) = (5, 10, 10, 2)$ . Note that the function can also be written as the “chain”:

$$x_0 \mapsto z_1 = W_0 x_0 \mapsto x_1 = \sigma(z_1) \mapsto z_2 = W_1 x_1 \mapsto x_2 = \sigma(z_2) \mapsto f = W_2 x_2$$

For this exercise, let’s focus on the Jacobian of  $f$ , which analytically is given as

$$\frac{\partial f}{\partial x_0} = \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial z_2} \frac{\partial z_2}{\partial x_1} \frac{\partial x_1}{\partial z_1} \frac{\partial z_1}{\partial x_0}$$

The Jacobians in this are:

$$\frac{\partial x_l}{\partial z_l} = \text{diag}(x_l \circ (1 - x_l)), \quad \frac{\partial z_{l+1}}{\partial x_l} = W_l, \quad \frac{\partial f}{\partial x_2} = W_2$$

a) Now, write code to implement  $f$  and  $\partial f$  for arbitrary  $(W_0, \dots, W_2)$ .

b) Sample random  $(W_0, \dots, W_2)$  using `@randn@` and a random  $x_0$ , and check the implemented Jacobien by comparing to the finite difference approximation.

Debugging Tip: If your first try does not work right away, the typical approach to debug is to “comment out” parts of your function  $f$  and  $df$ . For instance, start with testing  $f(x) = W_0 x_0$ ; then test  $f(x) = \sigma(W_0 x_0)$ ; then  $f(x) = W_1 \sigma(W_0 x_0)$ ; then I’m sure all bugs are found.

### 3.8.6 Logistic Regression Gradient & Hessian

Consider the function

$$L : \mathbb{R}^d \rightarrow \mathbb{R} : L(\beta) = - \sum_{i=1}^n \left[ y_i \log \sigma(x_i^\top \beta) + (1 - y_i) \log [1 - \sigma(x_i^\top \beta)] \right] + \lambda \beta^\top \beta$$

where  $x_i \in \mathbb{R}^d$  is the  $i$ th row of a matrix  $X \in \mathbb{R}^{n \times d}$ , and  $y \in \{0, 1\}^n$  is a *binary vector* (There is the word I never wanted to say! ;-))

The gradient is

$$\frac{\partial L(\beta)}{\partial \beta} = (p - y)^\top X + 2\lambda \beta^\top, \quad p = \sigma(X\beta)$$

and the hessian

$$\nabla^2 L(\beta) = X^\top \text{diag}(p \circ (1 - p)) X + 2\lambda I$$

Implement the function, gradient and hessian. Generate a random matrix  $X$  (using `@randn@`) and random  $y$  (using `@randi(2,n,1)-1@`) and test using finite differences that  $\partial L$  is the Jacobian of  $L$ , and that  $\nabla^2 L$  is the Jacobian of  $\nabla L = \partial L^\top$ .

## 4 Optimization

### 4.1 The general optimization problem – a mathematical program

**Definition 4.1.** Let  $x \in \mathbb{R}^n$ ,  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $h : \mathbb{R}^n \rightarrow \mathbb{R}^l$ . An optimization problem, or *mathematical program*, is

$$\min_x f(x)$$

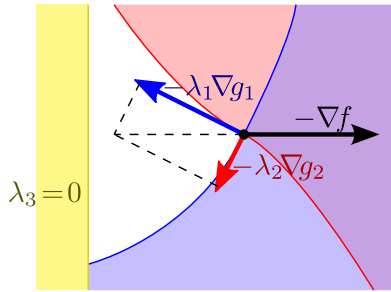


Figure 3: 2D example:  $f(x, y) = -x$ , pulling constantly to the right; three inequality constraints, two active, one inactive. The “pull/push” vectors fulfil the stationarity condition  $\nabla f + \lambda_1 \nabla g_1 + \lambda_2 \nabla g_2 = 0$ .

$$\text{s.t. } g(x) \leq 0, \quad h(x) = 0$$

We typically at least assume  $f, g, h$  to be differentiable or smooth (even if the solver does not have access to gradient or hessian, see below).

Get an intuition about this problem formulation by considering the following examples. Always discuss where is the optimum, and at the optimum, how the objective  $f$  pulls at the point, while the constraints  $g$  or  $h$  push against it.

- A 1D example:  $x \in \mathbb{R}$ ,  $h(x) = \sin(x)$ ,  $g(x) = x^2/4 - 1$ ,  $f$  some non-linear function.
- 2D example:  $f(x, y) = x$  (intuition: the objective is constantly pulling to the left),  $h(x, y) = 0$  is some non-linear path in the plane  $\rightarrow$  the optimum is at the left-most tangent-point of  $h$ . Tangent-point means that the tangent of  $h$  is vertical.  $h$  pushes to the right (always orthogonal to the zero-line of  $h$ ).
- 2D example:  $f(x, y) = x$ ,  $g(x, y) = y - x^2 - 1$ . The zero-line of  $g$  is a parabola towards the right. The objective  $f$  pulls into this parabola; the optimum is in the ‘bottom’ of the parabola, at  $(1, 0)$ .
- 2D example:  $f(x, y) = x$ ,  $g(x, y) = x^2 + y^2 - 1$ . The zero-line of  $g$  is a circle. The objective  $f$  pulls to the left; the optimum is at the left tangent-point of the circle, at  $(-1, 0)$ .
- Figure 3

## 4.2 The KKT conditions

**Theorem 4.1** (Karush-Kuhn-Tucker conditions). Given a mathematical program,

$$x \text{ optimal} \quad \Rightarrow \quad \exists \lambda \in \mathbb{R}^m, \kappa \in \mathbb{R}^l \quad \text{s.t.}$$

$$\nabla f(x) + \sum_{i=1}^m \lambda_i \nabla g_i(x) + \sum_{j=1}^l \kappa_j \nabla h_j(x) = 0 \quad (\text{stationarity}) \quad (47)$$

$$\forall_j : h_j(x) = 0, \quad \forall_i : g_i(x) \leq 0 \quad (\text{primal feasibility}) \quad (48)$$

$$\forall_i : \lambda_i \geq 0 \quad (\text{dual feasibility}) \quad (49)$$

$$\forall_i : \lambda_i g_i(x) = 0 \quad (\text{complementarity}) \quad (50)$$

Note that these are, in general, only necessary conditions. Only in special cases, e.g. convex, these are also sufficient.

These conditions should be intuitive in the previous examples:

- *The first condition describes the “force balance” of the objective pulling and the active constraints pushing back.* The existence of dual parameters  $\lambda, \kappa$  could implicitly be expressed by stating

$$\nabla f(x) \in \text{span}(\{\nabla g_{1..m}, \nabla h_{1..l}\})$$

The specific values of  $\lambda$  and  $\kappa$  tell us, how strongly the constraints push against the objective, e.g.,  $\lambda_i |\nabla g_i|$  is the force exerted by the  $i$ th inequality.

- *The fourth condition very elegantly describes the logic of inequality constraints being either active ( $\lambda_i > 0, g_i = 0$ ) or inactive ( $\lambda_i = 0, g_i \leq 0$ ).* Intuitively it says: An inequality can only push at the boundary, where  $g_i = 0$ , but not inside the feasible region, where  $g_i < 0$ . The trick of using the equation  $\lambda_i g_i = 0$  to express this logic is beautiful, especially when later we discuss a case which relaxes this strict logic to  $\lambda_i g_i = \mu$  for some small  $\mu$ —which roughly means that inequalities may push a little also inside the feasible region.
- Special case  $m = l = 0$  (no constraints). The first condition is just the usual  $\nabla f(x) = 0$ .
- Discuss the previous examples as special cases; and how the force balance is met.

### 4.3 Unconstrained problems to tackle a constrained problem

Assume you’d know about basic unconstrained optimization methods (like standard gradient descent or Newton method) but nothing about constrained optimization methods. How would you solve a constrained problem? Well, I think you’d very quickly have the idea to introduce extra cost terms for violation of constraints— a million people have had this idea and successfully applied it in practice.

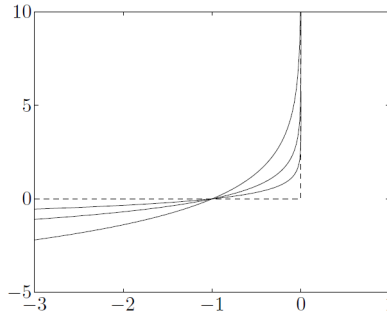


Figure 4: The function  $-\mu \log(-g)$  (with  $g$  on the “ $x$ -axis”) for various  $\mu$ . This is always undefined (“ $\infty$ ”) for  $g > 0$ . For  $\mu \rightarrow 0$  this becomes the hard step function.

In the following we define a new cost function  $F(x)$ , which include the objective  $f(x)$  and some extra terms.

**Definition 4.2** (Log barrier, squared penalty, Lagrangian, Augmented Lagrangian).

$$F_{sp}(x; \nu, \mu) = f(x) + \nu \sum_j h_j(x)^2 + \mu \sum_i [g_i(x) > 0] g_i(x)^2 \quad (\text{sqr. penalty})$$

$$F_{lb}(x; \mu) = f(x) - \mu \sum_i \log(-g_i(x)) \quad (\text{log barrier})$$

$$L(x, \lambda, \kappa) = f(x) + \sum_j \kappa_j h_j(x) + \sum_i \lambda_i g_i(x) \quad (\text{Lagrangian})$$

$$\begin{aligned} \hat{L}(x) = f(x) + \sum_j \kappa_j h_j(x) + \sum_i \lambda_i g_i(x) + \\ + \nu \sum_j h_j(x)^2 + \mu \sum_i [g_i(x) > 0] g_i(x)^2 \quad (\text{Aug. Lag.}) \end{aligned}$$

- The squared penalty method is straight-forward if we have an algorithm to minimize  $F(x)$ . We initialize  $\nu = \mu = 1$ , minimize  $F(x)$ , then increase  $\nu, \mu$  (multiply with a number  $> 1$ ) and iterate. For  $\nu, \mu \rightarrow \infty$  we retrieve the optimum.
- The log barrier method (see Fig. 4) does exactly the same, except that we decrease  $\mu$  towards zero (multiply with a number  $< 1$  in each iteration). Note that we need a feasible initialization  $x_0$ , because otherwise the barriers are ill-defined! The whole algorithm will keep the temporary solutions always *inside* the feasible regions (because the barriers push away from the constraints). That’s why it is also called **interior point method**.

- The Lagrangian is a function  $L(x, \lambda, \kappa)$  which has the gradient

$$\nabla L(x, \lambda, \kappa) = \nabla f(x) + \sum_i \lambda_i \nabla g_i(x) + \sum_j \kappa_j \nabla h_j(x) .$$

That is,  $\nabla L(x, \lambda, \kappa) = 0$  is our first KKT condition! In that sense, the additional terms in the Lagrangian generate the push forces of the constraints. If we knew the correct  $\lambda$ 's and  $\kappa$ 's beforehand, then we could indeed find the optimal  $x$  simply by the unconstrained problem  $\min_x L(x, \lambda, \kappa)$ .

- The Augmented Lagrangian  $\hat{L}$  is a function that includes both, squared penalties, and Lagrangian terms that push proportional to  $\lambda, \kappa$ . The Augmented Lagrangian method is exactly such an iterative algorithm that, while running, figures out how strongly we need to push to ensure that the final solution is *exactly* on the constraints, where all squared penalties will anyway be zero. It never needs to increase  $\nu, \mu$  and still converges to the correct solution.

### 4.3.1 Log Barrier

[[TODO: include equality constraints]]

Given

$$F(x; \mu) = f(x) - \mu \sum_i \log(-g_i(x))$$

the optimality condition is

$$\nabla F(x) = 0 \quad \Rightarrow \quad \nabla f(x) - \sum_i \frac{\mu}{g_i(x)} \nabla g_i(x) = 0 \tag{51}$$

$$\Leftrightarrow \quad \nabla f(x) + \sum_i \lambda_i \nabla g_i(x) = 0, \quad \lambda_i g_i(x) = -\mu \tag{52}$$

where we defined  $\lambda_i = -\mu/g_i(x)$ , which is guaranteed  $\geq 0$  if  $g_i \leq 0$ .

So,  $\nabla F(x) = 0$  is equivalent to the **modified (=approximate) KKT conditions**. For  $\mu \rightarrow 0$  we converge to the exact KKT conditions with strict complementarity. For finite  $\mu$ , the complementarity is relaxed: inequalities may push also inside the feasible region.

- Central path

### 4.3.2 Augmented Lagrangian\*

This is not a main-stream algorithm, but I love it. See ?.

In the Augmented Lagrangian  $\hat{L}$ , the solver has two types of knobs to tune: the strengths of the penalties  $\nu, \mu$  and the strengths of the Lagrangian forces  $\lambda, \kappa$ . The trick is conceptually easy:

- Initially we set  $\lambda, \kappa = 0$  and  $\nu, \mu = 1$  (or some other constant). In the first iteration, the solver will find  $x' = \min_x \hat{L}(x)$ ; the objective  $f$  will pull into the penalizations.
- For the second iteration we then choose parameters  $\lambda, \kappa$  that try to avoid that we will be pulled into penalizations the next time. Let's update

$$\kappa_j \leftarrow \kappa_j + 2\mu h_j(x'), \quad \lambda_i \leftarrow \max(\lambda_i + 2\mu g_i(x'), 0).$$

Note that  $2\mu h_j(x')$  is the force (gradient) of the equality penalty at  $x'$ ; and  $\max(\lambda_i + 2\mu g_i(x'), 0)$  is the force of the inequality constraint at  $x'$ . So what this update does is: it analyzes the forces exerted by the penalties, and translates them to forces exerted by the Lagrange terms for the next iteration. It tries to trade the penalizations for the Lagrange terms.

More rigorously, observe that, if  $f, g, h$  are linear and the same constraints are active in two consecutive iterations, then this update will guarantee that all penalty terms are zero in the second iteration, and therefore the solution fulfils the first KKT condition. ?

## 4.4 The Lagrangian

### 4.4.1 How the Lagrangian relates to the KKT conditions

The Lagrangian  $L(x, \kappa, \lambda) = f + \kappa^\top h + \lambda^\top g$  has a number of properties that relates it to the KKT conditions:

- (i) Requiring a zero- $x$ -gradient,  $\nabla_x L = 0$ , implies the *1st KKT condition*.
- (ii) Requiring a zero- $\kappa$ -gradient,  $0 = \nabla_\kappa L = h$ , implies primal feasibility (the *2nd KKT condition*) w.r.t. the equality constraints.
- (iii) Requiring that  $L$  is maximized w.r.t.  $\lambda \geq 0$  is related to the remaining 2nd and 4th KKT conditions:

$$\max_{\lambda \geq 0} L(x, \lambda) = \begin{cases} f(x) & \text{if } g(x) \leq 0 \\ \infty & \text{otherwise} \end{cases} \quad (53)$$

$$\lambda = \operatorname{argmax}_{\lambda \geq 0} L(x, \lambda) \Rightarrow \begin{cases} \lambda_i = 0 & \text{if } g_i(x) < 0 \\ 0 = \nabla_{\lambda_i} L(x, \lambda) = g_i(x) & \text{otherwise} \end{cases} \quad (54)$$

This implies either  $(\lambda_i = 0 \wedge g_i(x) < 0)$  or  $g_i(x) = 0$ , which is equivalent to the *complementarity* and *primal feasibility* for inequalities.

These three facts show how tightly the Lagrangian is related to the KKT conditions. To simplify the discussion let us assume only inequality constraints from now on.



Fact (i) tells us that if we  $\min_x L(x, \lambda)$ , we reproduce the 1st KKT condition. Fact (iii) tells us that if we  $\max_{\lambda \geq 0} L(x, \lambda)$ , we reproduce the remaining KKT conditions. Therefore, the optimal primal-dual solution  $(x^*, \lambda^*)$  can be characterized as a **saddle point of the Lagrangian**. Finding the saddle point can be written in two ways:

**Definition 4.3** (primal and dual problem).

$$\min_x \max_{\lambda \geq 0} L(x, \lambda) \qquad \text{(primal problem)} \qquad (55)$$

$$\max_{\lambda \geq 0} \underbrace{\min_x L(x, \lambda)}_{l(\lambda)} \qquad \text{(dual problem)} \qquad (56)$$

Convince yourself, using (53), that the first expression is indeed the original primal problem  $\min_x f(x) \text{ s.t. } g(x) \leq 0$ .

**What can we draw out of this?** The KKT conditions state that, at an optimum, *there exist some*  $\lambda, \kappa$ . This existence statement is not very helpful to actually find them. In contrast, the Lagrangian tells us directly how the dual parameters can be found: by maximizing w.r.t. them. This can be exploited in several ways:

#### 4.4.2 Solving mathematical programs analytically, on paper.

Consider the problem

$$\min_{x \in \mathbb{R}^2} x^2 \text{ s.t. } x_1 + x_2 = 1.$$

We can find the solution analytically via the Lagrangian:

$$\begin{aligned} L(x, \kappa) &= x^2 + \kappa(x_1 + x_2 - 1) \\ 0 = \nabla_x L(x, \kappa) &= 2x + \kappa \begin{pmatrix} 1 \\ 1 \end{pmatrix} \Rightarrow x_1 = x_2 = -\kappa/2 \\ 0 = \nabla_\kappa L(x, \kappa) &= x_1 + x_2 - 1 = -\kappa/2 - \kappa/2 - 1 \Rightarrow \kappa = -1 \\ \Rightarrow x_1 = x_2 &= 1/2 \end{aligned}$$

Here we first formulated the Lagrangian. In this context,  $\kappa$  is often called **Lagrange multiplier**, but I prefer the term *dual variable*. Then we find a saddle point of  $L$  by requiring  $0 = \nabla_x L(x, \kappa)$ ,  $0 = \nabla_\kappa L(x, \kappa)$ . If we want to solve a problem with an inequality constrained, we do the same calculus for both cases: 1) the constraint is active (handled like an equality constrained), and 2) the constrained is inactive. Then we check if the inactive case solution is feasible. Note that if we have  $m$  inequality constraints we have to analytically evaluate every combination of constraints being active/inactive—which are  $2^m$  cases. This already hints at the fact that a real difficulty in solving mathematical programs is to find out which inequality constraints are active or inactive. In fact, if we knew this a priori, everything would reduce to an equality constrained problem, which is much easier to solve.

#### 4.4.3 Solving the dual problem ,instead of the primal.

In some cases the dual function  $l(\lambda) = \min_x L(x, \lambda)$  can analytically be derived. In this case it makes very much sense to try solving the dual problem instead of the primal, especially because 1) the dual problem is guaranteed to be convex, 2) and the inequality constraints of the dual problem are very simple: just  $\lambda \geq 0$ . Such inequality constraints are called **bound constraints** and can be handled with specialized methods, e.g. simple clipping or projection of steps.

#### 4.4.4 Finding the “saddle point” directly with a primal-dual Newton method.

In basic unconstrained optimization an efficient way to find an optimum (minimum or maximum) is to find a point where the gradient is zero with a Newton method. At saddle points all gradients are also zero. So, to find a saddle point of the Lagrangian we can equally use a Newton method that seeks for roots of the gradient. Note that such a Newton method optimizes in the joint **primal-dual space** of  $(x, \kappa, \lambda)$ .

