

Chapter 29

Mathematical background

ap,math

Contents

29.1 Introduction (s,math,intro)	29.2
29.2 Fourier transforms (s,math,four)	29.2
29.2.1 Properties	29.2
29.2.1.1 Linearity	29.2
29.2.1.2 Convolution property	29.2
29.2.1.3 Shift property	29.2
29.2.1.4 Scaling property	29.2
29.2.1.5 Affine scaling property	29.3
29.2.1.6 Rotation property	29.3
29.2.1.7 Parseval's theorem	29.3
29.2.1.8 Differentiation property	29.3
29.2.1.9 Laplacian property	29.3
29.2.1.10 Circular symmetry	29.3
29.2.1.11 Separability	29.3
29.2.1.12 Hermitian symmetry property	29.3
29.3 Discrete Fourier transform (s,math,dsp)	29.3
29.3.1 Properties of DFT	29.3
29.4 Laplace transform (s,math,laplace)	29.4
29.4.1 Laplace transform pairs	29.4
29.4.2 Laplace transform properties	29.4
29.5 Differential equations (s,math,diff)	29.5
29.6 Z transform (s,math,z)	29.6
29.6.1 Z transform pairs	29.6
29.6.2 Z transform properties	29.6
29.7 B-splines (s,math,spline)	29.6
29.8 Gradients (s,math,gradient)	29.7
29.8.1 Gradients of linear and quadratic forms	29.7
29.8.2 Taylor series expansions	29.7
29.8.3 Lipschitz continuity	29.8
29.9 Convexity (s,math,convex)	29.9
29.9.1 Convex sets	29.9
29.9.2 Convex projections	29.9
29.9.3 Convex functions (s,math,convex,fun)	29.10
29.9.3.1 Minimizers of convex functions	29.10
29.9.3.2 Assessing convexity using properties	29.11
29.9.3.3 Properties of convex functions	29.11
29.9.3.4 Properties of differentiable convex functions	29.12

29.9.3.5	Convex conjugate	29.14
29.9.3.6	Proximal mapping	29.14
29.9.3.7	Moreau envelope	29.14
29.9.3.8	Minimization of convex functions over convex sets	29.15
29.9.3.9	Exchanging order of minimization and maximization	29.15
29.10	Minimizers with nonnegativity constraints (s,math,min,nonneg)	29.15
29.11	Augmented Lagrangian methods (s,math,al)	29.16
29.12	Special functions (s,math,floor)	29.16
29.13	Convergence rates of iterations (s,math,rate)	29.17
29.14	Problems (s,math,prob)	29.18
29.15	Bibliography	29.18

s,math,intro

29.1 Introduction (s,math,intro)

This appendix reviews some basic mathematical tools and notation used in this book.

s,math,four

29.2 Fourier transforms (s,math,four)

The \bar{d} -dimensional Fourier transform $F(\vec{\nu})$ of a function $f(\vec{x})$, where $\vec{x} \in \mathbb{R}^{\bar{d}}$ and $\vec{\nu} \in \mathbb{R}^{\bar{d}}$, and the inverse Fourier transform are given by:

$$F(\vec{\nu}) = \int f(\vec{x}) e^{-i2\pi\vec{\nu}\cdot\vec{x}} d\vec{x} \quad (29.2.1)$$

$$f(\vec{x}) = \int F(\vec{\nu}) e^{i2\pi\vec{\nu}\cdot\vec{x}} d\vec{\nu}. \quad (29.2.2)$$

e,math,four,ft

29.2.1 Properties

29.2.1.1 Linearity

$$\sum_j \alpha_j f_j(\vec{x}) \xleftrightarrow{\text{FT}} \sum_j \alpha_j F_j(\vec{\nu})$$

29.2.1.2 Convolution property

The convolution of two \bar{d} -dimensional signals $f(\vec{x})$ and $g(\vec{x})$ is denoted

$$(f * g)(\vec{x}) = \int_{\mathbb{R}^{\bar{d}}} f(\vec{x}') g(\vec{x} - \vec{x}') d\vec{x}'. \quad (29.2.3)$$

e,math,four,conv

The **convolution property** of the **Fourier transform** is:

$$h(\vec{x}) = (f * g)(\vec{x}) \xleftrightarrow{\text{FT}} H(\vec{\nu}) = F(\vec{\nu}) G(\vec{\nu}). \quad (29.2.4)$$

e,math,four,,conv

29.2.1.3 Shift property

$$g(\vec{x}) = f(\vec{x} - \vec{x}_0) \xleftrightarrow{\text{FT}} G(\vec{\nu}) = e^{-i2\pi\vec{\nu}\cdot\vec{x}_0} F(\vec{\nu}).$$

29.2.1.4 Scaling property

$$g(\vec{x}) = f(\alpha\vec{x}) \xleftrightarrow{\text{FT}} G(\vec{\nu}) = \frac{1}{|\alpha|^{\bar{d}}} F\left(\frac{\vec{\nu}}{\alpha}\right), \quad \alpha \neq 0.$$

29.2.1.5 Affine scaling property

For an invertible $\bar{d} \times \bar{d}$ matrix \mathbf{B} : (see Problem 29.1):

$$g(\vec{x}) = f(\mathbf{B}^{-1}(\vec{x} - \vec{x}_0)) \xleftrightarrow{\text{FT}} G(\vec{\nu}) = |\det\{\mathbf{B}\}| e^{-i2\pi\vec{\nu} \cdot \vec{x}_0} F(\mathbf{B}'\vec{\nu}). \quad (29.2.5) \quad \text{e,math,four,affine}$$

29.2.1.6 Rotation property

For an orthonormal $\bar{d} \times \bar{d}$ matrix \mathbf{U} :

$$g(\vec{x}) = f(\mathbf{U}\vec{x}) \xleftrightarrow{\text{FT}} G(\vec{\nu}) = F(\mathbf{U}\vec{\nu}). \quad (29.2.6) \quad \text{e,math,four,rotate}$$

29.2.1.7 Parseval's theorem

$$\int f(\vec{x}) g^*(\vec{x}) d\vec{x} = \int F(\vec{\nu}) G^*(\vec{\nu}) d\vec{\nu}. \quad (29.2.7) \quad \text{e,math,four,parseval}$$