In the case of inequalities, the zero-gradients view is over-simplified: While facts (i) and (ii) characterize a saddle point in terms of zero gradients, fact (iii) makes this more precise to handle the inequality case. For this reason it is actually easier to describe the primal-dual Newton method directly in terms of the KKT conditions: We seek a point  $(x, \lambda)$ , with  $\lambda \geq 0$ , that solves the equation system

$$\nabla_x f(x) + \lambda^\top \partial_x g = 0 \quad (57)$$

$$\text{diag}(\lambda)g(x) + \mu \mathbf{1}_m = 0 \quad (58)$$

Note that the first equation is the 1st KKT, and the second is the *approximate* 4th KKT with log barrier parameter  $\mu$ . These two equations correctly reflect the saddle point properties (facts (i) and (iii) above). We define

$$r(x, \lambda) = \begin{pmatrix} \nabla f(x) + \lambda^\top \partial_x g \\ \text{diag}(\lambda)g(x) + \mu \mathbf{1}_m \end{pmatrix} \quad (59)$$

and use the Newton method

$$\begin{pmatrix} x \\ \lambda \end{pmatrix} \leftarrow \begin{pmatrix} x \\ \lambda \end{pmatrix} - \alpha \nabla r(x, \lambda)^{-1} r(x, \lambda)$$

to find the root  $r(x, \lambda) = 0$  ( $\alpha$  is the stepsize). We have

$$\nabla r(x, \lambda) = \begin{pmatrix} \nabla^2 f(x) + \sum_i \lambda_i \nabla^2 g_i(x) & \nabla g(x)^\top \\ \text{diag}(\lambda) \nabla g(x) & \text{diag}(g(x)) \end{pmatrix} \in \mathbb{R}^{(n+m) \times (n+m)} \quad (60)$$

Note that this method uses the Hessians  $\nabla^2 f$  and  $\nabla^2 g$ .

The above formulation allows for a duality gap  $\mu$ ; choose  $\mu = 0$  or consult Boyd how to update on the fly (sec 11.7.3)

The **feasibility constraints**  $g_i(x) \leq 0$  and  $\lambda_i \geq 0$  need to be handled explicitly by the root finder (the line search needs to ensure these constraints).

#### 4.4.5 Properties of the dual problem\*

Note that in general  $\min_x \max_y f(x, y) \neq \max_y \min_x f(x, y)$ . For example, in discrete domain  $x, y \in \{1, 2\}$ , let  $f(1, 1) = 1, f(1, 2) = 3, f(2, 1) = 4, f(2, 2) = 2$ , and  $\min_x f(x, y) = (1, 2)$  and  $\max_y f(x, y) = (3, 4)$ . Therefore, the dual problem is in general not equivalent to the primal.

We note, without proof, that the dual problem  $\max_{\lambda \geq 0} l(\lambda)$  is convex even if the primal is non-convex. (The dual function  $l(\lambda)$  is concave, and the constraint  $\lambda \geq 0$  convex.) But note that  $l(\lambda)$  is itself defined as the result of a generally non-convex optimization problem  $\min_x L(x, \lambda)$ .

The dual function is, for  $\lambda \geq 0$ , a lower bound

$$l(\lambda) = \min_x L(x, \lambda) \leq \left[ \min_x f(x) \quad \text{s.t.} \quad g(x) \leq 0 \right].$$

And consequently

$$\text{(dual)} \quad \max_{\lambda \geq 0} \min_x L(x, \lambda) \leq \min_x \max_{\lambda \geq 0} L(x, \lambda) \quad \text{(primal)}$$

We say **strong duality** holds iff

$$\max_{\lambda \geq 0} \min_x L(x, \lambda) = \min_x \max_{\lambda \geq 0} L(x, \lambda)$$

If the primal is convex, and there exist an interior point

$$\exists x : \forall_i : g_i(x) < 0$$

(which is called **Slater condition**), then we have *strong duality*

#### 4.4.6 Log barrier again

Let  $x^*(\mu) = \min_x F(x; \mu)$  be the central path. At such a  $x^*$  we may define, as above,  $\lambda_i = -\mu/g_i(x)$ . We note that  $\lambda \geq 0$  (dual feasible), as well that  $x^*(\mu)$  minimizes the Lagrangian  $L(x, \lambda)$ !! This is because,

$$0 = \nabla F(x, \mu) = \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla g_i(x) = \nabla L(x, \lambda). \quad (61)$$

Therefore,  $x^*$  is actually the solution to  $\min_x L(x, \lambda)$ , which defines the dual function. We have

$$l(\lambda) = \min_x L(x, \lambda) = f(x^*) + \sum_{i=1}^m \lambda_i g_i(x^*) = f(x^*) - m\mu. \quad (62)$$

That is,  $m\mu$  is the duality gap between the (suboptimal) primal value  $f(x^*)$  and the dual value  $l(\lambda)$ . Further, given that the dual function is a lower bound,  $l(\lambda) \leq p^*$ , where  $p^* = \min_x f(x)$  s.t.  $g(x) \leq 0$  is the optimal primal value, we have

$$f(x^*) - p^* \leq m\mu$$

which again verifies that  $\mu \rightarrow 0$  will let the primal value converge to the optimal.

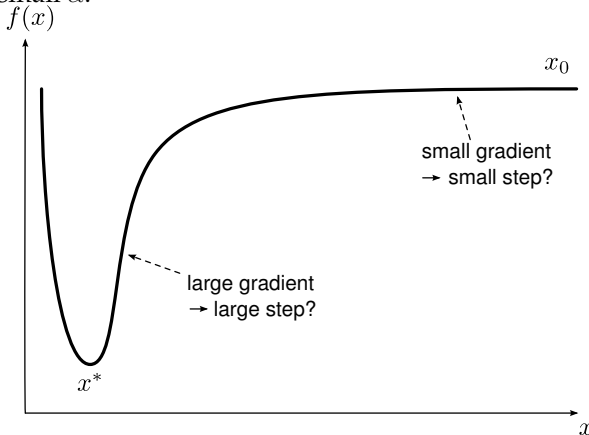
## 4.5 Downhill algorithms for unconstrained optimization

We discuss here algorithms that have one goal: walk downhill as quickly as possible. These target at efficiently finding local optima—in contrast to *global optimization* methods, which try to find the global optimum.

For such downhill walkers, there are two essential things to discuss: the stepsize and the step direction. When discussing the stepsize we'll hit on topics like backtracking line search, the Wolfe conditions and its implications in a basic convergence proof. The discussion of the step direction will very much circle around Newton's method and thereby also cover topics like quasi-Newton methods (BFGS), Gauss-Newton, covariant and conjugate gradients.

### 4.5.1 Why you shouldn't trust the magnitude of the gradient

Consider the following 1D function and naive gradient descent  $x \leftarrow x - \alpha \nabla f$  for some fixed and small  $\alpha$ .



- In plateaus we'd make small steps, at steep slopes (here close to the optimum) we make huge steps, very likely overstepping the optimum. In fact, for some  $\alpha$  the algorithm might indefinitely loop a non-sensical sequence of very slowly walking left on the plateau, then accelerating, eventually overstepping the

optimum, then being thrown back far to the right again because of the huge negative gradient on the left.

- Generally, never trust an algorithm that depends on the scaling—or choice of units—of a function! An optimization algorithm should be invariant on whether you measure costs in cents or Euros! Naive gradient descent  $x \leftarrow x - \alpha \nabla f$  is not!

As a conclusion, the gradient  $\nabla f$  gives a reasonable descent direction, but its magnitude is really arbitrary and no good indication of a good stepsize. Therefore, it often makes sense to just compute the step direction

$$\delta = -\frac{1}{|\nabla f(x)|} \nabla f(x)$$

and iterate  $x \leftarrow x - \alpha \delta$  for some appropriate stepsize.

#### 4.5.2 Ensuring monotone and sufficient decrease: Backtracking line search, Wolfe conditions, & convergence

The first idea is simple: If a step would increase the objective value, reduce the stepsize. We typically use multiplicative stepsize adaptations: Reduce  $\alpha \leftarrow \varrho_{\alpha}^{-} \alpha$  with  $\varrho_{\alpha}^{-} \approx 0.5$ ; and increase  $\alpha \leftarrow \varrho_{\alpha}^{+} \alpha$  with  $\varrho_{\alpha}^{+} \approx 1.2$ . A simple monotone gradient descent algorithm reads as follows (the blue part is explained later). Here, the step

---

##### Algorithm 2 Plain gradient descent

---

**Input:** initial  $x \in \mathbb{R}^n$ , functions  $f(x)$  and  $\nabla f(x)$ , tolerance  $\theta$ , parameters (defaults:  $\varrho_{\alpha}^{+} = 1.2, \varrho_{\alpha}^{-} = 0.5, \delta_{\max} = \infty, \varrho_{\text{ls}} = 0.01$ )

```

1: initialize stepsize  $\alpha = 1$ 
2: repeat
3:    $\delta \leftarrow -\frac{\nabla f(x)}{|\nabla f(x)|}$  // (alternative:  $\delta = -\nabla f(x)$ )
4:   while  $f(x + \alpha\delta) > f(x) + \varrho_{\text{ls}} \nabla f(x)^{\top}(\alpha\delta)$  do // line search
5:      $\alpha \leftarrow \varrho_{\alpha}^{-} \alpha$  // decrease stepsize
6:   end while
7:    $x \leftarrow x + \alpha\delta$ 
8:    $\alpha \leftarrow \min\{\varrho_{\alpha}^{+} \alpha, \delta_{\max}\}$  // increase stepsize
9: until  $|\alpha\delta| < \theta$  // perhaps for 10 iterations in sequence
```

---

vector  $\delta$  is always normalized and  $\alpha$  is adapted on the fly; decreasing when  $f(x + \alpha\delta)$  is not sufficiently smaller than  $f(x)$ .

This *sufficiently smaller* is described by the blue part and is called the (1st) **Wolfe condition**

$$f(x + \alpha\delta) > f(x) + \varrho_{\text{ls}} \nabla f(x)^{\top}(\alpha\delta).$$

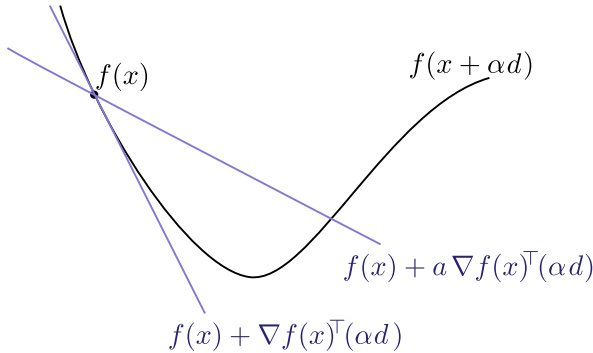


Figure 5: The 1st Wolfe condition:  $f(x) + \nabla f(x)^T(\alpha d)$  is the tangent, which describes the expected decrease of  $f(x + \alpha d)$  if  $f$  was linear. We cannot expect this to be the case; so  $f(x + \alpha d) > f(x) + \varrho_{ls} \nabla f(x)^T(\alpha d)$  weakens this condition.

Figure 5 illustrates this. Note that  $\nabla f(x)^T(\alpha d)$  is a negative value and describes how much the objective would decrease if  $f$  was linear. But  $f$  is of course not linear; we cannot expect that a step would really decrease  $f$  that much. Instead we require that it decreases by a fraction of this expectation.  $\varrho_{ls}$  describes this fraction and is typically chosen very moderate, e.g.  $\varrho_{ls} \in [0.01, 0.1]$ . So, the Wolfe conditions requires that  $f$  decreases by the  $\varrho_{ls}$ -fraction of what it would decrease if  $f$  was linear. Note that for  $\alpha \rightarrow 0$  the Wolfe condition will always be fulfilled for smooth functions, because  $f$  “becomes locally linear”.

In the exercises you prove that any algorithms that, in each iteration, finds a new point that fulfills the Wolfe condition will converge with a convergence rate  $(1 - \frac{2m\varrho_{ls}\varrho_{\alpha}}{M})$  if  $m > 0$  and  $M > m$  are lower and upper bounds on the hessian eigenvalues.

There is a second Wolfe condition,

$$|\nabla f(x + \alpha d)^T \delta| \leq b |\nabla f(x)^T \delta|,$$

which states that the gradient magnitude should have decreased sufficiently. We do not use it much.

### 4.5.3 The Newton direction

We already discussed the *steepest descent direction*  $-G^{-1} \nabla f(x)$  if  $G$  is a metric tensor. Let’s keep this in mind!

The original Newton method is a method to find the **root** (that is, zero point) of a function  $f(x)$ . In 1D it iterates  $x \leftarrow x - \frac{f(x)}{f'(x)}$ , that is, it uses the gradient  $f'$  to estimate where the function might cross the  $x$ -axis. To find an minimum of maximum of  $f$

we want to find the root of its gradient. For  $x \in \mathbb{R}^n$  the Newton method iterates

$$x \leftarrow x - \nabla^2 f(x)^{-1} \nabla f(x).$$

Note that the Newton step  $\delta = -\nabla^2 f(x)^{-1} \nabla f(x)$  is the solution to

$$\min_{\delta} \left[ f(x) + \nabla f(x)^{\top} \delta + \frac{1}{2} \delta^{\top} \nabla^2 f(x) \delta \right].$$

So the Newton method can also be viewed as 1) compute the 2nd-order Taylor approximation to  $f$  at  $x$ , and 2) jump to the optimum of this approximation.

Note:

- If  $f$  is just a 2nd-order polynomial, the Newton method will jump to the optimum in just one step.
- *Unlike the gradient magnitude*  $|\nabla f(x)|$ , the magnitude of the Newton step  $\delta$  is very meaningful. It is scale invariant! If you'd rescale  $f$  (trade cents by Euros),  $\delta$  is unchanged.  $|\delta|$  is the distance to the optimum of the 2nd-order Taylor.
- *Unlike the gradient*  $\nabla f(x)$ , the Newton step  $\delta$  is truly a contra-variant vector! The vector itself is invariant under coordinate transformations; the coordinate vector  $\delta$  transforms contra-variant, as it is supposed to for vector coordinates.
- **The hessian as metric, and the Newton step as steepest descent:** Assume that the hessian  $H = \nabla^2 f(x)$  is pos-def. Then it fulfils all necessary conditions to define a scalar product  $\langle v, w \rangle = \sum_{ij} v_i w_j H_{ij}$ , where  $H$  plays the role of the metric tensor. If  $H$  was the space's metric, then the steepest descent direction is  $-H^{-1} \nabla f(x)$ —which is the Newton direction!

Another way to understand the same: In the 2nd-order Taylor approximation  $f(x + \delta) \approx f(x) + \nabla f(x)^{\top} \delta + \frac{1}{2} \delta^{\top} H \delta$  the Hessian plays the role of a metric tensor. Or: we may think of the function  $f$  as being an isometric parabola  $f(x + \delta) \propto \langle \delta, \delta \rangle$ , but we've chosen coordinates where  $\langle v, v \rangle = v^{\top} H v$  and the parabola seems squeezed.

Note that this discussion only holds for pos-dev hessian.

Here a full Newton method

Some notes on the  $\lambda$ :

- The first line chooses  $\lambda$  to ensure that  $(\nabla^2 f(x) + \lambda \mathbf{I})$  is indeed pos-dev—and a Newton step actually decreases  $f$  instead of seeking for a maximum. Clearly there would be other options: instead of adding to all eigenvalues we could only set the negative ones to some  $\lambda > 0$ .

**Algorithm 3** Newton method

---

**Input:** initial  $x \in \mathbb{R}^n$ , functions  $f(x), \nabla f(x), \nabla^2 f(x)$ , tolerance  $\theta$ , parameters (defaults:  $\varrho_\alpha^+ = 1.2, \varrho_\alpha^- = 0.5, \varrho_\lambda^+ = \varrho_\lambda^- = 1, \varrho_{ls} = 0.01$ )

- 1: initialize stepsize  $\alpha = 1$
- 2: **repeat**
- 3:   choose  $\lambda > -\max \text{eig}(\nabla^2 f(x))$
- 4:   compute  $\delta$  to solve  $(\nabla^2 f(x) + \lambda \mathbf{I}) \delta = -\nabla f(x)$
- 5:   **while**  $f(x + \alpha\delta) > f(x) + \varrho_{ls} \nabla f(x)^\top (\alpha\delta)$  **do** // line search
- 6:      $\alpha \leftarrow \varrho_\alpha^- \alpha$  // decrease stepsize
- 7:     optionally:  $\lambda \leftarrow \varrho_\lambda^+ \lambda$  and recompute  $d$  // increase damping
- 8:   **end while**
- 9:    $x \leftarrow x + \alpha\delta$  // step is accepted
- 10:    $\alpha \leftarrow \min\{\varrho_\alpha^+ \alpha, 1\}$  // increase stepsize
- 11:   optionally:  $\lambda \leftarrow \varrho_\lambda^- \lambda$  // decrease damping
- 12: **until**  $\|\alpha\delta\|_\infty < \theta$

---

- $\delta$  solves the problem

$$\min_{\delta} \left[ \nabla f(x)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 f(x) \delta + \frac{1}{2} \lambda \delta^2 \right].$$

So, we added a squared potential  $\lambda \delta^2$  to the local 2nd-order Taylor approximation. This is like introducing a squared penalty for large steps!

- **Trust region method:** Let's consider a different mathematical program over the step:

$$\min_{\delta} \nabla f(x)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 f(x) \delta \quad \text{s.t.} \quad \delta^2 \leq \beta$$

This problem wants to find the minimum of the 2nd-order Taylor (like the Newton step), but constrained to a stepsize no larger than  $\beta$ . This  $\beta$  defines the *trust region*: The region in which we trust the 2nd-order Taylor to be a reasonable enough approximation.

Let's solve this as we learned to: Let's assume the inequality constraint is active. Then we have

$$L(\delta, \lambda) = \nabla f(x)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 f(x) \delta + \lambda(\delta^2 - \beta) \quad (63)$$

$$\nabla_{\delta} L(\delta, \lambda) = \nabla f(x)^\top + \delta^\top (\nabla^2 f(x) + 2\lambda \mathbf{I}) \quad (64)$$

Setting this to zero gives the step  $\delta = -(\nabla^2 f(x) + 2\lambda \mathbf{I})^{-1} \nabla f(x)$ .

Therefore, the  $\lambda$  can be viewed as **dual variable** of a trust region method. There is no analytic relation between  $\beta$  and  $\lambda$ ; we can't determine  $\lambda$  directly from  $\beta$ . We could use a constrained optimization method, like primal-dual Newton, or Augmented Lagrangian approach to solve for  $\lambda$  and  $\delta$ . Alternatively, we can always increase  $\lambda$  when the computed steps are too large, and



decrease if they are smaller than  $\beta$ —the Augmented Lagrangian update equations could be used to do such an update of  $\lambda$ .

- The  $\lambda\delta^2$  term can be interpreted as penalizing the *velocity* (stepsizes) of the algorithm. This is in analogy to a damping, such as “honey damping”, in physical dynamic systems. The parameter  $\lambda$  is therefore also called **damping** parameter. Such a damped (or regularized) least-squares methods is also called **Levenberg-Marquardt** method.
- For  $\lambda \rightarrow \infty$  the step direction  $\delta$  becomes aligned with the plain gradient direction  $-\nabla f(x)$ . This shows that for  $\lambda \rightarrow \infty$  the Hessian (and metric deformation of the space) becomes less important and instead we’ll walk orthogonal to the iso-lines of the function.
- The  $\lambda$  term makes the  $\delta$  non-scale-invariant!  $\delta$  is not anymore a proper contra-variant vector!

#### 4.5.4 Gauss-Newton: a super important special case

A special case that appears really a lot in intelligent systems is the Gauss-Newton case: Consider an objective function of the form

$$f(x) = \phi(x)^\top \phi(x) = \sum_i \phi_i(x)^2 \tag{65}$$

where we call  $\phi(x)$  the **cost features**. This is also called a **sum-of-squares** problem, here a sum of cost feature squares. We have

$$\nabla f(x) = 2 \frac{\partial}{\partial x} \phi(x)^\top \phi(x) \tag{66}$$

$$\nabla^2 f(x) = 2 \frac{\partial}{\partial x} \phi(x)^\top \frac{\partial}{\partial x} \phi(x) + 2 \phi(x)^\top \frac{\partial^2}{\partial x^2} \phi(x) \tag{67}$$

The Gauss-Newton method is the Newton method for  $f(x) = \phi(x)^\top \phi(x)$  while approximating  $\nabla^2 \phi(x) \approx 0$ . That is, it computes approximate Newton steps

$$\delta = -\left( \frac{\partial}{\partial x} \phi(x)^\top \frac{\partial}{\partial x} \phi(x) + \lambda \mathbf{I} \right)^{-1} \frac{\partial}{\partial x} \phi(x)^\top \phi(x) .$$

Note:

- The approximate Hessian  $2 \frac{\partial}{\partial x} \phi(x)^\top \frac{\partial}{\partial x} \phi(x)$  is **always semi-pos-def!** Therefore, no problems arise with negative hessian eigenvalues.
- The approximate Hessian only requires the first-order derivatives of the cost features. There is no need for computationally expensive Hessians of  $\phi$ .

- The objective  $f(x)$  can be interpreted as just the Euclidean norm  $f(\phi) = \phi^\top \phi$  but pulled back into the  $x$ -space. More precisely: Consider a mapping  $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and a general scalar product  $\langle \cdot, \cdot \rangle_\phi$  in the output space. In differential geometry there is the notion of a pull-back of a metric, that is, we define a scalar product  $\langle \cdot, \cdot \rangle_x$  in the input space as

$$\langle x, y \rangle_x = \langle d\phi(x), d\phi(y) \rangle_\phi \quad (68)$$

where  $d\phi$  is the differential of  $\phi$  (a  $\mathbb{R}^m$ -valued 1-form). Assuming  $\phi$ -coordinates such that the metric tensor of  $\langle \cdot, \cdot \rangle_\phi$  is Euclidean, we have

$$\langle x, x \rangle_x = \langle d\phi(x), d\phi(x) \rangle_\phi = \frac{\partial}{\partial x} \phi(x)^\top \frac{\partial}{\partial x} \phi(x) \quad (69)$$

and therefore, **the approximate Hessian is the pullback of a Euclidean cost feature metric**, and  $\langle x, x \rangle_x$  approximates the 2-order polynomial term of  $f(x)$ , with the non-constant (i.e., Riemannian) pull-back metric  $\langle \cdot, \cdot \rangle_x$ .

#### 4.5.5 Quasi-Newton & BFGS: approximating the hessian from gradient observations

To apply full Newton methods we need to be able to compute  $f(x)$ ,  $\nabla f(x)$ , and  $\nabla^2 f(x)$  for any  $x$ . However, sometimes, computing  $\nabla^2 f(x)$  is not possible, e.g., because we can't derive an analytic expression for  $\nabla^2 f(x)$ , or it would be too expensive to compute the hessian exactly, or even to store it in memory—especially in very high-dimensional spaces. In such cases it makes sense to approximate  $\nabla^2 f(x)$  or  $\nabla^2 f(x)^{-1}$  with a low-rank approximation.

Assume we have computed  $\nabla f(x_1)$  and  $\nabla f(x_2)$  at two different points  $x_1, x_2 \in \mathbb{R}^n$ . We define

$$y = \nabla f(x_2) - \nabla f(x_1), \quad \delta = x_2 - x_1.$$

From this we may wish to find some  $\nabla^2 f(x)$  or  $\nabla^2 f(x)^{-1}$  that fulfils

$$\nabla^2 f(x) \delta \stackrel{!}{=} y \quad \text{or} \quad \delta \stackrel{!}{=} \nabla^2 f(x)^{-1} y$$

The first equation is called *secant equation*. Here are guesses of  $\nabla^2 f(x)$  and  $\nabla^2 f(x)^{-1}$ :

$$\nabla^2 f(x) = \frac{y y^\top}{y^\top \delta} \quad \text{or} \quad f(x)^{-1} = \frac{\delta \delta^\top}{\delta^\top y}$$

Convince yourself that these choices fulfil the respective desired relation above. However, these choices are under-determined. There exist many alternative  $H$  or  $H^{-1}$  that would be consistent with the observed change in gradient. However, given our understanding of the structure of matrices it is clear that these choices are the lowest rank solutions, namely rank 1.

**Broyden-Fletcher-Goldfarb-Shanno (BFGS):** An optimization algorithm computes  $\nabla f(x_k)$  at a series of points  $x_{0:K}$ . We incrementally update our guess of  $H^{-1}$  with an update equation

$$H^{-1} \leftarrow \left( \mathbf{I} - \frac{y\delta^\top}{\delta^\top y} \right)^\top H^{-1} \left( \mathbf{I} - \frac{y\delta^\top}{\delta^\top y} \right) + \frac{\delta\delta^\top}{\delta^\top y}.$$

Note:

- If  $H^{-1}$  is initially zero, this update will assign  $H^{-1} \leftarrow \frac{\delta\delta^\top}{\delta^\top y}$ , which is the minimal rank 1 update we discussed above.
- If  $H^{-1}$  is previously non-zero, the red part “deletes certain dimensions” from  $H^{-1}$ . More precisely, it assigns zero eigenvalues to projections  $\text{TODO}...$ . Discuss the alternative updates. optimality principle behind this:

$$H_{k+1} = \underset{H}{\operatorname{argmin}} \|H - H_k\| \quad \text{s.t.} \quad H = H^\top, \quad \delta = Hy \quad (70)$$

7

The **BFGS** algorithms uses this  $H^{-1}$  instead of a precise  $\nabla^2 f(x)^{-1}$  to compute the steps in a Newton method. All we said about line search and Levenberg-Marquardt damping is unchanged.

In very high-dimensional spaces we do not want to store  $H^{-1}$  densely. Instead we use a compressed storage for low-rank matrices, e.g., storing vectors  $\{v_i\}$  such that  $H^{-1} = \sum_i v_i v_i^\top$ . **Limited memory BFGS (L-BFGS)** makes this more memory efficient: it limits the rank of the  $H^{-1}$  and thereby the used memory. I do not know the details myself, but I assume that with every update it might aim to delete the lowest eigenvalue to keep the rank constant.

#### 4.5.6 Conjugate Gradient\*

The “Conjugate Gradient Method” is a method for solving large linear eqn. systems  $Ax + b = 0$ . We only mention its extension for optimizing nonlinear functions  $f(x)$ .

As above, assume that we evaluated  $\nabla f(x_1)$  and  $\nabla f(x_2)$  at two different points  $x_1, x_2 \in \mathbb{R}^n$ . But now we make one more assumption: The point  $x_2$  is the minimum of a line search from  $x_1$  along the direction  $\delta_1$ . This latter assumption is quite optimistic: it assumes we did perfect line search. But it gives great information: The iso-lines of  $f(x)$  at  $x_2$  are tangential to  $\delta_1$ .

<sup>7</sup>Taken from Peter Blomgren’s lecture slides: [terminus.sdsu.edu/SDSU/Math693a\\_f2013/Lectures/18/lecture.pdf](http://terminus.sdsu.edu/SDSU/Math693a_f2013/Lectures/18/lecture.pdf) This is the original Davidon-Fletcher-Powell (DFP) method suggested by W.C. Davidon in 1959. The original paper describing this revolutionary idea the first quasi-Newton method was not accepted for publication. It later appeared in 1991 in the first issue the the SIAM Journal on Optimization.

In this setting, convince yourself of the following: Ideally each search direction should be orthogonal to the previous one—but not orthogonal in the conventional Euclidean sense, but orthogonal w.r.t. the Hessian  $H$ . Two vectors  $\delta$  and  $\delta'$  are called **conjugate** w.r.t. a metric  $H$  iff  $d'^T H d = 0$ . Therefore, we one subsequent search directions to be conjugate to each other.

Conjugate gradient descent does the following:

---

#### Algorithm 4 Conjugate gradient descent

---

**Input:** initial  $x \in \mathbb{R}^n$ , functions  $f(x)$ ,  $\nabla f(x)$ , tolerance  $\theta$

**Output:**  $x$

```

1: initialize descent direction  $d = g = -\nabla f(x)$ 
2: repeat
3:    $\alpha \leftarrow \operatorname{argmin}_{\alpha} f(x + \alpha d)$  // line search
4:    $x \leftarrow x + \alpha d$ 
5:    $g' \leftarrow g$ ,  $g = -\nabla f(x)$  // store and compute grad
6:    $\beta \leftarrow \max \left\{ \frac{g'^T (g - g')}{g'^T g'}, 0 \right\}$ 
7:    $d \leftarrow g + \beta d$  // conjugate descent direction
8: until  $|\Delta x| < \theta$ 

```

---

- The equation for  $\beta$  is by Polak-Ribière: On a quadratic function  $f(x) = x^T H x$  this leads to **conjugate** search directions,  $d'^T H d = 0$ .
- Intuitively,  $\beta > 0$  implies that the new descent direction always adds a bit of the old direction. This essentially provides 2nd order information.
- For arbitrary quadratic functions CG converges in  $n$  iterations. But this really only works with **perfect line search**.

#### 4.5.7 Rprop\*

Read through Algorithm 5. Notes on this:

- Stepsize adaptation is done in each coordinate *separately!*
- The algorithm not only ignores  $|\nabla f|$  but also its exact direction! Only the gradient signs in each coordinate are relevant. Therefore, the step directions may differ up to  $< 90^\circ$  from  $-\nabla f$ .
- It often works surprisingly efficient and robust.
- If you like, have a look at:

**Algorithm 5** Rprop**Input:** initial  $x \in \mathbb{R}^n$ , function  $f(x)$ ,  $\nabla f(x)$ , initial stepsize  $\alpha$ , tolerance  $\theta$ **Output:**  $x$ 


---

```

1: initialize  $x = x_0$ , all  $\alpha_i = \alpha$ , all  $g_i = 0$ 
2: repeat
3:    $g \leftarrow \nabla f(x)$ 
4:    $x' \leftarrow x$ 
5:   for  $i = 1 : n$  do
6:     if  $g_i g'_i > 0$  then                                     // same direction as last time
7:        $\alpha_i \leftarrow 1.2\alpha_i$ 
8:        $x_i \leftarrow x_i - \alpha_i \text{sign}(g_i)$ 
9:        $g'_i \leftarrow g_i$ 
10:    else if  $g_i g'_i < 0$  then                                 // change of direction
11:       $\alpha_i \leftarrow 0.5\alpha_i$ 
12:       $x_i \leftarrow x_i - \alpha_i \text{sign}(g_i)$ 
13:       $g'_i \leftarrow 0$                                        // force last case next time
14:    else
15:       $x_i \leftarrow x_i - \alpha_i \text{sign}(g_i)$ 
16:       $g'_i \leftarrow g_i$ 
17:    end if
18:    optionally: cap  $\alpha_i \in [\alpha_{\min} x_i, \alpha_{\max} x_i]$ 
19:  end for
20: until  $|x' - x| < \theta$  for 10 iterations in sequence

```

---

Christian Igel, Marc Toussaint, W. Weishui (2005): Rprop using the natural gradient compared to Levenberg-Marquardt optimization. In Trends and Applications in Constructive Approximation. International Series of Numerical Mathematics, volume 151, 259-272.

## 4.6 Convex Problems

We do not put much emphasis on discussing convex problems in this lectures. The algorithms we discussed so far equally apply on general non-linear programs as well as on convex problems—of course, only on convex problems we have convergence guarantees to the optimum, as we can see from the convergence rate analysis of Wolfe steps based on the assumption of positive upper and lower bounds on the hessian's eigenvalues.

Nevertheless, we at least define standard LPs, QPs, etc. Perhaps the most interesting part is the discussion of the Simplex algorithm—not because the algorithms is nice or particularly efficient, but rather because one gains a lot of insights in what actually makes (inequality) constrained problems hard.

### 4.6.1 Convex sets, functions, problems

**Definition 4.4** (Convex set, convex function). A set  $X \subseteq V$  (a subset of some vector space  $V$ ) is **convex** iff

$$\forall x, y \in X, a \in [0, 1] : ax + (1-a)y \in X \quad (71)$$

A function is defined

$$\mathbf{convex} \iff \forall x, y \in \mathbb{R}^n, a \in [0, 1] : f(ax + (1-a)y) \leq a f(x) + (1-a) f(y) \quad (72)$$

$$\mathbf{quasiconvex} \iff \forall x, y \in \mathbb{R}^n, a \in [0, 1] : f(ax + (1-a)y) \leq \max\{f(x), f(y)\} \quad (73)$$

Note: quasiconvex  $\iff$  for any  $\alpha \in \mathbb{R}$  the sublevel set  $\{x | f(x) \leq \alpha\}$  is convex. Further, I call a function **unimodal** if it has only one local minimum, which is the global minimum.

**Definition 4.5** (Convex program).

*Variant 1:* A mathematical program  $\min_x f(x)$  s.t.  $g(x) \leq 0, h(x) = 0$  is convex if  $f$  is convex and the feasible set is convex.

*Variant 2:* A mathematical program  $\min_x f(x)$  s.t.  $g(x) \leq 0, h(x) = 0$  is convex if  $f$  and every  $g_i$  are convex and  $h$  is linear.

Variant 2 is the stronger and usual definition. Concerning variant 1, if the feasible set is convex the zero-level sets of all  $g$ 's need to be convex and the zero-level sets of  $h$ 's needs to be linear. Above these zero levels the  $g$ 's and  $h$ 's could in principle be arbitrarily non-linear, but these non-linearities are irrelevant for the mathematical program itself. We could replace such  $g$ 's and  $h$ 's by convex and linear functions and get the same problem.

### 4.6.2 Linear and quadratic programs

**Definition 4.6** (Linear program (LP), Quadratic program (QP)). Special case mathematical programs are

$$\mathbf{Linear Program (LP):} \quad \min_x c^\top x \quad \text{s.t.} \quad Gx \leq h, Ax = b \quad (74)$$

$$\mathbf{LP in standard form:} \quad \min_x c^\top x \quad \text{s.t.} \quad x \geq 0, Ax = b \quad (75)$$

$$\mathbf{Quadratic Program (QP):} \quad \min_x \frac{1}{2} x^\top Qx + c^\top x \quad \text{s.t.} \quad Gx \leq h, Ax = b, Q \text{ pos-def} \quad (76)$$

Rarely, also a **Quadratically Constrained QP (QCQP)** is considered.

An important example for LP are **relaxations** of integer linear programs,

$$\min_x c^\top x \quad \text{s.t.} \quad Ax = b, x_i \in \{0, 1\},$$

which includes Travelling Salesman, MaxSAT or MAP inference problems. Relaxing such a problem means to instead solve the continuous LP

$$\min_x c^\top x \quad \text{s.t.} \quad Ax = b, x_i \in [0, 1].$$

If one is lucky and the continuous LP problem converges to a fully integer solution, where all  $x_i \in \{0, 1\}$ , this is also the solution to the integer problem. Typically, the solution of the continuous LP will be partially integer (some values converge to the extreme  $x_i \in \{0, 1\}$ , while others are inbetween  $x_i \in (0, 1)$ ). This continuous valued solution gives a lower bound on the integer problem, and provides very efficient heuristics for backtracking or branch-and-bound search for a fully integer solution.

The standard example for a QP are Support Vector Machines. The primal problem is

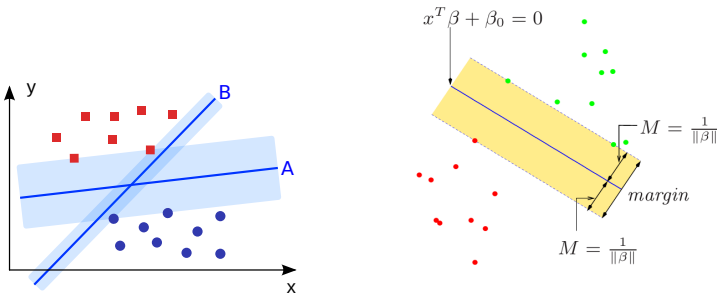
$$\min_{\beta, \xi} \|\beta\|^2 + C \sum_{i=1}^n \xi_i \quad \text{s.t.} \quad y_i(x_i^\top \beta) \geq 1 - \xi_i, \quad \xi_i \geq 0$$

the dual

$$l(\alpha, \mu) = \min_{\beta, \xi} L(\beta, \xi, \alpha, \mu) = -\frac{1}{4} \sum_{i=1}^n \sum_{i'=1}^n \alpha_i \alpha_{i'} y_i y_{i'} \hat{x}_i^\top \hat{x}_{i'} + \sum_{i=1}^n \alpha_i$$

$$\max_{\alpha, \mu} l(\alpha, \mu) \quad \text{s.t.} \quad 0 \leq \alpha_i \leq C$$

(See ML lecture 5:13 for a derivation.)



### 4.6.3 The Simplex Algorithm

Consider an LP. We make the following observations:

- First, in LPs the equality constraints could be resolved simply by introducing new coordinates along the zero-hyperplane of  $h$ . Therefore, for the conceptual discussion we neglect equality constraints.

- The objective constantly pulls in the direction  $-c = -\nabla f(x)$ .
- If the solution is bounded there need to be some inequality constraints that keep the solution from travelling to  $\infty$  in the  $-c$  direction.
- It follows: **The solution will always be located at a vertex**, that is, an intersection point of several zero-hyperplanes of inequality constraints.
- In fact, we should think of the feasible region as a **polytope** that is defined by all the zero-hyperplanes of the inequalities. The inside the polytope is the feasible region. The polytope has edges (intersections of two constraint planes), faces, etc. A solution will always be located at a vertex of the polytope; more precisely, there could be a whole set of optimal points (on a face orthogonal to  $c$ ), but at least one vertex is also optimal.
- An idea for finding the solution is to **walk on the edges of the polytope** until an optimal vertex is found. This is the simplex algorithm of Georg Dantzig, 1947. In practise this procedure is done by “pivoting on the simplex tableaux”—but we fully skip such details here.
- The simplex algorithm is often efficient, but in worst case it is exponential in both,  $n$  and  $m$ ! This is hard to make intuitive, because the effects of high dimensions are not intuitive. But roughly, consider that in high dimensions there is a combinatorial number of ways of how constraints may intersect and form edges and vertices.

Here is a view that much more relates to our discussion of the log barrier method: Sitting on an edge/face/vertex is equivalent to temporarily deciding which constraints are active. If we knew which constraints are eventually active, the problem would be solved: all inequalities become equalities or void. So, walking on the edges of the polytope is equivalent to sequentially making decisions on which constraints might be active. Note though that there are  $2^m$  configurations of active/non-active constraints. The simplex algorithm therefore walks through this combinatorial space.