29.2.1.8 Differentiation property

$$g(\vec{x}) = \frac{\partial}{\partial x_j} f(\vec{x}) \xleftrightarrow{\text{FT}} G(\vec{\nu}) = i2\pi\nu_j F(\vec{\nu}).$$

29.2.1.9 Laplacian property

$$g(\vec{x}) = \Delta^2 f(\vec{x}) = \sum_{j=1}^{\bar{d}} \frac{\partial^2}{\partial x_j^2} f(\vec{x}) \xleftrightarrow{\text{FT}} G(\vec{\nu}) = \sum_{j=1}^{\bar{d}} -(2\pi\nu_j)^2 F(\vec{\nu}) = -(2\pi)^2 \|\vec{\nu}\|^2 F(\vec{\nu}). \quad (29.2.8) \quad \text{e,math,four,laplace}$$

29.2.1.10 Circular symmetry

If $f(\vec{x}) = f_0(\|\vec{x}\|)$ then $F(\vec{\nu}) = F_0(\|\vec{\nu}\|)$.

29.2.1.11 Separability

If $f(\vec{x}) = \prod_{i=1}^{\bar{d}} f_i(x_i)$ then $F(\vec{\nu}) = \prod_{i=1}^{\bar{d}} F_i(\nu_i)$, where $f_i(x) \xleftrightarrow{\text{FT}} F_i(\nu)$.

29.2.1.12 Hermitian symmetry property

If $f(\vec{x})$ is **real**, then its Fourier transform is **Hermitian symmetric**: $F(-\vec{\nu}) = F^*(\vec{\nu})$.

29.3 Discrete Fourier transform (s,math,dsp)

The **1D N -point discrete Fourier transform (DFT)** of a signal $x[n]$ is defined by

$$X[k] = \sum_{n=0}^{N-1} x[n] e^{-i\frac{2\pi}{N}nk}, \quad k = 0, \dots, N-1, \quad (29.3.1) \quad \text{e,math,dsp,dft}$$

and the inverse DFT is given by

$$x[n] = \frac{1}{N} \sum_{k=0}^{N-1} X[k] e^{i\frac{2\pi}{N}nk}, \quad n = 0, \dots, N-1. \quad (29.3.2) \quad \text{e,math,dsp,idft}$$

29.3.1 Properties of DFT

- **Shift (periodic)**

$$x[(n - n_0) \bmod N] \xleftrightarrow{\text{DFT}} e^{-i\frac{2\pi}{N}n_0} X[k].$$

- **Parseval's theorem**

$$\sum_{n=0}^{N-1} |x[n]|^2 = \frac{1}{N} \sum_{k=0}^{N-1} |X[k]|^2. \quad (29.3.3) \quad \text{e,math,dsp,parseval}$$

29.4 Laplace transform (s,math,laplace)

The **Laplace transform** is useful for analyzing the properties of 1D linear time-invariant systems. For signal processing applications, the **bilateral Laplace transform** is usually more relevant than the **unilateral Laplace transform**. Therefore we focus on the bilateral, or **two-sided Laplace transform** here.

Although there exists a formula for the **inverse Laplace transform**, typically all that is needed is to combine Laplace transform properties with known transform pairs such as those summarized below.

29.4.1 Laplace transform pairs

signal	transform	ROC	Notes
$h(t)$	$H(s) = \int_{-\infty}^{\infty} h(t) e^{-st} dt$		
$\delta(t)$	1	\mathbb{C}	
$e^{-at} \text{step}(t)$	$\frac{1}{s+a}$	$-\text{real}\{a\} < \text{real}\{s\}$	
$e^{-a t }$	$\frac{a^2 - s^2}{s^2 + a^2}$	$ \text{real}\{s\} < \text{real}\{a\}$	
$\cos(bt) e^{-at} \text{step}(t)$	$\frac{s}{(s+a)^2 + b^2}$	$-\text{real}\{a\} < \text{real}\{s\}$	Problem 29.2
$\sin(bt) e^{-at} \text{step}(t)$	$\frac{b}{(s+a)^2 + b^2}$	$-\text{real}\{a\} < \text{real}\{s\}$	Problem 29.3
$\cos(bt) e^{-a t }$	$\frac{2a(a^2 + b^2 - s^2)}{(a^2 + b^2 - s^2)^2 + (2bs)^2}$	$ \text{real}\{s\} < \text{real}\{a\}$	Problem 29.4
$\sin(b t) e^{-a t }$	$\frac{2b(a^2 + b^2 + s^2)}{(a^2 + b^2 - s^2)^2 + (2bs)^2}$	$ \text{real}\{s\} < \text{real}\{a\}$	Problem 29.5
$(\cos(at) + \sin(a t)) e^{-a t }$	$\frac{8a^3}{s^4 + 4a^4}$	$ \text{real}\{s\} < \text{real}\{a\}$	

29.4.2 Laplace transform properties

property	signal	transform	ROC
linearity	$ah(t) + bg(t)$	$aH(s) + bG(s)$	$\text{ROC} \supseteq \text{ROC}_1 \cap \text{ROC}_2$
differentiation	$\frac{d}{dt} h(t)$	$sH(s)$	$\text{ROC} \supseteq \text{ROC}_h$
convolution	$h(t) * g(t)$	$H(s)G(s)$	$\text{ROC} \supseteq \text{ROC}_h \cap \text{ROC}_g$
time shift	$h(t - \tau)$	$e^{-\tau s} H(s)$	ROC unchanged
time scale	$h(at), a \neq 0$	$\frac{1}{ a } H(s/a)$	$\text{ROC} = \text{ROC}_h/a$
modulation	$e^{bt} h(t)$	$H(s - b)$	$\text{ROC} = \text{ROC}_h + \text{real}\{b\}$
differentiation in s	$-th(t)$	$\frac{d}{ds} H(s)$	ROC unchanged
running integration	$\int_{-\infty}^t h(t') dt'$	$\frac{1}{s} H(s)$	$\text{ROC} \supseteq \text{ROC}_h \cap \{\text{real}\{s\} > 0\}$

29.5 Differential equations (s,math,diff)

If the time-varying state vector $\mathbf{x}(t)$ evolves according to the linear, constant-coefficient differential equation

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{u}(t),$$

with initial condition $\mathbf{x}(0) = \mathbf{x}_0$, then it is readily verified that the solution for $t \geq 0$ is

$$\mathbf{x}(t) = e^{\mathbf{A}t} \mathbf{x}_0 + \int_0^t e^{\mathbf{A}(t-s)} \mathbf{u}(s) \, ds. \quad (29.5.1)$$

If \mathbf{A} has eigen-decomposition $\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^{-1}$, then $e^{\mathbf{A}s} = \mathbf{V} \mathbf{\Lambda}^s \mathbf{V}^{-1}$. This is useful, for example, when solving the **Bloch equation** in NMR.