Interior point methods do exactly the opposite: Recall that the **4th KKT condition** is  $\lambda_i g_i(x) = 0$ . The log barrier method (for instance) instead *relaxes* this hard logic of active/non-active constraints and finds in each iteration a solution to the relaxed 4th KKT condition  $\lambda_i g_i(x) = -\mu$ , which intuitively means that every constraint may be “somewhat active”. In fact, every constraint contributes somewhat to the stationarity condition via the log barrier’s gradients. Thereby interior point methods

- post-pone the hard decisions about active/non-active constraints
- approach the optimal vertex from the inside of the polytope; avoiding the polytope surface (and its hard decisions)



- thereby avoids the need to search through a combinatorial space of constraint activities and instead continuously converges to a decision
- has polynomial worst-case guaranteed

Historically, penalty and barrier methods were standard before the Simplex Algorithm. When SA was discovered in the 50ies, it was quickly considered great. But then, later in the 70-80ies, a lot more theory was developed for interior point methods, which now again have become somewhat more popular than the simplex algorithm.

#### 4.6.4 Sequential Quadratic Programming

Just for reference, SQP is another standard approach to solving non-linear mathematical programs. In each iteration we compute all coefficients of the 2nd order Taylor  $f(x + \delta) \approx f(x) + \nabla f(x)^\top \delta + \frac{1}{2} \delta^\top H \delta$  and 1st-order Taylor  $g(x + \delta) \approx g(x) + \nabla g(x)^\top \delta$  and then solve the QP

$$\min_{\delta} f(x) + \nabla f(x)^\top \delta + \frac{1}{2} \delta^\top H \delta \quad \text{s.t.} \quad g(x) + \nabla g(x)^\top \delta \leq 0$$

The optimal  $\delta^*$  of this problem should be seen analogous to the optimal Newton step: If  $f$  were a 2nd-order polynomial and  $g$  linear, then  $\delta^*$  would jump directly to the optimum. However, as this is generally not the case  $\delta^*$  only gives us a very good direction for line search. In SQP, we need to backtrack until we found a feasible point and  $f$  decreases sufficiently.

### 4.7 Blackbox & Global Optimization: It's all about learning

Even if  $f, g, h$  are smooth, the solver might not have access to analytic equations or efficient numeric methods to evaluate the gradients or Hessians of these. Therefore we distinguish (here neglecting the constraint functions  $g$  and  $h$ ):

#### Definition 4.7.

- *Blackbox optimization*: Only  $f(x)$  can be evaluated.
- *1st-order/gradient optimization*: Only  $f(x)$  and  $\nabla f(x)$  can be evaluated.
- *Quasi-Newton optimization*: Only  $f(x)$  and  $\nabla f(x)$  can be evaluated, but the solver does tricks to estimate  $\nabla^2 f(x)$ . (So this is a special case of 1st-order optimization.)
- *Gauss-Newton type optimization*:  $f$  is of the special form  $f(x) = \phi(x)^\top \phi(x)$  and  $\frac{\partial}{\partial x} \phi(x)$  can be evaluated.

- *2nd order optimization*:  $f(x)$ ,  $\nabla f(x)$  and  $\nabla^2 f(x)$  can be evaluated.

In this lecture I very briefly want to add comments on **global** blackbox optimization. Global means that we now, for the first time, really aim to find the global optimum (within some pre-specified bounded range). In essence, to address such a problem we need to explicitly know what we know about  $f^8$ , and an obvious way to do this is to use Bayesian learning.

#### 4.7.1 A sequential decision problem formulation

From now on, let's neglect constraints and focus on the mathematical program

$$\min_x f(x)$$

for a blackbox function  $f$ . The optimization process can be viewed as a Markov Decision Process that describes the interaction of the solver (agent) with the function (environment):

- At step  $t$ ,  $D_t = \{(x_i, y_i)\}_{i=1}^{t-1}$  is the data that the solver has collected from previous samples. This  $D_t$  is the state of the MDP.
- At step  $t$ , the solver may choose a new decision  $x_t$  about where to sample next.
- Given state  $D_t$  and decision  $x_t$ , the next state is  $D_{t+1} = D \cup \{(x_t, f(x_t))\}$ , which is a deterministic transition *given* the function  $f$ .
- A solver policy is a mapping  $\pi : D_t \mapsto x_t$  that maps any state (of knowledge) to a new decision.
- We may define an **optimal solver policy** as

$$\pi^* = \operatorname{argmin}_{\pi} \langle y_T \rangle = \operatorname{argmin}_{\pi} \int_f P(f) P(D_T | \pi, f) y_T$$

where  $P(D_T | \pi, f)$  is deterministic, and  $P(f)$  is a **prior over functions**.

This objective function cares only about the *last* value  $y_T$  sampled by the solver for a fixed time horizon (budget)  $T$ . Alternatively, we may choose objectives  $\sum_{t=1}^T y_t$  or  $\sum_{t=1}^T \gamma^t y_t$  for some discounting  $\gamma \in [0, 1]$

The above defined what is an optimal solver! Something we haven't touched at all before. The transition dynamics of this MDP is deterministic, given  $f$ . However, from the perspective of the solver, we do not know  $f$  a priori. But we can always

<sup>8</sup>Cf. the KWIK (knows what it knows) framework.

compute a **posterior belief**  $P(f|D_t) = P(D_t|f) P(f)/P(D_t)$ . This posterior belief defines a **belief MDP** with stochastic transitions

$$P(D_{t+1}) = \int_{D_t} \int_f \int_{x_t} [D_{t+1} = D \cup \{(x_t, f(x_t))\}] \pi(x_t|D_t) P(f|D_t) P(D_t) .$$

The belief MDP’s state space is  $P(D_t)$  (or equivalently,  $P(f|D_t)$ , the current belief over  $f$ ). This belief MDP is something that the solver can, in principle, forward simulate—it has all information about it. One can prove that, if the solver could solve its own belief MDP (find an optimal policy for its belief MDP), then this policy is the optimal solver policy for the original problem given a prior distribution  $P(f)$ ! So, in principle we not only defined what is an optimal solver policy, but can also provide an algorithm to compute it (Dynamic programming in the belief MDP)! However, this is so expensive to compute that heuristics need to be used in practise.

One aspect we should learn from this discussion: The solver’s optimal decision is based on its current belief  $P(f|D_t)$  over the function. This belief is the Bayesian representation of everything one could possibly have learned about  $f$  from the data  $D_t$  collected so far. Bayesian Global Optimization methods compute  $P(f|D_t)$  in every step and, based on this, use a heuristic to choose the next decision.

#### 4.7.2 Acquisition Functions for Bayesian Global Optimization\*

In practice one typically uses a Gaussian Process representation of  $P(f|D_t)$ . This means that in every iteration we have an estimate  $\hat{f}(x)$  of the function mean and a variance estimate  $\hat{\sigma}(x)^2$  that describes our uncertainty about the mean estimate. Based on this we may define the following acquisition functions

**Definition 4.8. Probability of Improvement (MPI)**

$$\alpha_t(x) = \int_{-\infty}^{y^*} \mathcal{N}(y|\hat{f}(x), \hat{\sigma}(x))$$

**Expected Improvement (EI)**

$$\alpha_t(x) = \int_{-\infty}^{y^*} \mathcal{N}(y|\hat{f}(x), \hat{\sigma}(x)) (y^* - y)$$

**Upper Confidence Bound (UCB)**

$$\alpha_t(x) = -\hat{f}(x) + \beta_t \hat{\sigma}(x)$$

**Predictive Entropy Search ?**

$$\alpha_t(x) = H[p(x^*|D_t)] - \mathbb{E}\{p(y|D_t; x)\} H[p(x^*|D_t \cup \{(x, y)\})] \tag{77}$$

$$= I(x^*, y|D_t) = H[p(y|D_t, x)] - \mathbb{E}\{p(x^*|D_t)\} H[p(y|D_t, x, x^*)] \tag{78}$$

The last one is special; we’ll discuss it below.

These acquisition functions are heuristics that define how valuable it is to acquire data from the site  $x$ . The solver then makes the decision

$$x_t = \operatorname{argmax}_x \alpha_t(x) .$$

MPI is hardly being used in practise anymore. EI is classical, originating way back in the 50ies or earlier, I think ?. [[TODO]]. UCB received a lot of attention recently due to the underlying bandit theory and bounded regret theorems due to the sub-modularity ?. But I think that in practise EI and UCB perform about equally. As UCB is somewhat easier to implement and intuitive.

In all cases, note that the solver policy  $x_t = \operatorname{argmax}_x \alpha_t(x)$  requires to internally solve another non-linear optimization problem. However,  $\alpha_t$  is an analytic function for which we can compute gradients and Hessians which ensures every efficient *local* convergence. But again,  $x_t = \operatorname{argmax}_x \alpha_t(x)$  needs to be solved *globally*—otherwise the solver will also not solve the original problem properly and globally. As a consequence, the optimization of the acquisition function needs to be restarted from many many potential start points close to potential local minima; typically from grid(!) points over the full domain range. The number of grid points is exponential in the problem dimension  $n$ . Therefore, this inner loop can be very expensive.

And a subjective note: This all sounds great, but be aware that Gaussian Processes with standard squared-exponential kernels do not generalize much in high dimensions: one roughly needs exponentially many data points to fully cover the domain and reduce belief uncertainty globally, almost as if we were sampling from a grid with grid size equal to the kernel width. So, the whole approach is not magic. It just does what is possible given a belief  $P(f)$ . It would be interesting to have much more structured (and heteroscedastic) beliefs specifically for optimization.

The last acquisition function is called **Predictive Entropy Search** . This formulation is beautiful: We sample at places  $x$  where the (expected) observed value  $y$  informs us as much as possible about the optimum  $x^*$  of the function! Formally, this means to maximize the mutual information between  $y$  and  $x^*$ , in expectation over  $y|x$ .

### 4.7.3 Classical model-based blackbox optimization (non-global)\*

A last method very worth mentioning: Classical model-based blackbox optimization simply fits a local polynomial model to the recent data and takes this a basis for search. This is similar to BFGS, but now for the blackbox case where we not even observe gradients. See Algorithm 6.

The local fitting of a polynomial model is again a Machine Learning method. Whether this gives a function approximation for optimization depends on the quality of the data  $D_t$  used for this approximation. Classical model-based optimization has interesting heuristics to evaluate the data quality as well as sample new points to improve the data quality. Here is a rough algorithm (following Nodet et al.'s section

on “Derivative-free optimization”):

---

**Algorithm 6** Classical model-based optimization
 

---

```

1: Initialize  $D$  with at least  $\frac{1}{2}(n+1)(n+2)$  data points
2: repeat
3:   Compute a regression  $\hat{f}(x) = \phi_2(x)^\top \beta$  on  $D$ 
4:   Compute  $\delta = \operatorname{argmin}_\delta \hat{f}(\hat{x} + \delta)$  s.t.  $|\delta| < \alpha$ 
5:   if  $f(\hat{x} + \delta) < f(\hat{x}) - \varrho_{\text{ls}}[f(\hat{x}) - \hat{f}(\hat{x} + \delta)]$  then           // test sufficient decrease
6:     Increase the stepsize  $\alpha$ 
7:     Accept  $\hat{x} \leftarrow \hat{x} + \delta$ 
8:     Add to data,  $D \leftarrow D \cup \{(\hat{x}, f(\hat{x}))\}$ 
9:   else                                                                 // no sufficient decrease
10:    if  $\det(D)$  is too small then                                     // blame the data quality
11:      Compute  $x^+ = \operatorname{argmax}_{x'} \det(D \cup \{x'\})$  s.t.  $|x - x'| < \alpha$ 
12:      Add to data,  $D \leftarrow D \cup \{(x^+, f(x^+))\}$ 
13:    else                                                                 // blame the stepsize
14:      Decrease the stepsize  $\alpha$ 
15:    end if
16:  end if
17:  Perhaps prune the data, e.g., remove  $\operatorname{argmax}_{x \in \Delta} \det(D \setminus \{x\})$ 
18: until  $x$  converges

```

---

Some notes:

- Line 4 implement an explicit trust region approach, which hard bound  $\alpha$  on the step size.
- Line 5 is like the Wolfe condition. But here, the expected decrease is  $[f(\hat{x}) - \hat{f}(\hat{x} + \delta)]$  instead of  $-\alpha \delta \nabla f(x)$ .
- If there is no sufficient decrease we may blame it on two reasons: bad data or a too large stepsize.
- Line 10 uses the *data determinant* as a measure of quality! This is meant in the sense of linear regression on polynomial features. Note that, with data matrix  $X \in \mathbb{R}^{n \times \dim(\beta)}$ ,  $\hat{\beta}^{\text{ls}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top y$  is the optimal regression. The determinant  $\det(\mathbf{X}^\top \mathbf{X})$  or  $\det(\mathbf{X}) = \det(D)$  is a measure for well the data supports the regression. If the determinant is zero, the regression problem is ill-defined. The larger the determinant, the lower the variance of the regression estimator.
- Line 11 is an explicit exploration approach: We add a data point solely for the purpose of increasing the data determinant (increasing the data spread). Interesting. Nocedal describes in more detail a geometry-improving procedure to update  $D$ .

### 4.7.4 Evolutionary Algorithms\*

There are interesting and theoretically well-grounded evolutionary algorithms for optimization, such as Estimation-of-Distribution Algorithms (EDAs). But generally, don't use them as first choice.

## 4.8 Examples and Exercises

### 4.8.1 Optimize a constrained problem

Consider the following constrained problem

$$\min_x \sum_{i=1}^n x_i \quad \text{s.t.} \quad g(x) \leq 0 \quad (79)$$

$$g(x) = \begin{pmatrix} x^\top x - 1 \\ -x_1 \end{pmatrix} \quad (80)$$

a) First, assume  $x \in \mathbb{R}^2$  is 2-dimensional, and draw on paper what the problem looks like and where you expect the optimum.

b) Find the optimum analytically using the Lagrangian. Here, assume that you know apriori that all constraints are active! What are the dual parameters  $\lambda = (\lambda_1, \lambda_2)$ ?

Note: Assuming that you know a priori which constraints are active is a huge assumption! In real problems, this is the actual hard (and combinatorial) problem. More on this later in the lecture.

c) Implement a simple the Log Barrier Method. Tips:

- Initialize  $x = (\frac{1}{2}, \frac{1}{2})$  and  $\mu = 1$
- First code an inner loop:
  - In each iteration, first compute the gradient of the log-barrier function. Recall that

$$F(x; \mu) = f(x) - \mu \sum_i \log(-g_i(x)) \quad (81)$$

$$\nabla F(x; \mu) = \nabla f - \mu \sum_i (1/g_i(x)) \nabla g_i(x) \quad (82)$$

- Then make a small step along the negative gradient  $x \leftarrow x - \alpha \nabla F(x, \mu)$ , for  $\alpha = 1/100$ .
- Iterate until convergence; let's call the result  $x^*(\mu)$ . Further, compute  $\lambda^*(m) = -(\mu/g_1(x), \mu/g_2(x))$  at convergence.
- Decrease  $\mu \leftarrow \mu/2$ , recompute  $x^*(\mu)$  (with the previous  $x^*$  as initialization) and iterate this.

Does  $x^*$  and  $\lambda^*$  converge to the expected solution?

Note: The path  $x^*(\mu) = \operatorname{argmin}_x F(x; \mu)$  (the optimum in dependence of  $\mu$ ) is called *central path*.

### 4.8.2 Lagrangian and dual function

(Taken roughly from ‘Boyd et al: Convex Optimization’, Ex. 5.1)

A simple example. Consider the optimization problem

$$\min x^2 + 1 \quad \text{s.t.} \quad (x - 2)(x - 4) \leq 0$$

with variable  $x \in \mathbb{R}$ .

- Derive the optimal solution  $x^*$  and the optimal value  $p^* = f(x^*)$  by hand.
- Write down the Lagrangian  $L(x, \lambda)$ . Plot (using gnuplot or so)  $L(x, \lambda)$  over  $x$  for various values of  $\lambda \geq 0$ . Verify the lower bound property  $\min_x L(x, \lambda) \leq p^*$ .
- Derive the dual function  $l(\lambda) = \min_x L(x, \lambda)$  and plot it (for  $\lambda \geq 0$ ). Derive the dual optimal solution  $\lambda^* = \operatorname{argmax}_\lambda l(\lambda)$ . Is  $\max_\lambda l(\lambda) = p^*$  (strong duality)?

### 4.8.3 Convergence proof

a) Given a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  with  $f_{\min} = \min_x f(x)$ . Assume that its Hessian—that is, the eigenvalues of  $\nabla^2 f$ —are lower bounded by  $m > 0$  and upper bounded by  $M > m$ . Proof that for any  $x \in \mathbb{R}^n$  it holds

$$f(x) - \frac{1}{2m} |\nabla f(x)|^2 \leq f_{\min} \leq f(x) - \frac{1}{2M} |\nabla f(x)|^2.$$

Tip: Start with bounding the 2nd-order Taylor expansion. Then consider the minima of these bounds. Note, it also follows:

$$|\nabla f(x)|^2 \geq 2m(f(x) - f_{\min}).$$

b) Consider backtracking line search with Wolfe parameter  $\varrho_{\text{ls}} \leq \frac{1}{2}$ , and step decrease factor  $\varrho_{\alpha}^-$ . First prove that line search terminates the latest when  $\frac{\varrho_{\alpha}^-}{M} \leq \alpha \leq \frac{1}{M}$ , and then it found a new point  $y$  for which

$$f(y) \leq f(x) - \frac{\varrho_{\text{ls}} \varrho_{\alpha}^-}{M} |\nabla f(x)|^2.$$

From this, using the result from a), prove the convergence equation

$$f(y) - f_{\min} \leq \left[ 1 - \frac{2m \varrho_{\text{ls}} \varrho_{\alpha}^-}{M} \right] (f(x) - f_{\min}).$$

#### 4.8.4 Robust unconstrained optimization

A ‘flattened’ variant of the Rosenbrock function is defined as

$$f(x) = \log[1 + (x_2 - x_1^2)^2 + \frac{1}{100}(1 - x_2)^2]$$

and has the minimum at  $x^* = (1, 1)$ . For reference, the gradient and hessian are

$$g(x) := 1 + (x_2 - x_1^2)^2 + \frac{1}{100}(1 - x_2)^2 \quad (83)$$

$$\partial_{x_1} f(x) = \frac{1}{g(x)} \left[ -4(x_2 - x_1^2)x_1 \right] \quad (84)$$

$$\partial_{x_2} f(x) = \frac{1}{g(x)} \left[ 2(x_2 - x_1^2) - \frac{2}{100}(1 - x_2) \right] \quad (85)$$

$$\partial_{x_1}^2 f(x) = -\left[ \partial_{x_1} f(x) \right]^2 + \frac{1}{g(x)} \left[ 8x_1^2 - 4(x_2 - x_1^2) \right] \quad (86)$$

$$\partial_{x_2}^2 f(x) = -\left[ \partial_{x_2} f(x) \right]^2 + \frac{1}{g(x)} \left[ 2 + \frac{2}{100} \right] \quad (87)$$

$$\partial_{x_1} \partial_{x_2} f(x) = -\left[ \partial_{x_1} f(x) \right] \left[ \partial_{x_2} f(x) \right] + \frac{1}{g(x)} \left[ -4x_1 \right] \quad (88)$$

a) Use gnuplot to display the function copy-and-pasting the following lines:

```
set isosamples 50,50
set contour
f(x,y) = log(1+(y-(x**2))**2 + .01*(1-x)**2) - 0.01
splot [-3:3][-3:4] f(x,y)
```

(The ‘-0.01’ ensures that you can see the contour at the optimum.) List and discuss at least three properties of the function (at different locations) that may raise problems to naive optimizers.

b) Use  $x = (-3, 3)$  as starting point for an optimization algorithm. Try to code an optimization method that uses all ideas mentioned in the lecture. Try to tune it to be efficient on this problem (without cheating, e.g. by choosing a perfect initial stepsize.)