If the scalar function $x(t)$ evolves according to the linear, *time-varying* differential equation

$$\dot{x}(t) = a(t)x(t) + u(t),$$

with initial condition $x(0) = x_0$, then it is readily verified that the solution for $t \geq 0$ is

$$x(t) = \exp\left(\int_0^t a(s) \, ds\right) x_0 + \int_0^t \exp\left(\int_s^t a(\tau) \, d\tau\right) u(s) \, ds.$$

Unfortunately, there is no simple generalization of this result to the vector state case:

$$\dot{\mathbf{x}}(t) = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{u}(t).$$

Define $\mathbf{B}(t) \triangleq \int_0^t \mathbf{A}(s) \, ds$ and consider the function

$$\mathbf{f}(t) = \exp(\mathbf{B}(t)) \mathbf{f}_0 = \left(\mathbf{I} + \sum_{k=1}^{\infty} \frac{1}{k!} [\mathbf{B}(t)]^k \right) \mathbf{f}_0.$$

Its derivative is

$$\dot{\mathbf{f}}(t) = \left(\sum_{k=1}^{\infty} \frac{1}{k!} \left(\frac{d}{dt} [\mathbf{B}(t)]^k \right) \right) \mathbf{f}_0.$$

This simplifies if the matrix $\frac{d}{dt} \mathbf{B}(t)$ commutes with $\mathbf{B}(t)$, which happens in the scalar case but not in general.

29.6 Z transform (s,math,z)

The **Z transform** is useful for analyzing the properties of discrete-time 1D linear time-invariant systems.

Although there exists a formula for the **inverse Z transform**, typically all that is needed is to combine Z transform properties with known transform pairs such as those summarized below.

29.6.1 Z transform pairs (s,math,z,pair)

In the following table, $p = r e^{ib} \in \mathbb{C}$, where $r, b \in \mathbb{R}$.

signal	transform	ROC
$h[n]$	$\sum_{n=-\infty}^{\infty} h[n] z^{-n}$	
$\delta[n]$	1	\mathbb{C}
$p^n \text{step}(n)$	$\frac{1}{1 - pz^{-1}}$	$ p < z $
$p^{ n }$	$\frac{(p - p^{-1})z^{-1}}{1 - (p + p^{-1})z^{-1} + z^{-2}} = \frac{(p - p^{-1})z^{-1}}{(1 - pz^{-1})(1 - p^*z^{-1})}$	$ p < z < \frac{1}{ p }$
$\cos(bn) r^n \text{step}(n)$	$\frac{1 - r \cos(b) z^{-1}}{1 - 2r \cos(b) z^{-1} + r^2 z^{-2}} = \frac{1 - \text{real}\{p\} z^{-1}}{(1 - pz^{-1})(1 - p^*z^{-1})}$	$ r < z $
$\sin(bn) r^n \text{step}(n)$	$\frac{r \sin(b) z^{-1}}{1 - 2r \cos(b) z^{-1} + r^2 z^{-2}} = \frac{\text{imag}\{p\} z^{-1}}{(1 - pz^{-1})(1 - p^*z^{-1})}$	$ r < z $
$\cos(bn) r^{ n }$	$z^{-1} \frac{(1 + z^{-2})(r - 1/r) \cos(b) - (r^2 - 1/r^2)z^{-1}}{(1 - pz^{-1})(1 - 1/pz^{-1})(1 - p^*z^{-1})(1 - 1/p^*z^{-1})}$	$ r < z < \frac{1}{ r }$
$\sin(b n) r^{ n }$	$z^{-1} \frac{(1 + z^{-2})(r + 1/r) \sin(b) - 2 \sin(2b) z^{-1}}{(1 - pz^{-1})(1 - 1/pz^{-1})(1 - p^*z^{-1})(1 - 1/p^*z^{-1})}$	$ r < z < \frac{1}{ r }$

29.6.2 Z transform properties (s,math,z,prop)

property	signal	transform	ROC
linearity	$ah[n] + bg[n]$	$aH(z) + bG(z)$	$\text{ROC} \supseteq \text{ROC}_1 \cap \text{ROC}_2$
differencing	$h[n] - h[n-1]$	$(1 - z^{-1})H(z)$	$\text{ROC} \supseteq \text{ROC}_h$
convolution	$h[n] * g[n]$	$H(z)G(z)$	$\text{ROC} \supseteq \text{ROC}_h \cap \text{ROC}_g$
time shift	$h[n-m]$	$z^{-m}H(z)$	ROC unchanged, except origin
time reversal	$h[-n]$	$H(z^{-1})$	ROC inverted
modulation	$b^n h[n]$	$H(z/b)$	$\text{ROC} = \text{ROC}_h / b $
time expansion (upsampling), $k \in \mathbb{N}$	$h[n/k] \mathbb{I}_{\{n/k \in \mathbb{Z}\}}$	$H(z^k)$	$\text{ROC} = \text{ROC}_h^{1/k}$
differentiation in z	$-nh[n]$	$-\frac{d}{dz}H(z)$	ROC unchanged
accumulation	$\sum_{k=-\infty}^n h[k]$	$\frac{1}{1-z^{-1}}H(z)$	$\text{ROC} \supseteq \text{ROC}_h \cap \{ z > 1\}$

29.7 B-splines (s,math,spline)

B-spline functions of order $k \geq 0$ with **knots** $t_1, t_2, \dots, t_{N+k+1}$ are **defined recursively** as follows for $i = 1, \dots, N$:

$$B_{i,0}(t) = \mathbb{I}_{\{t_i \leq t < t_{i+1}\}}$$

$$B_{i,k}(t) = \frac{t - t_i}{t_{i+k} - t_i} B_{i,k-1}(t) + \frac{t_{i+1+k} - t}{t_{i+1+k} - t_{i+1}} B_{i+1,k-1}(t). \quad (29.7.1)$$

In image processing, often the knots are spaced equally.

MIRT See `ir_bspline_k.m`.

29.8 Gradients (s,math,gradient)

If $f(\mathbf{x})$ is a (differentiable) function from \mathbb{R}^{n_p} to \mathbb{R} , then the **row gradient** of f is defined as

$$\nabla f(\mathbf{x}) = \left[\frac{\partial}{\partial x_1} f(\mathbf{x}), \dots, \frac{\partial}{\partial x_{n_p}} f(\mathbf{x}) \right]. \quad (29.8.1)$$

The **column gradient**, denoted ∇f is the transpose of the row gradient.

29.8.1 Gradients of linear and quadratic forms

If $f(\mathbf{x}) = \text{real}\{\mathbf{v}' \mathbf{x}\}$ for some vector $\mathbf{v} \in \mathbb{C}^{n_p}$, then (see §29.2):

$$\nabla f(\mathbf{x}) = \mathbf{v}.$$

If $f(\mathbf{x}) = \mathbf{x}' \mathbf{M} \mathbf{x}$ for some matrix $\mathbf{M} \in \mathbb{R}^{n_p \times n_p}$, then

$$\nabla f(\mathbf{x}) = (\mathbf{M} + \mathbf{M}') \mathbf{x}. \quad (29.8.2)$$

See Appendix 29 for the complex case.

Leibniz's rule:

$$G(x) = \int_{a(x)}^{b(x)} h(x, y) dy \implies \frac{d}{dx} G(x) = h(x, b(x)) \frac{d}{dx} b(x) - h(x, a(x)) \frac{d}{dx} a(x) + \int_{a(x)}^{b(x)} \frac{\partial}{\partial x} h(x, y) dy.$$

29.8.2 Taylor series expansions

If $g : \mathbb{R} \rightarrow \mathbb{R}$ is an n -times differentiable function, then **Taylor's theorem** is:

$$g(x) = g(a) + \sum_{k=1}^{n-1} \frac{1}{k!} g^{(k)}(a) (x-a)^k + \int_a^x \frac{1}{(n-1)!} (x-t)^{n-1} g^{(n)}(t) dt.$$

In particular, for $n = 2$ we have

$$\begin{aligned} g(t) &= g(s) + \dot{g}(s)(t-s) + \int_s^t (t-\tau) \ddot{g}(\tau) d\tau \\ &= g(s) + \dot{g}(s)(t-s) + (t-s)^2 \int_0^1 (1-\alpha) \ddot{g}(\alpha t + (1-\alpha)s) d\alpha. \end{aligned}$$

Conversely, if

$$g(t) = g(s) + \dot{g}(s)(t-s) + (t-s)^2 \int_0^1 (1-\alpha) \ddot{g}(\alpha t + (1-\alpha)s) d\alpha,$$

then $\ddot{g}(t) = \ddot{g}(s)$.

If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable, then the **1st-order Taylor series** expansion of f around a point \mathbf{z} is

$$f(\mathbf{x}) = f(\mathbf{z}) + \left[\int_0^1 \nabla f(\alpha \mathbf{x} + (1-\alpha)\mathbf{z}) d\alpha \right] (\mathbf{x} - \mathbf{z}). \quad (29.8.3)$$

If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice differentiable, then the **2nd-order Taylor series** expansion of f around a point \mathbf{z} is

$$f(\mathbf{x}) = f(\mathbf{z}) + \nabla f(\mathbf{z})(\mathbf{x} - \mathbf{z}) + (\mathbf{x} - \mathbf{z})' \left[\int_0^1 (1-\alpha) \nabla^2 f(\alpha \mathbf{x} + (1-\alpha)\mathbf{z}) d\alpha \right] (\mathbf{x} - \mathbf{z}). \quad (29.8.4)$$

For functions with complex arguments, see §29.4.

29.8.3 Lipschitz continuity

A function $g : \mathbb{C}^n \rightarrow \mathbb{C}^m$ is called **Lipschitz continuous** if there exists a finite real number \mathcal{L} such that

$$\|g(x) - g(z)\|_2 \leq \mathcal{L} \|x - z\|_2, \quad \forall x, z \in \mathbb{C}^n. \quad (29.8.5)$$

Such an \mathcal{L} is called a **Lipschitz constant**. If g is **Lipschitz continuous**, then it is natural to seek the smallest **Lipschitz constant**:

$$\mathcal{L}_* \triangleq \sup_{x \neq z} \frac{\|g(x) - g(z)\|_2}{\|x - z\|_2},$$

which one might call “the” **Lipschitz constant** of g .

A function need *not* be differentiable to be Lipschitz continuous. For example, $g(x) = |x|$ is Lipschitz continuous with $\mathcal{L} = 1$. However, any **Lipschitz continuous** function is **absolutely continuous** and thus differentiable **almost everywhere**. We generalize this definition in Definition 29.9.16.

29.9 Convexity (s,math,convex)

Convex sets and convex functions have important roles in optimization. This section reviews some relevant properties.

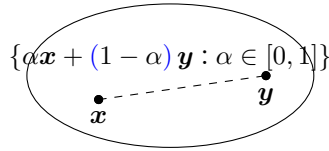
29.9.1 Convex sets (s,math,convex,set)

Definition 29.9.1 A set \mathcal{C} is **convex** iff

$$\mathbf{x}, \mathbf{z} \in \mathcal{C} \implies \alpha \mathbf{x} + (1 - \alpha) \mathbf{z} \in \mathcal{C}, \quad \forall \alpha \in [0, 1].$$

In words, for any two points in a convex set, all points on the **line segment** between those points also lie in the set.

Example 29.9.2 An ellipse is a convex set, as illustrated in the following figure.