Solving real-world problems involves 2 subproblems:

- 1) formulating the problem as an optimization problem (conform to a standard optimization problem category) ( $\rightarrow$  human)
- 2) the actual optimization problem ( $\rightarrow$  algorithm)



These exercises focus on the first type, which is just as important as the second, as it enables the use of a wider range of solvers. Exercises from Boyd et al [http://www.stanford.edu/~boyd/cvxbook/bv\\_cvxbook.pdf](http://www.stanford.edu/~boyd/cvxbook/bv_cvxbook.pdf):

#### 4.8.5 Network flow problem

Solve Exercise 4.12 (pdf page 207) from Boyd & Vandenberghe, *Convex Optimization*.

#### 4.8.6 Minimum fuel optimal control

Solve Exercise 4.16 (pdf page 208) from Boyd & Vandenberghe, *Convex Optimization*.

#### 4.8.7 Reformulating Norms

Solve Exercise 4.11 (pdf page 207) from Boyd & Vandenberghe, *Convex Optimization*.

#### 4.8.8 Solve one of these only, the others are optional

a) Grocery Shopping: You're at the market and you find  $n$  offers, each represented by a set of items  $A_i$  and the respective price  $c_i$ . Your goal is to buy at least one of each item for as little as possible. Formulate as an integer LP.

b) Facility Location: There are  $n$  facilities with which to satisfy the needs of  $m$  clients. The cost for opening facility  $j$  is  $f_j$ , and the cost for servicing client  $i$  through facility  $j$  is  $c_{ij}$ . You have to find an optimal way to open facilities and to associate clients to facilities. Formulate as an ILP.

c) Taxicab Driver: You're a taxicab driver in hyper-space ( $\mathbb{R}^d$ ) and have to service  $n$  clients. Each client  $i$  has a known initial position  $c_i \in \mathbb{R}^d$  and a destination  $d_i \in \mathbb{R}^d$ . You start out at position  $p_0 \in \mathbb{R}^d$  and have to service all the clients while minimizing fuel use, which is proportional to covered distance. Hyper-space is funny, so the geometry is not Euclidean and distances are Manhattan distances. Formulate as an LP.

## 5 Probabilities & Information

It is beyond the scope of these notes to give a detailed introduction to probability theory. There are excellent books:

- Thomas & Cover

- Bishop
- MacKay

Instead, we first recap very basics of probability theory, that I assume the reader has already seen before. The next section will cover this. Then we focus on specific topics that, in my opinion, deepen the understanding of the basics, such as the relation between optimization and probabilities, log-probabilities & energies, maxEntropy and maxLikelihood, minimal description length and learning.

## 5.1 Basics

First, in case you wonder about justifications of the use of (Bayesian) probabilities versus fuzzy sets or alike, here some pointers to look up: 1) Cox's theorem, which derives from basic assumptiond about "rationality and consistency" the standard probability axioms; 2) t-norms, which generalize probability and fuzzy calculus; and 3) read about objective vs. subjecte Bayesian probability.

### 5.1.1 Axioms, definitions, Bayes rule

**Definition 5.1** (set-theoretic axioms of probabilities).

- An experiment can have multiple outcomes; we call the set of possible outcomes **sample space** or **domain**  $S$
- A mapping  $P : A \subseteq S \mapsto [0, 1]$ , that maps any subset  $A \subseteq S$  to a real number is called **probability measure** on  $S$  iff
  - $P(A) \geq 0$  for any  $A \subseteq S$  (non-negativity)
  - $P(\bigcup_i A_i) = \sum_i P(A_i)$  if  $A_i \cap A_j = \emptyset$  (additivity)
  - $P(S) = 1$  (normalization)
- Implications are:
  - $0 \leq P(A) \leq 1$
  - $P(\emptyset) = 0$
  - $A \subseteq B \Rightarrow P(A) \leq P(B)$
  - $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
  - $P(S \setminus A) = 1 - P(A)$

Formally, a **random variable**  $X$  is mapping  $X : S \rightarrow \Omega$  from a measureable space  $S$  (that is, a sample space  $S$  that has a probability measure  $P$ ) to another sample space  $\Omega$ , which I typically call the **domain**  $\text{dom}(X)$  of the random variable. Thereby, the mapping  $X : S \rightarrow \Omega$  now also defines a probability measure over the domain  $\Omega$ :

$$P(B \subseteq \Omega) = P(\{s : X(s) \in B\})$$

In practise we just use the following notations:

**Definition 5.2** (Random Variable).

- Let  $X$  be a random variable with discrete domain  $\text{dom}(X) = \Omega$
- $P(X=x) \in \mathbb{R}$  denotes the specific probability that  $X = x$  for some  $x \in \Omega$
- $P(X)$  denotes the **probability distribution** (function over  $\Omega$ )
- $\forall_{x \in \Omega} : 0 \leq P(X=x) \leq 1$
- $\sum_{x \in \Omega} P(X=x) = 1$
- We often use the short hand  $\sum_X P(X) \dots = \sum_{x \in \text{dom}(X)} P(X=x) \dots$  when summing over possible values of a RV

If we have two or more random variables, we have

**Definition 5.3** (Joint, marginal, conditional, independence, Bayes' Theorem).

- We denote the **joint** distribution of two RVs as  $P(X, Y)$
- The **marginal** is defined as  $P(X) = \sum_Y P(X, Y)$
- The **conditional** is defined as  $P(X|Y) = \frac{P(X, Y)}{P(Y)}$ , which fulfils  $\forall_Y : \sum_X P(X|Y) = 1$ .
- $X$  is **independent** of  $Y$  iff  $P(X, Y) = P(X) P(Y)$ , or equivalently,  $P(X|Y) = P(X)$ .
- The definition of a conditional implies the **product rule**

$$P(X, Y) = P(X|Y) P(Y) = P(Y|X) P(X)$$

and **Bayes' Theorem**

$$P(X|Y) = \frac{P(Y|X) P(X)}{P(Y)}$$

The individual terms in Bayes' Theorem are typically given names:

$$\text{posterior} = \frac{\text{likelihood} \cdot \text{prior}}{\text{normalization}}$$

(Sometimes, the normalization is also called *evidence*.)

- $X$  is *conditionally independent* of  $Y$  given  $Z$  iff  $P(X|Y, Z) = P(X|Z)$  or  $P(X, Y|Z) = P(X|Z) P(Y|Z)$

### 5.1.2 Standard discrete distributions

	RV	parameter	distribution
Bernoulli	$x \in \{0, 1\}$	$\mu \in [0, 1]$	$\text{Bern}(x   \mu) = \mu^x (1 - \mu)^{1-x}$
Beta	$\mu \in [0, 1]$	$\alpha, \beta \in \mathbb{R}^+$	$\text{Beta}(\mu   a, b) = \frac{1}{B(a, b)} \mu^{a-1} (1 - \mu)^{b-1}$
Multinomial	$x \in \{1, \dots, K\}$	$\mu \in [0, 1]^K, \ \mu\ _1 = 1$	$P(x = k   \mu) = \mu_k$
Dirichlet	$\mu \in [0, 1]^K, \ \mu\ _1 = 1$	$\alpha_1, \dots, \alpha_K \in \mathbb{R}^+$	$\text{Dir}(\mu   \alpha) \propto \prod_{k=1}^K \mu_k^{\alpha_k - 1}$

Clearly, the Multinomial is a generalization of the Bernoulli, as the Dirichlet is of the Beta. The mean of the Dirichlet is  $\langle \mu_i \rangle = \frac{\alpha_i}{\sum_j \alpha_j}$ , its mode is  $\mu_i^* = \frac{\alpha_i - 1}{\sum_j \alpha_j - K}$ . The **mode** of a distribution  $p(x)$  is defined as  $\text{argmax}_x p(x)$ .

### 5.1.3 Conjugate distributions

**Definition 5.4** (Conjugacy). Let  $p(D|x)$  be a likelihood conditional on a RV  $x$ . A family  $\mathcal{C}$  of distributions (i.e.,  $\mathcal{C}$  is a space of distributions, like the space of all Beta distributions) is called **conjugate** to the likelihood function  $p(D|x)$  iff

$$p(x) \in \mathcal{C} \quad \Rightarrow \quad p(x|D) = \frac{p(D|x) p(x)}{p(D)} \in \mathcal{C}.$$

The standard conjugates you should know:

RV	likelihood	conjugate
$\mu$	Binomial $\text{Bin}(D   \mu)$	Beta $\text{Beta}(\mu   a, b)$
$\mu$	Multinomial $\text{Mult}(D   \mu)$	Dirichlet $\text{Dir}(\mu   \alpha)$
$\mu$	Gauss $\mathcal{N}(x   \mu, \Sigma)$	Gauss $\mathcal{N}(\mu   \mu_0, A)$
$\lambda$	1D Gauss $\mathcal{N}(x   \mu, \lambda^{-1})$	Gamma $\text{Gam}(\lambda   a, b)$
$\Lambda$	$n$ D Gauss $\mathcal{N}(x   \mu, \Lambda^{-1})$	Wishart $\text{Wish}(\Lambda   W, \nu)$
$(\mu, \Lambda)$	$n$ D Gauss $\mathcal{N}(x   \mu, \Lambda^{-1})$	Gauss-Wishart $\mathcal{N}(\mu   \mu_0, (\beta\Lambda)^{-1}) \text{Wish}(\Lambda   W, \nu)$

### 5.1.4 Distributions over continuous domain

**Definition 5.5.** Let  $x$  be a continuous RV. The **probability density function**

**(pdf)**  $p(x) \in [0, \infty)$  defines the probability

$$P(a \leq x \leq b) = \int_a^b p(x) dx \in [0, 1]$$

The **cumulative probability distribution**  $F(y) = P(x \leq y) = \int_{-\infty}^y dx p(x) \in [0, 1]$  is the cumulative integral with  $\lim_{y \rightarrow \infty} F(y) = 1$

However, I and most others say **probability distribution** to refer to probability density function.

One comment about integrals. If  $p(x)$  is a probability density function and  $f(x)$  some arbitrary function, typically one writes

$$\int_x f(x) p(x) dx ,$$

where  $dx$  denotes the (Borel) measure we integrate over. However, some authors (correctly) think of a distribution  $p(x)$  as being a measure over the space  $\text{dom}(x)$  (instead of just a function). So the above notation is actually “double” w.r.t. the measures. So they might (also correctly) write

$$\int_x p(x) f(x) ,$$

and take care that there is exactly one measure to the right of the integral.

### 5.1.5 Gaussian

**Definition 5.6.** We define an  $n$ -dim Gaussian in *normal form* as

$$\mathcal{N}(x | \mu, \Sigma) = \frac{1}{|2\pi\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu)^\top \Sigma^{-1} (x - \mu)\right\}$$

with **mean**  $\mu$  and **covariance** matrix  $\Sigma$ . In *canonical form* we define

$$\mathcal{N}[x | a, A] = \frac{\exp\left\{-\frac{1}{2}a^\top A^{-1}a\right\}}{|2\pi A^{-1}|^{1/2}} \exp\left\{-\frac{1}{2}x^\top A x + x^\top a\right\} \quad (89)$$

with **precision** matrix  $A = \Sigma^{-1}$  and coefficient  $a = \Sigma^{-1}\mu$  (and mean  $\mu = A^{-1}a$ ).

Gaussians are used all over—below we explain in what sense they are the probabilistic analogue to a parabola (or a 2nd-order Taylor expansions). The most important properties are:

- **Symmetry:**  $\mathcal{N}(x | a, A) = \mathcal{N}(a | x, A) = \mathcal{N}(x - a | 0, A)$

- **Product:**

$$\mathcal{N}(x | a, A) \mathcal{N}(x | b, B) = \mathcal{N}[x | A^{-1}a + B^{-1}b, A^{-1} + B^{-1}] \mathcal{N}(a | b, A + B)$$

$$\mathcal{N}[x | a, A] \mathcal{N}[x | b, B] = \mathcal{N}[x | a + b, A + B] \mathcal{N}(A^{-1}a | B^{-1}b, A^{-1} + B^{-1})$$

- **“Propagation”:**

$$\int_y \mathcal{N}(x | a + Fy, A) \mathcal{N}(y | b, B) dy = \mathcal{N}(x | a + Fb, A + FBF^T)$$

- **Transformation:**

$$\mathcal{N}(Fx + f | a, A) = \frac{1}{|F|} \mathcal{N}(x | F^{-1}(a - f), F^{-1}AF^{-T})$$

- **Marginal & conditional:**

$$\mathcal{N}\left(\begin{matrix} x \\ y \end{matrix} \middle| \begin{matrix} a \\ b \end{matrix}, \begin{matrix} A & C \\ C^T & B \end{matrix}\right) = \mathcal{N}(x | a, A) \cdot \mathcal{N}(y | b + C^T A^{-1}(x - a), B - C^T A^{-1}C)$$

More Gaussian identities are found at <http://ipvs.informatik.uni-stuttgart.de/mlr/marc/notes/gaussians.pdf>

pdf

**Example 5.1** (ML estimator of the mean of a Gaussian). Assume we have data  $D = \{x_1, \dots, x_n\}$ , each  $x_i \in \mathbb{R}^n$ , with likelihood

$$P(D | \mu, \Sigma) = \prod_i \mathcal{N}(x_i | \mu, \Sigma)$$

$$\operatorname{argmax}_{\mu} P(D | \mu, \Sigma) = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\operatorname{argmax}_{\Sigma} P(D | \mu, \Sigma) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$$

Assume we are initially uncertain about  $\mu$  (but know  $\Sigma$ ). We can express this uncertainty using again a Gaussian  $\mathcal{N}[\mu | a, A]$ . Given data we have

$$P(\mu | D) \propto P(D | \mu, \Sigma) P(\mu) = \prod_i \mathcal{N}(x_i | \mu, \Sigma) \mathcal{N}[\mu | a, A]$$

$$= \prod_i \mathcal{N}[\mu | \Sigma^{-1}x_i, \Sigma^{-1}] \mathcal{N}[\mu | a, A] \propto \mathcal{N}[\mu | \Sigma^{-1} \sum_i x_i, n\Sigma^{-1} + A]$$

Note: in the limit  $A \rightarrow 0$  (uninformative prior) this becomes

$$P(\mu | D) = \mathcal{N}\left(\mu \middle| \frac{1}{n} \sum_i x_i, \frac{1}{n} \Sigma\right)$$

which is consistent with the Maximum Likelihood estimator

### 5.1.6 “Particle distribution”

Usually, “particles” are not listed as standard continuous distribution. However I think they should be. They’re heavily used in several contexts, especially as approximating other distributions in Monte Carlo methods and particle filters.

**Definition 5.7** (Dirac or  $\delta$ -distribution). In *distribution theory* it is proper to define a distribution  $\delta(x)$  that is the derivative of the Heavyside step function  $H(x)$ ,

$$\delta(x) = \frac{\partial}{\partial x} H(x), \quad H(x) = [x \geq 0].$$

It is akward to think of  $\delta(x)$  as a normal function, as it'd be “infinite” at zero. But at least we understand that is has the properties

$$\delta(x) = 0 \text{ everywhere except at } x = 0, \quad \int \delta(x) dx = 1.$$

I sometimes call the Dirac distribution also a **point particle** : it has all its unit “mass” concentrated at zero.

**Definition 5.8** (Particle Distribution). We define a **particle distribution**  $q(x)$  as a **mixture** of Diracs,

$$q(x) := \sum_{i=1}^N w_i \delta(x - x_i),$$

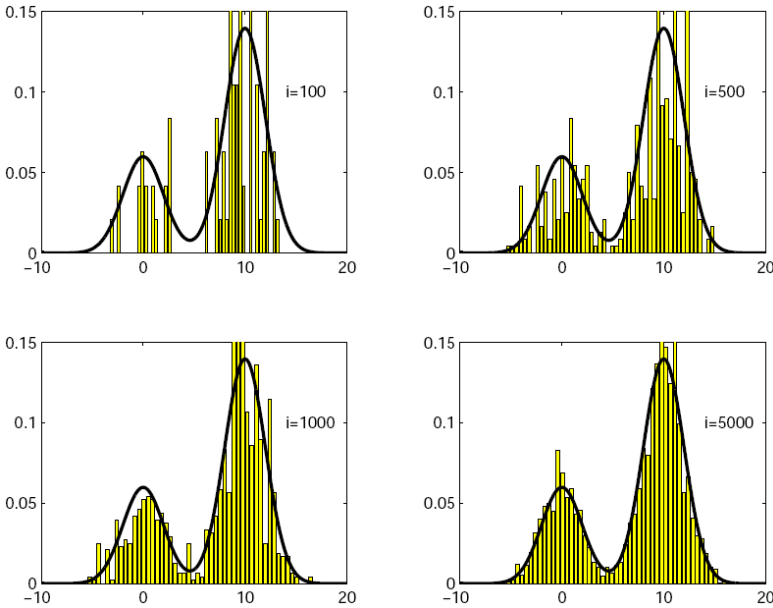
which is parameterized by the number  $N$ , the locations  $\{x_i\}_{i=1}^N$ ,  $x_i \in \mathbb{R}^n$ , and the normalized weights  $\{w_i\}_{i=1}^N$ ,  $w_i \in \mathbb{R}$ ,  $\|w\|_1 = 1$  of the  $N$  particles.

We say that a particle distribution  $q(x)$  approximates another distribution  $p(x)$  iff for any (smooth)  $f$

$$\langle f(x) \rangle_p = \int_x f(x)p(x)dx \approx \sum_{i=1}^N w_i f(x_i)$$

Note the generality of this statement!  $f$  could be anything, it could be any features of the variable  $x$ , like coordinates of  $x$ , or squares, or anything. So basically this statement says, whatever you might like to estimate about  $p$ , you can approximate it based on the particles  $q$ .

Computing particle approximations of complex (non-analytical, non-tracktable) distributions  $p$  is a core challenge in many fields. The true  $p$  could for instance be a distribution over games (action sequences). The approximation  $q$  could for instance be samples generated with Monte Carlo Tree Search (MCTS). The tutorial *An Introduction to MCMC for Machine Learning* [www.cs.ubc.ca/~nando/papers/mlintro.pdf](http://www.cs.ubc.ca/~nando/papers/mlintro.pdf) gives an excellent introduction. Here are some illustrations of what it means to approximate some  $p$  by particles  $q$ , taken from this tutorial. The black line is  $p$ , histograms illustrate the particles  $q$  by showing how many of (uniformly weighted) particles fall into a bin:



(from de Freitas et al.)

## 5.2 Between probabilities and optimization: neg-log-probabilities, exp-neg-energies, exponential family, Gibbs and Boltzmann

There is a natural relation between probabilities and “energy” (or “error”). Namely, if  $p(x)$  denotes a probability for every possible value of  $x$ , and  $E(x)$  denotes an energy for state  $x$ —or an error one assigns to choosing  $x$ —then a natural relation is

$$p(x) = e^{-E(x)}, \quad E(x) = -\log p(x).$$

Why is that? First, outside the context of physics it is perfectly fair to just define axiomatically an energy  $E(x)$  as neg-log-probability. But let me try to give some more arguments for why this is a useful definition.