Example 29.9.3 Any **box set** of the following form is convex:

$$\{\mathbf{x} \in \mathbb{R}^n : l_j \leq x_j \leq u_j, \quad j = 1, \dots, n\}. \quad (29.9.1)$$

Example 29.9.4 The **nonnegative orthant** is the most important convex set used in this book:

$$\{\mathbf{x} \in \mathbb{R}^n : 0 \leq x_j < \infty, \quad j = 1, \dots, n\}. \quad (29.9.2)$$

29.9.2 Convex projections (math,convex,projection)

A particularly important property of convex sets is that for any point \mathbf{x} outside the set, there is a unique point $\mathbf{x}^{(*)}$ within the set that is closest to \mathbf{x} , as established by the following theorem from functional analysis [1, p. 69].

Theorem 29.9.5 Let \mathcal{C} be a nonempty closed convex subset of a Hilbert space H (such as \mathbb{R}^n or \mathbb{C}^n) with associated inner product $\langle \cdot, \cdot \rangle$ and norm $\|\cdot\|$. For any $\mathbf{x} \in H$, there is a unique vector $\mathbf{x}^{(*)} \in \mathcal{C}$ such that

$$\|\mathbf{x} - \mathbf{x}^{(*)}\| \leq \|\mathbf{x} - \mathbf{z}\|, \quad \forall \mathbf{z} \in \mathcal{C}.$$

Furthermore, $\mathbf{x}^{(*)}$ is characterized (in a necessary and sufficient sense) by

$$\operatorname{real}\{\langle \mathbf{x} - \mathbf{x}^{(*)}, \mathbf{z} - \mathbf{x}^{(*)} \rangle\} \leq 0, \quad \forall \mathbf{z} \in \mathcal{C}. \quad (29.9.3)$$

This theorem is a close relative of the **projection theorem**. Because of its existence and uniqueness results, it is valid to define the following **projector** onto a convex set \mathcal{C} :

$$\mathcal{P}_{\mathcal{C}}(\mathbf{x}) \triangleq \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|. \quad (29.9.4)$$

This function gives the closest point in \mathcal{C} to \mathbf{x} .

Example 29.9.6 If \mathcal{C} is the box set defined in (29.9.1) and $\mathbf{x}^{(*)} = \mathcal{P}_{\mathcal{C}}(\mathbf{x})$, then

$$\mathbf{x}_j^* = \begin{cases} l_j, & x_j < l_j \\ x_j, & l_j \leq x_j \leq u_j \\ u_j, & x_j > u_j. \end{cases}$$

29.9.3 Convex functions (s,math,convex,fun)

Definition 29.9.7 For a convex set $\mathcal{D} \subset \mathbb{R}^n$, a real functional $f : \mathcal{D} \rightarrow \mathbb{R}$ is called a **convex function** on \mathcal{D} iff

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{z}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{z}), \quad \forall \alpha \in [0, 1], \quad \forall \mathbf{x}, \mathbf{z} \in \mathcal{D}.$$

The function f is called **strictly convex** iff the inequality is strict for all $\mathbf{x} \neq \mathbf{z}$ and $\alpha \in (0, 1)$.

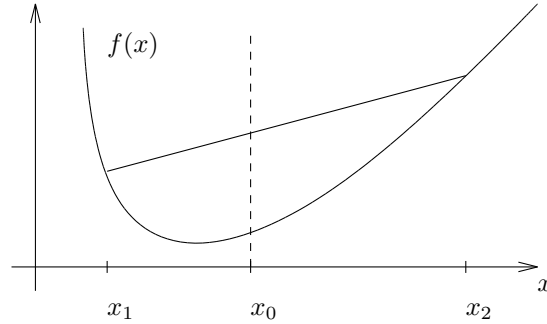


Figure 29.9.1: Illustration of a 1D convex function.

An alternative (equivalent) definition is that a function is convex iff its **epigraph** (the set of points on or above the graph of the function, i.e., $\{(\mathbf{x}, t) : \mathbf{x} \in \mathcal{D}, f(\mathbf{x}) \leq t \in \mathbb{R}\}$) is a **convex set** [2, p. 75].

Definition 29.9.8 A function f is called **concave** if its negative, $-f$, is convex.

A **strictly concave** function is defined by analogy.

Example 29.9.9 A parabola $f(x) = a + bx + cx^2$ with $c \geq 0$ is convex on \mathbb{R} , because for $\alpha \in [0, 1]$:

$$\alpha f(x) + (1 - \alpha) f(y) - f(\alpha x + (1 - \alpha) y) = c\alpha(1 - \alpha)(x - y)^2 \geq 0.$$

Example 29.9.10 The function $f(x) = |x|^p$ is convex for $p \geq 1$ [2, p. 71] and is strictly convex for $p > 1$. (See Problem 29.7.)

For generalizations of convexity, see [3] (**g-convex**) and [2] (**quasi-convex**).

29.9.3.1 Minimizers of convex functions

Definition 29.9.11

$\mathbf{x}^{(*)}$ is a **local minimizer** of $f(\mathbf{x})$ with respect to norm $\|\cdot\|$ iff $\exists \varepsilon > 0$ s.t. $\|\mathbf{x} - \mathbf{x}^{(*)}\| \leq \varepsilon \implies f(\mathbf{x}^{(*)}) \leq f(\mathbf{x})$ [4, p. 19].

$\mathbf{x}^{(*)} \in \mathcal{D}$ is a **global minimizer** of $f(\mathbf{x})$ over the domain \mathcal{D} iff $f(\mathbf{x}^{(*)}) \leq f(\mathbf{x}) \forall \mathbf{x} \in \mathcal{D}$.

Convex functions are particularly important for optimization problems because of the following properties.

- Any **local minimizer** of a convex function is a global minimizer.
- A global minimizer of a strictly convex function is **unique**.

The converse of uniqueness is not true; i.e., there are many functions that have unique global minimizers that are not strictly convex, such as the ℓ_1 norm: $f(\mathbf{x}) = \|\mathbf{x}\|_1$, and even non-convex, such $f(x) = \sqrt{|x|}$. Convexity is sufficient *but not necessary* for ensuring uniqueness of a global minimizer.

Neither of these properties ensures that local or global minimizers **exist** for convex functions. For example, the 1D function $f(x) = e^x$ is strictly convex but has no local (or global) minimizers on \mathbb{R} .

For differentiable convex functions on \mathbb{R}^n we can characterize the **global minimizer(s)** using the gradient.

- If $\Psi(\mathbf{x})$ is **convex** on \mathbb{R}^n , then $\nabla \Psi(\mathbf{x}) = \mathbf{0}$ iff \mathbf{x} is a global minimizer of Ψ .
- If $\Psi(\mathbf{x})$ is **strictly convex** on \mathbb{R}^n , then $\nabla \Psi(\mathbf{x}) = \mathbf{0}$ iff \mathbf{x} is the (unique) global minimizer of Ψ .

In other words, to find an **unconstrained minimizer** of a convex function over \mathbb{R}^n , it suffices to find a value of \mathbf{x} where the gradient is zero. In contrast, for general (possibly non-convex) functions, finding a point \mathbf{x} where the gradient is zero tells us very little because \mathbf{x} could be a local minimizer or a local maximizer or neither (a saddle point).

29.9.3.2 Assessing convexity using properties

One way to determine if a function is convex is to resort to the definition. Often it is simpler to use *properties* of the function such as the following.

- A function $f(\mathbf{x})$ is convex on \mathcal{D} iff its restriction $g(t) = f(\mathbf{x} + t\mathbf{z})$ to a line in \mathcal{D} is convex for any \mathbf{x} and \mathbf{z} in \mathcal{D} [2, p. 68].
- If f is twice differentiable, then f is convex if and only if its **Hessian matrix** is positive-semidefinite:

$$\nabla^2 f \succeq \mathbf{0}. \quad (29.9.5) \quad \text{e, math, convex, fun, hess, mgeq}$$

(See Lemma 29.5.7 and Lemma 29.5.8.)

- If f is twice differentiable, then f is **strictly convex** if its Hessian matrix is positive definite:

$$\nabla^2 f \succ \mathbf{0}. \quad (29.9.6) \quad \text{e, math, convex, fun, hess, mgt}$$

However, the converse is not true. Example: $f(x) = x^4$ is strictly convex, but its Hessian (second derivative) is not positive everywhere.

Often convex functions are constructed from other (convex) functions and the following properties apply.

- If f is convex on \mathbb{R}^n , then $g(\mathbf{x}) \triangleq f(\mathbf{M}\mathbf{x} + \mathbf{b})$ is convex on \mathbb{R}^m for any matrix $\mathbf{M} \in \mathbb{R}^{n \times m}$ and vector $\mathbf{b} \in \mathbb{R}^n$. In other words, convexity is preserved under affine transformations. (See Problem 29.6.)
- If $f(\mathbf{x})$ and $g(\mathbf{x})$ are convex, then so is their point-wise maximum: $h(\mathbf{x}) = \max \{f(\mathbf{x}), g(\mathbf{x})\}$.
- Sums of convex functions are convex.
- Sums of strictly convex functions are strictly convex.
- The sum of a convex function and a strictly convex function is strictly convex.

x, convex, fun, sum

Example 29.9.12 The following cost function is **strictly convex** for any $\beta > 0$: $\Psi(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \beta \|\mathbf{x}\|_2^2$. The Hessian of the Tikhonov regularizer $g(\mathbf{x}) = \beta \|\mathbf{x}\|_2^2$ is $\nabla^2 g = \beta \mathbf{I}$, which is positive definite, so g is strictly convex. The data-fit term $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ is convex by the affine property above, so the sum is strictly convex.

29.9.3.3 Properties of convex functions

- By induction, one can show that if f is convex and $\sum_{k=1}^K \alpha_k = 1$ for $\alpha_k \geq 0$, then f obeys the following **convexity inequality**

$$f\left(\sum_{k=1}^K \alpha_k \mathbf{x}_k\right) \leq \sum_{k=1}^K \alpha_k f(\mathbf{x}_k), \quad (29.9.7) \quad \text{e, math, convex, convexity, inequality}$$

for \mathbf{x}_k values in \mathcal{D} , the domain of f .

- A special case of the preceding result is **Jensen's inequality**. If $f : \mathcal{D} \rightarrow \mathbb{R}$ is a convex function and $\mathbf{X} \in \mathbb{R}^n$ is a random vector with $P\{\mathbf{X} \in \mathcal{D}\} = 1$, then

$$f(E[\mathbf{X}]) \leq E[f(\mathbf{X})]. \quad (29.9.8) \quad \text{e, jensen}$$

- The **sublevel sets** $\{\mathbf{x} \in \mathcal{D} : f(\mathbf{x}) < a\}$ of a convex function f are convex sets for any $a \in \mathbb{R}$.

x, convex, fun, mean

Example 29.9.13 Applying (29.9.7) to the function $f(x) = -\log(x)$ yields the weighted **arithmetic-geometric mean inequality** [5, p. 53a]:

$$\sum_k \alpha_k x_k \geq \prod_k x_k^{\alpha_k}, \quad \forall x_k \geq 0,$$

if $\sum_k \alpha_k = 1$ and $\alpha_k \geq 0$, with equality iff all x_k values are equal.

29.9.3.4 Properties of differentiable convex functions

- If f is convex on \mathcal{D} and differentiable at $z \in \mathcal{D}$, then it satisfies the following **support property**:

$$f(x) \geq f(z) + \langle \nabla f(z), x - z \rangle, \quad \forall x \in \mathcal{D}.$$

e,math,convex,above,tangent
(29.9.9)

In other words, a convex function lies above the tangent plane at any point.

- Conversely, if f is differentiable on \mathcal{D} and satisfies (29.9.9) everywhere, then f is convex [2, p. 70].

The following lemmas are used to prove convergence of the **gradient projection** algorithm in §12.2.