Let assume we have  $p(x)$ . We want to find a quantity, let’s call it *error*  $E(x)$ , which is a function of  $p(x)$ . Intuitively, if a certain value  $x_1$  is more likely than another,  $p(x_1) > p(x_2)$ , then picking  $x_1$  should imply less error,  $E(x_1) < E(x_2)$  (Axiom 1). Further, when we have two independent random variables  $x$  and  $y$ , **probabilities are multiplicative**,  $p(x, y) = p(x)p(y)$ . We require axiomatically that **error is additive**,  $E(x, y) = E(x) + E(y)$ . From both follows that  $E$  needs to be some logarithm of  $p$ !

The same argument, now more talking about *energy*: Assume we have two independent (physical) systems  $x$  and  $y$ .  $p(x, y) = p(x)p(y)$  is the probability to find



them in certain states. We axiomatically require that **energy is additive**,  $E(x, y) = E(x) + E(y)$ . Again,  $E$  needs to be some logarithm of  $p$ . In the context of physics, what could be questioned is “why is  $p(x)$  a function of  $E(x)$  in the first place?”. Well, that is much harder to explain and really is a question about statistical physics. Wikipedia under keywords “Maxwell-Boltzmann statistics” and “Derivation from microcanonical ensemble” gives an answer. Essentially the argument is as follows: Given many many molecules in a gas, each of which can have a different energy  $e_i$ . The total energy  $E = \sum_{i=1}^n e_i$  must be conserved. What is the distribution over energy levels that has the most microstates? The answer is the Boltzmann distribution. (And why do we, in nature, find energy distributions that have the most microstates? Because these are most likely.)

Bottom line is:  $p(x) = e^{-E(x)}$ , probabilities are multiplicative, energies or errors additive.

Let me state some fact just to underline how useful this way of thinking is:

- Given an energy function  $E(x)$ , its **Boltzmann distribution** is defined as

$$p(x) = e^{-E(x)}.$$

This is sometimes also called Gibbs distribution.

- In machine learning, when data  $D$  is given and we have some model  $\beta$ , we typically try to maximize the likelihood  $p(D|\beta)$ . This is equivalent to minimizing the neg-log-likelihood

$$L(\beta) = -\log p(D|\beta).$$

This neg-log-likelihood is a typical measure for *error* of the model. And this error is additive w.r.t. the data, whereas the likelihood is multiplicative, fitting perfectly to the above discussion.

- The Gaussian distribution  $p(x) \propto \exp\{-\frac{1}{2}\|x - \mu\|^2/\sigma^2\}$  is related to the error  $E(x) = \frac{1}{2}\|x - \mu\|^2/\sigma^2$ , which is nothing but the squared error with the precision matrix as metric. That’s why squared error measures (classical regression) and Gaussian distributions (e.g., Bayesian Ridge regression) are directly related.

A Gaussian is the probabilistic analogue to a parabola.

- The exponential family is defined as

$$p(x|\beta) = h(x)g(\beta) \exp\{\beta^\top \phi(x)\}$$

Often  $h(x) = 1$ , so let’s neglect this for now. The key point is that the energy is linear in the features  $\phi(x)$ . This is exactly how discriminative functions (for classification in Machine learning) are typically formulated.

In the continuous case, the features  $\phi(x)$  are often chosen as basis polynomials—just as in polynomial regression. Then,  $\beta$  are the coefficients of the energy polynomial and the exponential family is just the probabilistic analogue to the space of polynomials.

- When we have many variables  $x_1, \dots, x_n$ , the structure of a cost function over these variables can often be expressed as being additive in terms:  $f(x_1, \dots, x_n) = \sum_i \phi_i(x_{\partial i})$  where  $\partial i$  denotes the  $i$ th group of variables. The respective Boltzmann distribution is a **factor graph**  $p(x_1, \dots, x_n) \propto \prod_i f_i(x_{\partial i}) = \exp\{\sum_i \beta_i \phi_i(x_{\partial i})\}$  where  $\partial i$  denotes the

So, factor graphs are the probabilistic analogue to additive functions.

- $-\log p(x)$  is also the “optimal” coding length you should assign to a symbol  $x$ .

**Entropy is expected error**:  $H[p] = \sum_x -p(x) \log p(x) = \langle -\log p(x) \rangle_{p(x)}$ , where  $p$  itself it used to take the expectation.

Assume you use a “wrong” distribution  $q(x)$  to decide on the coding length of symbols drawn from  $p(x)$ . The expected length of an encoding is  $\int_x p(x) [-\log q(x)] \geq H(p)$ .

The **Kullback-Leibler divergence** is the difference:

$$D(p \| q) = \int_x p(x) \log \frac{p(x)}{q(x)} \geq 0$$

Proof of inequality, using the Jensen inequality:

$$-\int_x p(x) \log \frac{q(x)}{p(x)} \geq -\log \int_x p(x) \frac{q(x)}{p(x)} = 0$$

So, my message is that probabilities and error measures are naturally related. However, in the first case we typically do inference, in the second we optimize. Let’s discuss the relation between inference and optimization a bit more. For instance, given data  $D$  and parameters  $\beta$ , we may define

**Definition 5.9** (ML, MAP, and Bayes estimate). Given data  $D$  and a parametric model  $p(D|\beta)$ , we define

- **Maximum likelihood (ML) parameter estimate:**

$$\beta^{\text{ML}} := \operatorname{argmax}_{\beta} P(D|\beta)$$

- **Maximum a posteriori (MAP) parameter estimate:**

$$\beta^{\text{MAP}} = \operatorname{argmax}_{\beta} P(\beta|D)$$

- **Bayesian parameter estimate:**

$$P(\beta|D) \propto P(D|\beta) P(\beta)$$

used for **Bayesian prediction**:  $P(\text{prediction}|D) = \int_{\beta} P(\text{prediction}|\beta) P(\beta|D)$

Both, the MAP and the ML estimates are really just optimization problems.

The Bayesian parameter estimate  $P(\beta|D)$ , which can then be used to do fully Bayesian prediction, is in principle different. However, in practise also here optimization is a core tool for estimating such distributions if they cannot be given analytically. This is described next.

### 5.3 Information, Entropie & Kullback-Leibler

Consider the following problem. We have data drawn i.i.d. from  $p(x)$  where  $x \in X$  in some discrete space  $X$ . Let's call every  $x$  a *word*. The problem is to find a mapping from words to *codes*, e.g. binary codes  $c : X \rightarrow \{0, 1\}^*$ . The optimal solution is in principle simple: Sort all possible words in a list, ordered by  $p(x)$  with more likely words going first; write all possible binary codes in another list, with increasing code lengths. Match the two lists, and this is the optimal encoding.

Let's try to get a more analytical grip of this: Let  $l(x) = |c(x)|$  be the actual code length assigned to word  $x$ , which is an integer value. Let's define

$$q(x) = \frac{1}{Z} 2^{-l(x)}$$

with the normalization constraint  $Z = \sum_x 2^{-l(x)}$ . Then we have

$$\sum_{x \in X} p(x) [-\log_2 q(x)] = - \sum_x p(x) \log 2^{-l(x)} + \sum_x p(x) \log Z = \sum_x p(x) l(x) + \log Z .$$

What about  $\log Z$ ? Let  $l^{-1}(s)$  be the set of words that have been assigned codes of length  $l$ . There can only be a limited number of words encoded with a given length. For instance,  $|L^{-1}(1)|$  must not be greater than 2,  $|L^{-1}(2)|$  must not be greater than 4, and  $|l^{-1}(s)|$  must not be greater than  $2^l$ . We have

$$\forall_s : \sum_{x \in X} [l(x) = s] \leq 2^s$$

$$\forall_s : \sum_{x \in X} [l(x) = s] 2^{-s} \leq 1$$

$$\forall_s : \sum_{x \in X} 2^{-l(x)} \leq 1$$

However, this way of thinking is ok for separated codes. If such codes would be in a continuous stream of bits you'd never know where a code starts or ends. Prefix codes fix this problem by defining a code tree with leaves that clearly define when a code ends. For prefix codes it similarly holds

$$Z = \sum_{x \in X} 2^{-l(x)} \leq 1 ,$$

which is called *Kraft's inequality*. That finally gives

$$\sum_{x \in X} p(x) [-\log_2 q(x)] \leq \sum_x p(x) l(x)$$

## 5.4 The Laplace approximation: A 2nd-order Taylor of $\log p$

Assume we want to estimate some  $q(x)$  we cannot express analytically. E.g.,  $q(x) = p(x|D) \propto P(D|x)p(x)$  for some awkward likelihood function  $p(D|x)$ . An example from robotics is:  $x$  is stochastically controlled path of a robot.  $p(x)$  is a prior distribution over paths that includes how the robot can actually move and some Gaussian prior (squared costs!) over controlls. If the robot is “linear”,  $p(x)$  can be expressed nicely and analytically; if it non-linear, expressing  $p(x)$  is already hard. However,  $p(D|x)$  might indicate that we do *not* see collisions on the path—but collisions are a horrible function, usually computed by some black-box collision detection packages that computes distances between convex meshes, perhaps giving gradients but certainly not some analytic function. So  $q(x)$  can clearly not be expressed analytically.

One way to approximate  $q(x)$  is the Laplace approximation

**Definition 5.10** (Laplace approximation). Given a smooth distribution  $q(x)$ , we define its Laplace approximation as

$$\tilde{q}(x) = \exp\{-\tilde{E}(x)\},$$

where  $\tilde{E}(x)$  is the 2nd-order Taylor expansion

$$\tilde{E}(x) = E(x^*) + \frac{1}{2}(x - x^*)^\top \nabla^2 E(x^*)(x - x^*)$$

of the energy  $E(x) = -\log q(x)$  at the mode

$$x^* = \underset{x}{\operatorname{argmin}} E(x) = \underset{x}{\operatorname{argmax}} q(x).$$

First, we observe that the Laplace approximation is a Gaussian, because its energy is a parabola. Further, notice that in the Taylor expansion we skipped the linear term. That’s because we are at the mode  $x^*$  where  $\nabla E(x^*) = 0$ .

The Laplace approximation really is the probabilistic analogue of a local second-order approximation of a function, just as we used it in Newton methods. However, it is defined to be taken specifically at the mode of the distribution.

Now, computing  $x^*$  is a classical optimization problem  $x^* = \operatorname{argmin}_x E(x)$  which one might ideally solve using Newton methods. These Newton methods anyway compute the local Hessian of  $E(x)$  in every step—at the optimum we therefore have the Hessian already, which is then the precision matrix of our Gaussian.

The Laplace approximation is nice, very efficient to use, e.g., in the context of optimal control and robotics. While we can use the expressive power of probability theory to formalize the problem, the Laplace approximation brings us computationally back to efficient optimization methods.

## 5.5 Variational Inference

Another reduction of inference to optimization is variational inference.

**Definition 5.11** (variational inference). Given a distribution  $p(x)$ , and a parameterized family of distributions  $q(x|\beta)$ , the variational approximation of  $p(x)$  is defined as

$$\operatorname{argmin}_q D(q \| p)$$

## 5.6 The Fisher information metric: 2nd-order Taylor of the KLD

Recall our notion of steepest descent—it depends on the metric in the space!

Consider the space of probability distributions  $p(x; \beta)$  with parameters  $\beta$ . We think of every  $p(x; \beta)$  as a point in the space and wonder what metric is useful to compare two points  $p(x; \beta_1)$  and  $p(x; \beta_2)$ . Let's take the KLD

## 5.7 Examples and Exercises

Note: These exercises are for 'extra credits'. We'll discuss them on Thu, 21th Jan.

### 5.7.1 Maximum Entropy and Maximum Likelihood

(These are taken from MacKay's book *Information Theory...*, Exercise 22.12 & .13)

a) Assume that a random variable  $x$  with discrete domain  $\operatorname{dom}(x) = \mathcal{X}$  comes from a probability distribution of the form

$$P(x | w) = \frac{1}{Z(w)} \exp \left[ \sum_{k=1}^d w_k f_k(x) \right],$$

where the functions  $f_k(x)$  are given, and the parameters  $w \in \mathbb{R}^d$  are not known. A data set  $D = \{x_i\}_{i=1}^n$  of  $n$  points  $x$  is supplied. Show by differentiating the log likelihood  $\log P(D|w) = \sum_{i=1}^n \log P(x_i|w)$  that the maximum-likelihood parameters  $w^* = \operatorname{argmax}_w \log P(D|w)$  satisfy

$$\sum_{x \in \mathcal{X}} P(x | w^*) f_k(x) = \frac{1}{n} \sum_{i=1}^n f_k(x_i)$$

where the left-hand sum is over all  $x$ , and the right-hand sum is over the data points. A shorthand for this result is that each function-average under the fitted model must equal the function-average found in the data:

$$\langle f_k \rangle_{P(x|w^*)} = \langle f_k \rangle_D$$

b) When confronted by a probability distribution  $P(x)$  about which only a few facts are known, the maximum entropy principle (MaxEnt) offers a rule for choosing a distribution that satisfies those constraints. According to MaxEnt, you should select the  $P(x)$  that maximizes the entropy

$$H(P) = - \sum_x P(x) \log P(x)$$

subject to the constraints. Assuming the constraints assert that the averages of certain functions  $f_k(x)$  are known, i.e.,

$$\langle f_k \rangle_{P(x)} = F_k ,$$

show, by introducing Lagrange multipliers (one for each constraint, including normalization), that the maximum-entropy distribution has the form

$$P_{\text{MaxEnt}}(x) = \frac{1}{Z} \exp \left[ \sum_k w_k f_k(x) \right]$$

where the parameters  $Z$  and  $w_k$  are set such that the constraints are satisfied. And hence the maximum entropy method gives identical results to maximum likelihood fitting of an exponential-family model.

Note: The exercise will take place on Tue, 2nd Feb. Hung will also prepare how much ‘votes’ you collected in the exercises.

## 5.7.2 Maximum likelihood and KL-divergence

Assume we have a very large data set  $D = \{x_i\}_{i=1}^n$  of samples  $x_i \sim q(x)$  from some data distribution  $q(x)$ . Using this data set we can approximate any expectation

$$\langle f \rangle_q = \sum_x q(x) f(x) \approx \sum_{i=1}^n f(x_i) .$$

Assume we have a parametric family of distributions  $p(x|\beta)$  and would find the Maximum Likelihood (ML) parameter  $\beta^* = \operatorname{argmax}_{\beta} p(D|\beta)$ . Express this ML problem as a KL-divergence minimization.

### 5.7.3 Laplace Approximation

In the context of so-called “Gaussian Process Classification” the following problem arises (we neglect dependence on  $x$  here): We have a real-valued RV  $f \in \mathbb{R}$  with prior  $P(f) = \mathcal{N}(f | \mu, \sigma^2)$ . Further we have a Boolean RV  $y \in \{0, 1\}$  with conditional probability

$$P(y=1 | f) = \sigma(f) = \frac{e^f}{1 + e^f}.$$

The function  $\sigma$  is called sigmoid function, and  $f$  is a discriminative value which predicts  $y = 1$  if it is very positive, and  $y = 0$  if it is very negative. The sigmoid function has the property

$$\frac{\partial}{\partial f} \sigma(f) = \sigma(f) (1 - \sigma(f)).$$

Given that we observed  $y = 1$  we want to compute the posterior  $P(f | y=1)$ , which cannot be expressed analytically. Provide the Laplace approximation of this posterior.

(Bonus) As an alternative to the sigmoid function  $\sigma(f)$ , we can use the probit function  $\phi(z) = \int_{-\infty}^z \mathcal{N}(x|0, 1) dx$  to define the likelihood  $P(y=1 | f) = \phi(f)$ . Now how can the posterior  $P(f | y=1)$  be approximated?

### 5.7.4 Learning = Compression

In a very abstract sense, learning means to model the distribution  $p(x)$  for given data  $D = \{x_i\}_{i=1}^n$ . This is literally the case for unsupervised learning; regression, classification and graphical model learning could be viewed as specific instances of this where  $x$  factors in several random variables, like input and output.

Show in which sense the problem of learning is equivalent to the problem of compression.

### 5.7.5 A gzip experiment

Get three text files from the Web, approximately equal length, mostly text (no equations or stuff). Two of them should be in English, the third in French. (Alternatively, perhaps, not sure if it'd work, two of them on a very similar topic, the third on a very different.)

How can you use `gzip` (or some other compression tool) to estimate the mutual information between every pair of files? How can you ensure some “normalized” measures which do not depend too much on the absolute lengths of the text? Do

it and check whether in fact you find that two texts are similar while the third is different.

(Extra) Lempel-Ziv algorithms (like gzip) need to build a codebook on the fly. How does that fit into the picture?

### 5.7.6 Maximum Entropy and ML

(These are taken from MacKay's book *Information Theory...*, Exercise 22.12 & .13)

a) Assume that a random variable  $x$  with discrete domain  $\text{dom}(x) = \mathcal{X}$  comes from a probability distribution of the form

$$P(x|w) = \frac{1}{Z(w)} \exp \left[ \sum_{k=1}^d w_k f_k(x) \right],$$

where the functions  $f_k(x)$  are given, and the parameters  $w \in \mathbb{R}^d$  are not known. A data set  $D = \{x_i\}_{i=1}^n$  of  $n$  points  $x$  is supplied. Show by differentiating the log likelihood  $\log P(D|w) = \sum_{i=1}^n \log P(x_i|w)$  that the maximum-likelihood parameters  $w^* = \text{argmax}_w \log P(D|w)$  satisfy

$$\sum_{x \in \mathcal{X}} P(x|w^*) f_k(x) = \frac{1}{n} \sum_{i=1}^n f_k(x_i)$$

where the left-hand sum is over all  $x$ , and the right-hand sum is over the data points. A shorthand for this result is that each function-average under the fitted model must equal the function-average found in the data:

$$\langle f_k \rangle_{P(x|w^*)} = \langle f_k \rangle_D$$

b) When confronted by a probability distribution  $P(x)$  about which only a few facts are known, the maximum entropy principle (MaxEnt) offers a rule for choosing a distribution that satisfies those constraints. According to MaxEnt, you should select the  $P(x)$  that maximizes the entropy

$$H(P) = - \sum_x P(x) \log P(x)$$

subject to the constraints. Assuming the constraints assert that the averages of certain functions  $f_k(x)$  are known, i.e.,

$$\langle f_k \rangle_{P(x)} = F_k,$$



show, by introducing Lagrange multipliers (one for each constraint, including normalization), that the maximum-entropy distribution has the form

$$P_{\text{MaxEnt}}(x) = \frac{1}{Z} \exp \left[ \sum_k w_k f_k(x) \right]$$

where the parameters  $Z$  and  $w_k$  are set such that the constraints are satisfied. And hence the maximum entropy method gives identical results to maximum likelihood fitting of an exponential-family model.