Lemma 29.9.14 [6, p. 10] *If $f : \mathcal{D} \rightarrow \mathbb{R}$ is convex and differentiable on \mathcal{D} , then*

$$\langle \nabla f(x) - \nabla f(z), x - z \rangle \geq 0, \quad \forall x, z \in \mathcal{D}.$$

Proof:

f is convex $\implies f(z) \geq f(x) + \langle \nabla f(x), z - x \rangle$ and $f(x) \geq f(z) + \langle \nabla f(z), x - z \rangle$. Now add. □

Definition 29.9.15 *A function $g : \mathbb{C}^n \rightarrow \mathbb{C}^m$ is Lipschitz continuous w.r.t. norms $\|\cdot\|_m$ and $\|\cdot\|_n$ on \mathbb{C}^m and \mathbb{C}^n respectively iff*

$$\|g(x) - g(z)\|_m \leq \|x - z\|_n, \quad \forall x, z \in \mathbb{C}^n.$$

e,math,convex,fun,lip,mn
(29.9.10)

Definition 29.9.16 *A function $g : \mathbb{C}^n \rightarrow \mathbb{C}^m$ is (A, B) -Lipschitz continuous iff*

$$\|A(g(x) - g(z))\|_2 \leq \|B(x - z)\|_2, \quad \forall x, z \in \mathbb{C}^n,$$

e,math,convex,fun,lip,ab
(29.9.11)

where A has m columns and $B \in \mathbb{C}^{n \times n}$ is invertible.

Our primary interest in Lipschitz functions will be gradients of cost functions, which are mappings from \mathbb{C}^{n_p} into \mathbb{C}^{n_p} (or \mathbb{R}^{n_p} into \mathbb{R}^{n_p}). The following definition generalizes (29.8.5) and is a special case of Definition 29.9.16.

Definition 29.9.17 *A function $g : \mathbb{C}^{n_p} \rightarrow \mathbb{C}^{n_p}$ is S -Lipschitz continuous, for an invertible $n_p \times n_p$ matrix S , iff*

$$\|S^{-1}(g(x) - g(z))\|_2 \leq \|S'(x - z)\|_2, \quad \forall x, z \in \mathbb{C}^{n_p}.$$

e,math,convex,lip,s
(29.9.12)

This definition is a strict generalization of the usual definition of Lipschitz continuity because if (29.9.12) holds, then

$$\begin{aligned} \|g(x) - g(z)\|_2 &= \|SS^{-1}(g(x) - g(z))\|_2 \leq \|S\|_2 \|S^{-1}(g(x) - g(z))\|_2 \leq \|S\|_2 \|S'(x - z)\|_2 \\ &\leq \|S\|_2 \|S'\|_2 \|x - z\|_2, \end{aligned}$$

so $g(\cdot)$ is Lipschitz continuous with Lipschitz constant $L = \|S\|_2 \|S'\|_2$.

Theorem 29.9.18 *If Ψ is twice differentiable, then the Hessian bound $\|S^{-1} \nabla^2 \Psi(x) S^{-H}\|_2 \leq 1, \forall x \in \mathbb{R}^{n_p}$ holds iff its gradient is S -Lipschitz continuous per (29.9.12).*

Proof (extended from [7, p. 21]):

Using Taylor expansion with remainder (29.8.4) and the triangle inequality shows (29.9.12) as follows:

$$\begin{aligned} \|S^{-1}(\nabla \Psi(z) - \nabla \Psi(x))\| &= \left\| S^{-1} \int_0^1 \nabla^2 \Psi(x + \tau(z - x)) S^{-H} S'(z - x) d\tau \right\| \\ &\leq \left(\int_0^1 \|S^{-1} \nabla^2 \Psi(x + \tau(z - x)) S^{-H}\|_2 d\tau \right) \|S'(z - x)\| \\ &\leq \|S'(z - x)\|. \end{aligned}$$

Lemma 29.9.19 provides a (generalized) converse. □

Lemma 29.9.19 *If $g : \mathbb{C}^n \rightarrow \mathbb{C}^m$ is differentiable at $x \in \mathbb{C}^n$ and (A, B) -Lipschitz continuous per (29.9.11), then*

$$\|A \nabla g(x) B^{-1}\| \leq 1.$$

Proof:

Because \mathbf{g} is differentiable at \mathbf{x} , there exists $\mathbf{H}(\mathbf{x}) = \nabla \mathbf{g}(\mathbf{x}) \in \mathbb{C}^{m \times n}$ such that

$$\forall \varepsilon > 0, \exists \delta_\varepsilon > 0 \text{ s.t. } \|\mathbf{d}\| \leq \delta_\varepsilon \implies \|\mathbf{g}(\mathbf{x} + \mathbf{d}) - \mathbf{g}(\mathbf{x}) - \mathbf{H}(\mathbf{x})\mathbf{d}\| \leq \varepsilon \|\mathbf{d}\|.$$

Suppose $\|\mathbf{A}\mathbf{H}\mathbf{B}^{-1}\| > 1$. Then there exists $\mathbf{z} \neq \mathbf{0}$ and $c > 0$ such that $\|\mathbf{A}\mathbf{H}\mathbf{B}^{-1}\mathbf{z}\| = (1 + c)\|\mathbf{z}\|$. Consider $\varepsilon = c/(2\|\mathbf{A}\|\|\mathbf{B}^{-1}\|)$ and choose α such that $\alpha\|\mathbf{B}^{-1}\mathbf{z}\| < \delta_\varepsilon$. (The case where $\|\mathbf{A}\| = 0$ is trivial.) Defining $\mathbf{d} = \alpha\mathbf{B}^{-1}\mathbf{z}$ and using the triangle inequality:

$$\begin{aligned} (1 + c)\|\alpha\mathbf{z}\| &= \|\mathbf{A}\mathbf{H}\mathbf{B}^{-1}\alpha\mathbf{z}\| \leq \|\mathbf{A}(\mathbf{g}(\mathbf{x} + \mathbf{d}) - \mathbf{g}(\mathbf{x}) - \mathbf{H}\mathbf{d})\| + \|\mathbf{A}(\mathbf{g}(\mathbf{x} + \mathbf{d}) - \mathbf{g}(\mathbf{x}))\| \\ &\leq \|\mathbf{A}\|\varepsilon\|\mathbf{d}\| + \|\mathbf{B}\mathbf{d}\| \leq (c/2)\|\alpha\mathbf{z}\| + \|\alpha\mathbf{z}\|, \end{aligned}$$

a contradiction. So we must have $\|\mathbf{A}\mathbf{H}\mathbf{B}^{-1}\| \leq 1$. □

Corollary 29.9.20 *If Ψ is twice differentiable and convex, and \mathbf{S} is invertible, then the Lipschitz condition (29.9.12) holds iff $\nabla^2 \Psi(\mathbf{x}) \preceq \mathbf{S}\mathbf{S}'$, $\forall \mathbf{x}$. (Problem 29.9.)*

In other words, for a twice differentiable, convex cost function Ψ , the Lipschitz condition (29.9.12) is equivalent to a bound on the curvature (Hessian) of Ψ .