## A Gaussian identities

### Definitions

We define a Gaussian over  $x$  with mean  $a$  and covariance matrix  $A$  as the function

$$\mathcal{N}(x | a, A) = \frac{1}{|2\pi A|^{1/2}} \exp \left\{ -\frac{1}{2} (x-a)^\top A^{-1} (x-a) \right\} \quad (90)$$

with property  $\mathcal{N}(x | a, A) = \mathcal{N}(a | x, A)$ . We also define the canonical form with precision matrix  $A$  as

$$\mathcal{N}[x | a, A] = \frac{\exp \left\{ -\frac{1}{2} a^\top A^{-1} a \right\}}{|2\pi A^{-1}|^{1/2}} \exp \left\{ -\frac{1}{2} x^\top A x + x^\top a \right\} \quad (91)$$

with properties

$$\mathcal{N}[x | a, A] = \mathcal{N}(x | A^{-1}a, A^{-1}) \quad (92)$$

$$\mathcal{N}(x | a, A) = \mathcal{N}[x | A^{-1}a, A^{-1}]. \quad (93)$$

Non-normalized Gaussian

$$\bar{\mathcal{N}}(x, a, A) = |2\pi A|^{1/2} \mathcal{N}(x | a, A) \quad (94)$$

$$= \exp \left\{ -\frac{1}{2} (x-a)^\top A^{-1} (x-a) \right\} \quad (95)$$

**Matrices** [matrix cookbook: [http://www.imm.dtu.dk/pubdb/views/edoc\\_download.php/3274/pdf/imm3274.pdf](http://www.imm.dtu.dk/pubdb/views/edoc_download.php/3274/pdf/imm3274.pdf)]

$$(A^{-1} + B^{-1})^{-1} = A (A+B)^{-1} B = B (A+B)^{-1} A \quad (96)$$

$$(A^{-1} - B^{-1})^{-1} = A (B-A)^{-1} B \quad (97)$$

$$\partial_x |A_x| = |A_x| \text{tr}(A_x^{-1} \partial_x A_x) \quad (98)$$

$$\partial_x A_x^{-1} = -A_x^{-1} (\partial_x A_x) A_x^{-1} \quad (99)$$

$$(A + UBV)^{-1} = A^{-1} - A^{-1}U(B^{-1} + VA^{-1}U)^{-1}VA^{-1} \quad (100)$$

$$(A^{-1} + B^{-1})^{-1} = A - A(B + A)^{-1}A \quad (101)$$

$$(A + J^T B J)^{-1} J^T B = A^{-1} J^T (B^{-1} + J A^{-1} J^T)^{-1} \quad (102)$$

$$(A + J^T B J)^{-1} A = \mathbf{I} - (A + J^T B J)^{-1} J^T B J \quad (103)$$

(100)=Woodbury; (102,103) holds for pos def  $A$  and  $B$

## Derivatives

$$\partial_x \mathcal{N}(x|a, A) = \mathcal{N}(x|a, A) (-h^T), \quad h := A^{-1}(x-a) \quad (104)$$

$$\partial_\theta \mathcal{N}(x|a, A) = \mathcal{N}(x|a, A) \cdot$$

$$\left[ -h^T(\partial_\theta x) + h^T(\partial_\theta a) - \frac{1}{2} \text{tr}(A^{-1} \partial_\theta A) + \frac{1}{2} h^T(\partial_\theta A) h \right] \quad (105)$$

$$\begin{aligned} \partial_\theta \mathcal{N}[x|a, A] &= \mathcal{N}[x|a, A] \left[ -\frac{1}{2} x^T \partial_\theta A x + \frac{1}{2} a^T A^{-1} \partial_\theta A A^{-1} a \right. \\ &\quad \left. + x^T \partial_\theta a - a^T A^{-1} \partial_\theta a + \frac{1}{2} \text{tr}(\partial_\theta A A^{-1}) \right] \end{aligned} \quad (106)$$

$$\partial_\theta \bar{\mathcal{N}}_x(a, A) = \bar{\mathcal{N}}_x(a, A) \cdot$$

$$\left[ h^T(\partial_\theta x) + h^T(\partial_\theta a) + \frac{1}{2} h^T(\partial_\theta A) h \right] \quad (107)$$

## Product

The product of two Gaussians can be expressed as

$$\begin{aligned} \mathcal{N}(x|a, A) \mathcal{N}(x|b, B) \\ = \mathcal{N}[x|A^{-1}a + B^{-1}b, A^{-1} + B^{-1}] \mathcal{N}(a|b, A+B), \end{aligned} \quad (108)$$

$$= \mathcal{N}(x|B(A+B)^{-1}a + A(A+B)^{-1}b, A(A+B)^{-1}B) \mathcal{N}(a|b, A+B), \quad (109)$$

$$\begin{aligned} \mathcal{N}[x|a, A] \mathcal{N}[x|b, B] \\ = \mathcal{N}[x|a+b, A+B] \mathcal{N}(A^{-1}a|B^{-1}b, A^{-1} + B^{-1}) \end{aligned} \quad (110)$$

$$= \mathcal{N}[x|\dots] \mathcal{N}[A^{-1}a|A(A+B)^{-1}b, A(A+B)^{-1}B] \quad (111)$$

$$= \mathcal{N}[x|\dots] \mathcal{N}[A^{-1}a|(1-B(A+B)^{-1})b, (1-B(A+B)^{-1})B], \quad (112)$$

$$\begin{aligned} \mathcal{N}(x|a, A) \mathcal{N}[x|b, B] \\ = \mathcal{N}[x|A^{-1}a + b, A^{-1} + B] \mathcal{N}(a|B^{-1}b, A+B^{-1}) \end{aligned} \quad (113)$$

$$= \mathcal{N}[x|\dots] \mathcal{N}[a|(1-B(A^{-1}+B)^{-1})b, (1-B(A^{-1}+B)^{-1})B] \quad (114)$$

## Convolution

$$\int_x \mathcal{N}(x|a, A) \mathcal{N}(y-x|b, B) dx = \mathcal{N}(y|a+b, A+B) \quad (115)$$

## Division

$$\mathcal{N}(x|a, A) / \mathcal{N}(x|b, B) = \mathcal{N}(x|c, C) / \mathcal{N}(c|b, C+B)$$

$$C^{-1}c = A^{-1}a - B^{-1}b$$

$$C^{-1} = A^{-1} - B^{-1} \quad (116)$$

$$\mathcal{N}[x|a, A] / \mathcal{N}[x|b, B] \propto \mathcal{N}[x|a - b, A - B] \quad (117)$$

## Expectations

Let  $x \sim \mathcal{N}(x | a, A)$ ,

$$\mathbb{E}\{x\}g(x) := \int_x \mathcal{N}(x | a, A) g(x) dx \quad (118)$$

$$\mathbb{E}\{x\}x = a, \quad \mathbb{E}\{x\}xx^\top = A + aa^\top \quad (119)$$

$$\mathbb{E}\{x\}f + Fx = f + Fa \quad (120)$$

$$\mathbb{E}\{x\}x^\top x = a^\top a + \text{tr}(A) \quad (121)$$

$$\mathbb{E}\{x\}(x-m)^\top R(x-m) = (a-m)^\top R(a-m) + \text{tr}(RA) \quad (122)$$

**Transformation** Linear transformations imply the following identities,

$$\mathcal{N}(x | a, A) = \mathcal{N}(x + f | a + f, A), \quad \mathcal{N}(x | a, A) = |F| \mathcal{N}(Fx | Fa, FAF^\top) \quad (123)$$

$$\mathcal{N}(Fx + f | a, A) = \frac{1}{|F|} \mathcal{N}(x | F^{-1}(a - f), F^{-1}AF^{-\top}) = \frac{1}{|F|} \mathcal{N}[x | F^\top A^{-1}(a - f), F^\top A^{-1}] \quad (124)$$

$$\mathcal{N}[Fx + f | a, A] = \frac{1}{|F|} \mathcal{N}[x | F^\top(a - Af), F^\top AF] \quad (125)$$

**“Propagation”** (propagating a message along a coupling, using eqs (108) and (114), respectively)

$$\int_y \mathcal{N}(x | a + Fy, A) \mathcal{N}(y | b, B) dy = \mathcal{N}(x | a + Fb, A + FBF^\top) \quad (126)$$

$$\int_y \mathcal{N}(x | a + Fy, A) \mathcal{N}[y | b, B] dy = \mathcal{N}[x | (F^{-\top} - K)(b + BF^{-1}a), (F^{-\top} - K)BF^{-1}], \quad K = \quad (127)$$

**marginal & conditional:**

$$\mathcal{N}(x | a, A) \mathcal{N}(y | b + Fx, B) = \mathcal{N}\left(x \mid \begin{array}{c} a \\ b + Fa \end{array}, \begin{array}{cc} A & A^\top F^\top \\ FA & B + FA^\top F^\top \end{array}\right) \quad (128)$$

$$\mathcal{N}\left(x \mid \begin{array}{c} a \\ b \end{array}, \begin{array}{cc} A & C \\ C^\top & B \end{array}\right) = \mathcal{N}(x | a, A) \cdot \mathcal{N}(y | b + C^\top A^{-1}(x-a), B - C^\top A^{-1}C) \quad (129)$$

$$\mathcal{N}[x | a, A] \mathcal{N}(y | b + Fx, B) = \mathcal{N} \left[ \begin{array}{c} x \\ y \end{array} \middle| \begin{array}{cc} a + F^\top B^{-1} b & A + F^\top B^{-1} F \\ B^{-1} b & -B^{-1} F \end{array}, \begin{array}{cc} -F^\top B^{-1} \\ B^{-1} \end{array} \right] \quad (130)$$

$$\mathcal{N}[x | a, A] \mathcal{N}[y | b + Fx, B] = \mathcal{N} \left[ \begin{array}{c} x \\ y \end{array} \middle| \begin{array}{cc} a + F^\top B^{-1} b & A + F^\top B^{-1} F \\ b & -F \end{array}, \begin{array}{cc} -F^\top \\ B \end{array} \right] \quad (131)$$

$$\mathcal{N} \left[ \begin{array}{c} x \\ y \end{array} \middle| \begin{array}{cc} a & A \\ b & C^\top \end{array}, \begin{array}{cc} A & C \\ C^\top & B \end{array} \right] = \mathcal{N}[x | a - CB^{-1}b, A - CB^{-1}C^\top] \cdot \mathcal{N}[y | b - C^\top x, B] \quad (132)$$

$$\left| \begin{array}{cc} A & C \\ D & B \end{array} \right| = |A| |\hat{B}| = |\hat{A}| |B|, \text{ where } \begin{array}{l} \hat{A} = A - CB^{-1}D \\ \hat{B} = B - DA^{-1}C \end{array} \quad (133)$$

$$\left[ \begin{array}{cc} A & C \\ D & B \end{array} \right]^{-1} = \left[ \begin{array}{cc} \hat{A}^{-1} & -A^{-1}C\hat{B}^{-1} \\ -\hat{B}^{-1}DA^{-1} & \hat{B}^{-1} \end{array} \right] = \left[ \begin{array}{cc} \hat{A}^{-1} & -\hat{A}^{-1}CB^{-1} \\ -B^{-1}D\hat{A}^{-1} & \hat{B}^{-1} \end{array} \right] \quad (134)$$

**pair-wise belief** We have a message  $\alpha(x) = \mathcal{N}[x|s, S]$ , transition  $P(y|x) = \mathcal{N}(y|Ax + a, Q)$ , and a message  $\beta(y) = \mathcal{N}[y|v, V]$ , what is the belief  $b(y, x) = \alpha(x)P(y|x)\beta(y)$ ?

$$b(y, x) = \mathcal{N}[x|s, S] \mathcal{N}(y|Ax + a, Q^{-1}) \mathcal{N}[y|v, V] \quad (135)$$

$$= \mathcal{N} \left[ \begin{array}{c} x \\ y \end{array} \middle| \begin{array}{ccc} s & S & 0 \\ 0 & 0 & 0 \end{array}, \begin{array}{ccc} A^\top Q^{-1} a & A^\top Q^{-1} A & -A^\top Q^{-1} \\ Q^{-1} a & -Q^{-1} A & Q^{-1} \end{array} \right] \mathcal{N} \left[ \begin{array}{c} x \\ y \end{array} \middle| \begin{array}{ccc} 0 & 0 & 0 \\ v & 0 & V \end{array} \right] \quad (136)$$

$$\propto \mathcal{N} \left[ \begin{array}{c} x \\ y \end{array} \middle| \begin{array}{ccc} s + A^\top Q^{-1} a & S + A^\top Q^{-1} A & -A^\top Q^{-1} \\ v + Q^{-1} a & -Q^{-1} A & V + Q^{-1} \end{array} \right] \quad (137)$$

## Entropy

$$H(\mathcal{N}(a, A)) = \frac{1}{2} \log |2\pi e A| \quad (138)$$

## Kullback-Leibler divergence

$$p = \mathcal{N}(x|a, A), \quad q = \mathcal{N}(x|b, B), \quad n = \dim(x), \quad D(p \| q) = \sum_x p(x) \log \frac{p(x)}{q(x)} \quad (139)$$

$$2 D(p \| q) = \log \frac{|B|}{|A|} + \text{tr}(B^{-1}A) + (b - a)^\top B^{-1}(b - a) - n \quad (140)$$

$$4 D_{\text{sym}}(p \| q) = \text{tr}(B^{-1}A) + \text{tr}(A^{-1}B) + (b - a)^\top (A^{-1} + B^{-1})(b - a) - 2n \quad (141)$$

$\lambda$ -divergence

$$2 D_\lambda(p \parallel q) = \lambda D(p \parallel \lambda p + (1-\lambda)q) + (1-\lambda) D(p \parallel (1-\lambda)p + \lambda q) \quad (142)$$

For  $\lambda = .5$ : Jensen-Shannon divergence.

## Log-likelihoods

$$\log \mathcal{N}(x|a, A) = -\frac{1}{2} \left[ \log |2\pi A| + (x-a)^\top A^{-1} (x-a) \right] \quad (143)$$

$$\log \mathcal{N}[x|a, A] = -\frac{1}{2} \left[ \log |2\pi A^{-1}| + a^\top A^{-1} a + x^\top A x - 2x^\top a \right] \quad (144)$$

$$\sum_x \mathcal{N}(x|b, B) \log \mathcal{N}(x|a, A) = -D(\mathcal{N}(b, B) \parallel \mathcal{N}(a, A)) - H(\mathcal{N}(b, B)) \quad (145)$$

**Mixture of Gaussians** Collapsing a MoG into a single Gaussian

$$\operatorname{argmin}_{b, B} D\left(\sum_i p_i \mathcal{N}(a_i, A_i) \parallel \mathcal{N}(b, B)\right) = \left(b = \sum_i p_i a_i, B = \sum_i p_i (A_i + a_i a_i^\top - b b^\top)\right) \quad (146)$$

## B 3D geometry basics (for robotics)

This document introduces to some basic geometry, focussing on 3D transformations, and introduces proper conventions for notation. There exist one-to-one implementations of the concepts and equations in libORS.

### B.1 Rotations

There are many ways to represent rotations in  $SO(3)$ . We restrict ourselves to three basic ones: rotation matrix, rotation vector, and quaternion. The rotation vector is also the most natural representation for a “rotation velocity” (angular velocities). Euler angles or raw-pitch-roll are an alternative, but they have singularities and I don’t recommend using them in practice.

**A rotation matrix** is a matrix  $R \in \mathbb{R}^{3 \times 3}$  which is orthonormal (columns and rows are orthogonal unit vectors, implying determinant 1). While a  $3 \times 3$  matrix has 9 degrees of freedom (DoFs), the constraint of orthogonality and determinant 1 constrains this: The set of rotation matrices has only 3 DoFs ( $\sim$  the local Lie algebra is 3-dim).

The application of  $R$  on a vector  $x$  is simply the matrix-vector product  $Rx$ .

Concatenation of two rotations  $R_1$  and  $R_2$  is the normal matrix-matrix product  $R_1 R_2$ .

Inversion is the transpose,  $R^{-1} = R^T$ .

**A rotation vector** is an unconstrained vector  $w \in \mathbb{R}^3$ . The vector's direction  $\underline{w} = \frac{w}{|w|}$  determines the rotation axis, the vector's length  $|w| = \theta$  determines the rotation angle (in radians, using the right thumb convention).

The application of a rotation described by  $w \in \mathbb{R}^3$  on a vector  $x \in \mathbb{R}^3$  is given as (Rodrigues' formula)

$$w \cdot x = \cos \theta x + \sin \theta (\underline{w} \times x) + (1 - \cos \theta) \underline{w} (\underline{w}^T x) \quad (147)$$

where  $\theta = |w|$  is the rotation angle and  $\underline{w} = w/\theta$  the unit length rotation axis.

The inverse rotation is described by the negative of the rotation vector.

Concatenation is non-trivial in this representation and we don't discuss it here. In practice, a rotation vector is first converted to a rotation matrix or quaternion.

Conversion to a matrix: For every vector  $w \in \mathbb{R}^3$  we define its skew symmetric matrix as

$$\hat{w} = \begin{pmatrix} 0 & -w_3 & w_2 \\ w_3 & 0 & -w_1 \\ -w_2 & w_1 & 0 \end{pmatrix}. \quad (148)$$

Note that such skew-symmetric matrices are related to the cross product:  $w \times v = \hat{w} v$ , where the cross product is rewritten as a matrix product. The rotation matrix  $R(w)$  that corresponds to a given rotation vector  $w$  is:

$$R(w) = \exp(\hat{w}) \quad (149)$$

$$= \cos \theta I + \sin \theta \hat{w} / \theta + (1 - \cos \theta) w w^T / \theta^2 \quad (150)$$

The exp function is called exponential map (generating a group element (=rotation matrix) via an element of the Lie algebra (=skew matrix)). The other formular is called Rodrigues' formular: the first term is a diagonal matrix ( $I$  is the 3D identity matrix), the second terms the skew symmetric part, the last term the symmetric part ( $w w^T$  is also called outer product).

**Angular velocity & derivative of a rotation matrix:** We represent angular velocities by a vector  $w \in \mathbb{R}^3$ , the direction  $\underline{w}$  determines the rotation axis, the length  $|w|$  is the rotation velocity (in radians per second). When a body's orientation at time  $t$  is described by a rotation matrix  $R(t)$  and the body's angular velocity is  $w$ , then

$$\dot{R}(t) = \hat{w} R(t). \quad (151)$$

(That’s intuitive to see for a rotation about the  $x$ -axis with velocity 1.) Some insights from this relation: Since  $R(t)$  must always be a rotation matrix (fulfill orthogonality and determinant 1), its derivative  $\dot{R}(t)$  must also fulfill certain constraints; in particular it can only live in a 3-dimensional sub-space. It turns out that the derivative  $\dot{R}$  of a rotation matrix  $R$  must always be a skew symmetric matrix  $\hat{w}$  times  $R$  – anything else would be inconsistent with the constraints of orthogonality and determinant 1.

Note also that, assuming  $R(0) = I$ , the solution to the differential equation  $\dot{R}(t) = \hat{w} R(t)$  can be written as  $R(t) = \exp(t\hat{w})$ , where here the exponential function notation is used to denote a more general so-called exponential map, as used in the context of Lie groups. It also follows that  $R(w)$  from (149) is the rotation matrix you get when you rotate for 1 second with angular velocity described by  $w$ .

**Quaternion** (I’m not describing the general definition, only the “quaternion to represent rotation” definition.) A quaternion is a unit length 4D vector  $r \in \mathbb{R}^4$ ; the first entry  $r_0$  is related to the rotation angle  $\theta$  via  $r_0 = \cos(\theta/2)$ , the last three entries  $\hat{r} \equiv r_{1:3}$  are related to the unit length rotation axis  $\underline{w}$  via  $\hat{r} = \sin(\theta/2) \underline{w}$ . The inverse of a quaternion is given by negating  $\hat{r}$ ,  $r^{-1} = (r_0, -\hat{r})$  (or, alternatively, negating  $r_0$ ).

The concatenation of two rotations  $r, r'$  is given as the quaternion product

$$r \circ r' = (r_0 r'_0 - \hat{r}^\top \hat{r}', r_0 \hat{r}' + r'_0 \hat{r} + \hat{r}' \times \hat{r}) \tag{152}$$

The application of a rotation quaternion  $r$  on a vector  $x$  can be expressed by converting the vector first to the quaternion  $(0, x)$ , then computing

$$r \cdot x = (r \circ (0, x) \circ r^{-1})_{1:3} , \tag{153}$$

I think a bit more efficient is to first convert the rotation quaternion  $r$  to the equivalent rotation matrix  $R$ , as given by

$$R = \begin{pmatrix} 1 - r_{22} - r_{33} & r_{12} - r_{03} & r_{13} + r_{02} \\ r_{12} + r_{03} & 1 - r_{11} - r_{33} & r_{23} - r_{01} \\ r_{13} - r_{02} & r_{23} + r_{01} & 1 - r_{11} - r_{22} \end{pmatrix} \\ r_{ij} := 2r_i r_j . \tag{154}$$

(Note: In comparison to (149) this does not require to compute a sin or cos.) Inversely, the quaternion  $r$  for a given matrix  $R$  is

$$r_0 = \frac{1}{2} \sqrt{1 + \text{tr}R} \tag{155}$$

$$r_3 = (R_{21} - R_{12}) / (4r_0) \tag{156}$$

$$r_2 = (R_{13} - R_{31}) / (4r_0) \tag{157}$$

$$r_1 = (R_{32} - R_{23}) / (4r_0) . \tag{158}$$

**Angular velocity**  $\rightarrow$  **quaternion velocity** Given an angular velocity  $w \in \mathbb{R}^3$  and a current quaternion  $r(t) \in \mathbb{R}$ , what is the time derivative  $\dot{r}(t)$  (in analogy to Eq. (151))? For simplicity, let's first assume  $|w| = 1$ . For a small time interval  $\delta$ ,  $w$  generates a rotation vector  $\delta w$ , which converts to a quaternion

$$\Delta r = (\cos(\delta/2), \sin(\delta/2)w) . \quad (159)$$

That rotation is concatenated LHS to the original quaternion,

$$r(t + \delta) = \Delta r \circ r(t) . \quad (160)$$

Now, if we take the derivative w.r.t.  $\delta$  and evaluate it at  $\delta = 0$ , all the  $\cos(\delta/2)$  terms become  $-\sin(\delta/2)$  and evaluate to zero, all the  $\sin(\delta/2)$  terms become  $\cos(\delta/2)$  and evaluate to one, and we have

$$\dot{r}(t) = \frac{1}{2}(-w^\top \hat{r}, r_0 w + \hat{r} \times w) = \frac{1}{2}(0, w) \circ r(t) \quad (161)$$

Here  $(0, w) \in \mathbb{R}^4$  is a four-vector; for  $|w| = 1$  it is a normalized quaternion. However, due to the linearity the equation holds for any  $w$ .

**Quaternion velocity**  $\rightarrow$  **angular velocity** The following is relevant when taking the derivative w.r.t. the parameters of a quaternion, e.g., for a ball joint represented as quaternion. Given  $\dot{r}$ , we have

$$\dot{r} \circ r^{-1} = \frac{1}{2}(0, w) \circ r \circ r^{-1} = \frac{1}{2}(0, w) \quad (162)$$

which allows us to read off the angular velocity induced by a change of quaternion. However, the RHS zero will hold true only iff  $\dot{r}$  is orthogonal to  $r$  (where  $\dot{r}^\top r = \dot{r}_0 r_0 + \hat{r}^\top \dot{\hat{r}} = 0$ , see (152)). In case  $\dot{r}^\top r \neq 0$ , the change in length of the quaternion does not represent any angular velocity; in typical kinematics engines a non-unit length is ignored. Therefore one first orthogonalizes  $\dot{r} \leftarrow \dot{r} - r(\dot{r}^\top r)$ .

As a special case of application, consider computing the partial derivative w.r.t. quaternion coordinates, where  $\dot{r}$  is the unit vectors  $e_0, \dots, e_3$ . In this case, the orthogonalization becomes simply  $e_i \leftarrow e_i - r r_i$  and

$$(e_i - r_i r) \circ r^{-1} = e_i \circ r^{-1} - r_i(1, 0, 0, 0) \quad (163)$$

$$w_i = 2[e_i \circ r^{-1}]_{1:3} , \quad (164)$$

where  $w_i$  is the rotation vector implied by  $\dot{r} = e_i$ . In case the original quaternion  $r$  wasn't normalized (which could be, if a standard optimization algorithm searches in the quaternion configuration space), then  $r$  actually represents the normalized quaternion  $\hat{r} = \frac{1}{\sqrt{r^2}} r$ , and (due to linearity of the above), the rotation vector implied by  $\dot{r} = e_i$  is

$$w_i = \frac{2}{\sqrt{r^2}} [e_i \circ r^{-1}]_{1:3} . \quad (165)$$



## B.2 Transformations

We consider two types of transformations here: either static (translation+rotation), or dynamic (translation+velocity+rotation+angular velocity). The first maps between two static reference frames, the latter between moving reference frames, e.g. between reference frames attached to moving rigid bodies.

### B.2.1 Static transformations

Concerning the static transformations, again there are different representations:

A **homogeneous matrix** is a  $4 \times 4$ -matrix of the form

$$T = \begin{pmatrix} R & t \\ 0 & 1 \end{pmatrix} \quad (166)$$

where  $R$  is a  $3 \times 3$ -matrix (rotation in our case) and  $t$  a 3-vector (translation).

In homogeneous coordinates, vectors  $x \in \mathbb{R}^3$  are expanded to 4D vectors  $\begin{pmatrix} x \\ 1 \end{pmatrix} \in \mathbb{R}^4$  by appending a 1.

Application of a transform  $T$  on a vector  $x \in \mathbb{R}^3$  is then given as the normal matrix-vector product

$$x' = T \cdot x = T \begin{pmatrix} x \\ 1 \end{pmatrix} = \begin{pmatrix} R & t \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ 1 \end{pmatrix} = \begin{pmatrix} Rx + t \\ 1 \end{pmatrix}. \quad (167)$$

Concatenation is given by the ordinary 4-dim matrix-matrix product.

The inverse transform is

$$T^{-1} = \begin{pmatrix} R & t \\ 0 & 1 \end{pmatrix}^{-1} = \begin{pmatrix} R^{-1} & -R^{-1}t \\ 0 & 1 \end{pmatrix} \quad (168)$$

**Translation and quaternion:** A transformation can efficiently be stored as a pair  $(t, r)$  of a translation vector  $t$  and a rotation quaternion  $r$ . Analogous to the above, the application of  $(t, r)$  on a vector  $x$  is  $x' = t + r \cdot x$ ; the inverse is  $(t, r)^{-1} = (-r^{-1} \cdot t, r^{-1})$ ; the concatenation is  $(t_1, r_1) \circ (t_2, r_2) = (t_1 + r_1 \cdot t_2, r_1 \circ r_2)$ .

### B.2.2 Dynamic transformations

Just as static transformations map between (static) coordinate frames, dynamic transformations map between moving (inertial) frames which are, e.g., attached to moving bodies. A dynamic transformation is described by a tuple  $(t, r, v, w)$  with translation  $t$ , rotation  $r$ , velocity  $v$  and angular velocity  $w$ . Under a dynamic transform

$(t, r, v, w)$  a position and velocity  $(x, \dot{x})$  maps to a new position and velocity  $(x', \dot{x}')$  given as

$$x' = t + r \cdot x \quad (169)$$

$$\dot{x}' = v + w \times (r \cdot x) + r \cdot \dot{x} \quad (170)$$

(the second term is the additional linear velocity of  $\dot{x}'$  arising from the angular velocity  $w$  of the dynamic transform). The concatenation  $(t, r, v, w) = (t_1, r_1, v_1, w_1) \circ (t_2, r_2, v_2, w_2)$  of two dynamic transforms is given as

$$t = t_1 + r_1 \cdot t_2 \quad (171)$$

$$v = v_1 + w_1 \times (r_1 \cdot t_2) + r_1 \cdot v_2 \quad (172)$$

$$r = r_1 \circ r_2 \quad (173)$$

$$w = w_1 + r_1 \cdot w_2 \quad (174)$$

For completeness, the footnote<sup>9</sup> also describes how accelerations transform, including the case when the transform itself is accelerating. The inverse  $(t', r', v', w') = (t, r, v, w)^{-1}$  of a dynamic transform is given as

$$t' = -r^{-1} \cdot t \quad (181)$$

$$r' = r^{-1} \quad (182)$$

$$v' = r^{-1} \cdot (w \times t - v) \quad (183)$$

$$w' = -r^{-1} \cdot w \quad (184)$$

**Sequences of transformations** by  $T_{A \rightarrow B}$  we denote the transformation from frame  $A$  to frame  $B$ . The frames  $A$  and  $B$  can be thought of coordinate frames (tuples of an offset (in an affine space) and three local orthonormal basis vectors) attached to two bodies  $A$  and  $B$ . It holds

$$T_{A \rightarrow C} = T_{A \rightarrow B} \circ T_{B \rightarrow C} \quad (185)$$

where  $\circ$  is the concatenation described above. Let  $p$  be a point (rigorously, in the affine space). We write  $p^A$  for the coordinate vector of point  $p$  relative to

---

<sup>9</sup>Transformation of accelerations:

$$\begin{aligned} \dot{v} &= \dot{v}_1 + \dot{w}_1 \times (r_1 \cdot t_2) + w_1 \times (w_1 \times (r_1 \cdot t_2)) \\ &\quad + 2 w_1 \times (r_1 \cdot v_2) + r_1 \cdot \dot{v}_2 \end{aligned} \quad (175)$$

$$\dot{w} = \dot{w}_1 + w_1 \times (r_1 \cdot w_2) + r_1 \cdot \dot{w}_2 \quad (176)$$

Used identities: for any vectors  $a, b, c$  and rotation  $r$ :

$$r \cdot (a \times b) = (r \cdot a) \times (r \cdot b) \quad (177)$$

$$a \times (b \times c) = b(ac) - c(ab) \quad (178)$$

$$\partial_t(r \cdot a) = w \times (r \cdot a) + r \cdot \dot{a} \quad (179)$$

$$\partial_t(w \times a) = \dot{w} \times t + w \times \dot{a} \quad (180)$$

frame  $A$ ; and  $p^B$  for the coordinate vector of point  $p$  relative to frame  $B$ . It holds

$$p^A = T_{A \rightarrow B} p^B. \quad (186)$$

### B.2.3 A note on affine coordinate frames

Instead of the notation  $T_{A \rightarrow B}$ , other text books often use notations such as  $T_{AB}$  or  $T_B^A$ . A common question regarding notation  $T_{A \rightarrow B}$  is the following:

*The notation  $T_{A \rightarrow B}$  is confusing, since it transforms coordinates from frame  $B$  to frame  $A$ . Why not the other way around?*

I think the notation  $T_{A \rightarrow B}$  is intuitive for the following reasons. The core is to understand that a transformation can be thought of in two ways: as a transformation of the *coordinate frame itself*, and as transformation of the *coordinates relative to a coordinate frame*. I'll first give a non-formal explanation and later more formal definitions of affine frames and their transformation.

Think of  $T_{W \rightarrow B}$  as translating and rotating a real rigid body: First, the body is located at the world origin; then the body is moved by a translation  $t$ ; then the body is rotated (around its own center) as described by  $R$ . In that sense,  $T_{W \rightarrow B} = \begin{pmatrix} R & t \\ 0 & 1 \end{pmatrix}$  describes the “forward” transformation of the body. Consider that a coordinate frame  $B$  is attached to the rigid body and a frame  $W$  to the world origin. Given a point  $p$  in the world, we can express its coordinates relative to the world,  $p^W$ , or relative to the body  $p^B$ . You can convince yourself with simple examples that  $p^W = T_{W \rightarrow B} p^B$ , that is,  $T_{W \rightarrow B}$  also describes the “backward” transformation of body-relative-coordinates to world-relative-coordinates.

Formally: Let  $(A, V)$  be an affine space. A coordinate frame is a tuple  $(o, e_1, \dots, e_n)$  of an origin  $o \in A$  and basis vectors  $e_i \in V$ . Given a point  $p \in A$ , its coordinates  $p_{1:n}$  w.r.t. a coordinate frame  $(o, e_1, \dots, e_n)$  are given implicitly via

$$p = o + \sum_i p_i e_i. \quad (187)$$

A transformation  $T_{W \rightarrow B}$  is a (“forward”) transformation of the coordinate frame itself:

$$(o^B, e_1^B, \dots, e_n^B) = (o^W + t, R e_1^W, \dots, R e_n^W) \quad (188)$$

where  $t \in V$  is the affine translation in  $A$  and  $R$  the rotation in  $V$ . Note that the coordinates  $(e_i^B)_{1:n}^W$  of a basis vector  $e_i^B$  relative to frame  $W$  are the columns of  $R$ :

$$e_i^B = \sum_j (e_i^B)_j^W e_j^W = \sum_j R_{ji} e_j^W \quad (189)$$

Given this transformation of the coordinate frame itself, the coordinates transform as follows:

$$p = o^W + \sum_i p_i^W e_i^W \quad (190)$$

$$p = o^B + \sum_i p_i^B e_i^B \quad (191)$$

$$= o^W + t + \sum_i p_i^B (Re_i^W) \quad (192)$$

$$= o^W + \sum_i t_i^W e_i^W + \sum_j p_j^B (Re_j^W) \quad (193)$$

$$= o^W + \sum_i t_i^W e_i^W + \sum_j p_j^B \left( \sum_i R_{ij} e_i^W \right) \quad (194)$$

$$= o^W + \sum_i \left[ t_i^W + \sum_j R_{ij} p_j^B \right] e_i^W \quad (195)$$

$$\Rightarrow p_i^W = t_i^W + \sum_j R_{ij} p_j^B. \quad (196)$$

Another way to express this formally:  $T_{W \rightarrow B}$  maps *covariant* vectors (including “basis vectors”) forward, but *contra-variant* vectors (including “coordinate vectors”) backward.

## B.3 Kinematic chains

In this section we only consider static transformations composed of translation and rotation. But given the general concatenation rules above, everything said here generalizes directly to the dynamic case.

### B.3.1 Rigid and actuated transforms

A actuated kinematic chain with  $n$  joints is a series of transformations of the form

$$T_{W \rightarrow 1} \circ Q_1 \circ T_{1 \rightarrow 2} \circ Q_2 \circ T_{2 \rightarrow 3} \circ Q_3 \circ \dots \quad (197)$$

Each  $T_{i-1 \rightarrow i}$  describes so-called “links” (or bones) of the kinematic chain: the rigid (non-actuated) transformation from the  $i-1$ th joint to the  $i$ th joint. The first  $T_{W \rightarrow 1}$  is the transformation from *world* coordinates to the first joint.

Each  $Q_i$  is the actuated transformation of the joint – usually simply a rotation around the joint’s x-axis with a specific angle, the so-called *joint angle*. These joint angles (and therefore each  $Q_i$ ) are actuated and may change over time.

When we control the robot we essentially tell it to actuate its joint so as to change the joint angles. There are two fundamental computations necessary for control:

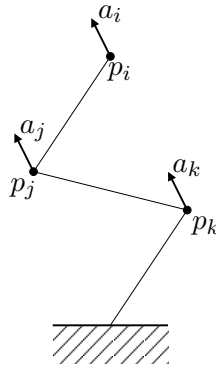


Figure 6: Illustration for the Jacobian and Hessian.

- (i) For a given  $n$ -dimensional vector  $q \in \mathbb{R}^n$  of joint angles, compute the absolute frames  $T_{W \rightarrow i}$  (*world* to link transformation) of each link  $i$ .
- (ii) For a given  $n$ -dimensional vector  $\dot{q} \in \mathbb{R}^n$  of joint angle velocities, compute absolute (*world*-relative) velocity and angular velocity of the  $i$ th link.

The first problem is solved by “forward chaining” the transformations: we can compute the absolute transforms  $T_{W \rightarrow i}$  (i.e., the transformation from *world* to the  $i$ th link) for each link, namely:

$$T_{W \rightarrow i} = T_{W \rightarrow i-1} \circ Q_i \circ T_{i-1 \rightarrow i}. \quad (198)$$

Iterating this for  $i = 2, \dots, n$  we get positions and orientations  $T_{W \rightarrow i}$  of all links in world coordinates.

The second problem is addressed in the next section.

### B.3.2 Jacobian & Hessian

Assume we have computed the absolute position and orientation of each link in an actuated kinematic chain. Then we want to know how a point  $p_i^W$  attached to the  $i$ th frame (and coordinates expressed w.r.t. the world frame  $W$ ) changes when rotating the  $j$ th joint. Equally we want to know how some arbitrary vector  $a_i^W$  attached to the  $i$ th frame (coordinates relative to the world frame  $W$ ) rotates when rotating the  $j$ th joint. Let  $a_j^W$  be the unit length rotation axis of the  $j$ th joint (which, by convention, is  $r_{W \rightarrow j} \cdot (1, 0, 0)$  if  $r_{W \rightarrow j}$  is the rotation in  $T_{W \rightarrow j}$ ). In the following we drop the superscript  $W$  because all coordinates are expressed in the world frame. Let  $q_j$  be the joint angle of the  $j$ th joint. For the purpose of computing a partial derivative of something attached to link  $i$  w.r.t. the actuation at joint  $j$  we can think

of everything inbetween as rigid. Hence, the point and vector Jacobian and Hessian are simply:

$$d_{ij} := p_i - p_j$$

$$\frac{\partial p_i}{\partial q_j} = a_j \times d_{ij} \quad (199)$$

$$\frac{\partial a_i}{\partial q_j} = \alpha_j \times a_i \quad (200)$$

$$\begin{aligned} \frac{\partial^2 p_i}{\partial q_j \partial q_k} &= \frac{\partial a_j}{\partial q_k} \times d_{ij} + a_j \times \frac{\partial d_{ij}}{\partial q_k} \\ &= (a_k \times a_j) \times d_{ij} + a_j \times [a_k \times (p_i - p_k) - a_k \times (p_j - p_k)] \\ &= (a_k \times a_j) \times d_{ij} + a_j \times (a_k \times d_{ij}) \\ &\quad \left[ \text{using } a \times (b \times c) + b \times (c \times a) + c \times (a \times b) = 0 \right] \\ &= a_k \times (a_j \times d_{ij}) \end{aligned} \quad (201)$$

$$\begin{aligned} \frac{\partial^2 a_i}{\partial q_j \partial q_k} &= (a_k \times \alpha_j) \times a_i + \alpha_j \times (a_k \times a_i) \\ &= a_k \times (\alpha_j \times a_i) \end{aligned} \quad (202)$$

Efficient computation: Assume we have an articulated kinematic *tree* of links and joints. We write  $j < i$  if joint  $j$  is inward from link (or joint)  $i$ . Then, for each body  $i$  consider all inward edges  $j < i$  and further inward edges  $k \leq j$  and compute  $H_{i,jk} = \partial_{\theta_j} \partial_{\theta_k} p_i$ . Note that  $H_{i,jk}$  is symmetric in  $j$  and  $k$  – so computing for  $k \leq j$  and copying the rest is sufficient.