The condition in Definition 29.9.17 leads to the following inequality related to Lemma 29.9.14. (See Problem 29.8.)

Lemma 29.9.21 *If Ψ is convex and differentiable and its gradient is \mathbf{S} -Lipschitz continuous per Definition 29.9.17, then*

$$\langle \nabla \Psi(\mathbf{x}) - \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle \leq \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|^2. \quad (29.9.13)$$

Lemma 29.9.22 *If Ψ is convex and differentiable and its gradient is \mathbf{S} -Lipschitz continuous per Definition 29.9.17, then we can strengthen Lemma 29.9.14 to*

$$\langle \nabla \Psi(\mathbf{x}) - \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle \geq \|\mathbf{S}^{-1}(\nabla \Psi(\mathbf{x}) - \nabla \Psi(\mathbf{z}))\|^2. \quad (29.9.14)$$

Proof (extended from [6, p. 24]):

Consider first the case where Ψ is twice differentiable. Then by the 2nd-order Taylor expansion (29.4.3): $\nabla \Psi(\mathbf{x}) = \nabla \Psi(\mathbf{z}) + \mathbf{H}(\mathbf{x}, \mathbf{z})(\mathbf{x} - \mathbf{z})$ where $\mathbf{H}(\mathbf{x}, \mathbf{z}) = \int_0^1 \nabla^2 \Psi(\tau \mathbf{x} + (1 - \tau)\mathbf{z}) d\tau \succeq \mathbf{0}$ by the convexity of Ψ . Defining $\mathbf{A} \triangleq \mathbf{H}(\mathbf{x}, \mathbf{z}) = \mathbf{S}^{-1} \mathbf{H}(\mathbf{x}, \mathbf{z}) \mathbf{S}^{-\text{H}} \succeq \mathbf{0}$, by Lemma 29.9.19 we have $\|\mathbf{A}\| \leq 1$. Note that $\mathbf{A}^2 \preceq \|\mathbf{A}\| \mathbf{A} \preceq \mathbf{A}$, so $\mathbf{w}' \mathbf{A} \mathbf{w} \geq \|\mathbf{A} \mathbf{w}\|^2$. Thus with $\mathbf{w} = \mathbf{S}'(\mathbf{x} - \mathbf{z})$:

$$\begin{aligned} \langle \nabla \Psi(\mathbf{x}) - \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle &= \langle \mathbf{H}(\mathbf{x}, \mathbf{z})(\mathbf{x} - \mathbf{z}), \mathbf{x} - \mathbf{z} \rangle = \langle \mathbf{S}^{-1} \mathbf{H} \mathbf{S}^{-\text{H}} \mathbf{w}, \mathbf{w} \rangle = \mathbf{w}' \mathbf{A} \mathbf{w} \\ &\geq \|\mathbf{A} \mathbf{w}\|^2 = \|\mathbf{S}^{-1} \mathbf{H}(\mathbf{x}, \mathbf{z})(\mathbf{x} - \mathbf{z})\|^2 = \|\mathbf{S}^{-1}(\nabla \Psi(\mathbf{x}) - \nabla \Psi(\mathbf{z}))\|^2. \end{aligned}$$

If Ψ is not twice differentiable¹, then consider the smoothed function

$$\Psi_\varepsilon(\mathbf{x}) = \int \Psi(\mathbf{x} - \mathbf{z}) p(\mathbf{z}/\varepsilon) \varepsilon^{-n_p} d\mathbf{z}$$

where $p(\cdot)$ is a smooth (i.e., C^∞) kernel with finite support and unit integral. Clearly Ψ_ε is convex and twice differentiable and its gradient satisfies the same Lipschitz condition as $\nabla \Psi$. Thus Ψ_ε satisfies the inequality (29.9.14), and taking the limit as $\varepsilon \rightarrow 0^+$ establishes (29.9.14) for Ψ . □

Lemma 29.9.23 *If Ψ is convex and differentiable and its gradient is \mathbf{S} -Lipschitz continuous per Definition 29.9.17, then [7, Thm. 2.1.5]:*

$$\Psi(\mathbf{x}) \leq \Psi(\mathbf{z}) + \langle \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle + \frac{1}{2} \|\mathbf{S}(\nabla \Psi(\mathbf{x}) - \nabla \Psi(\mathbf{z}))\|^2.$$

This property was used as a definition in [8].

¹ Thanks to Arkadi Nemirovski for help with this proof, in a 2001-08-20 email to Matt Jacobson.

29.9.3.5 Convex conjugate

The **convex conjugate** of a function, also known as the **Legendre-Fenchel transformation**, is used for deriving some optimization methods. (See §14.8.) The definition applies to general normed spaces but here we consider only real-valued functions on \mathbb{R}^n . For such functions, the **convex conjugate** is

$$f^*(z) \triangleq \sup_x (\langle z, x \rangle - f(x)), \quad \forall z \in \mathbb{R}^n. \quad (29.9.15)$$

Example 29.9.24 The convex conjugate of the absolute value function $f(x) = |x|$ is

$$f^*(z) = \begin{cases} 0, & |z| \leq 1 \\ \infty, & |z| > 1. \end{cases}$$

Example 29.9.25 The convex conjugate of a power function

$$f(x) = \frac{1}{p} |x|^p, \quad 1 < p < \infty$$

is

$$f^*(z) = \frac{1}{q} |z|^q, \quad \frac{1}{p} + \frac{1}{q} = 1.$$

Properties of convex conjugates include the following.

- $f^*(x)$ is convex on \mathbb{R}^n .
- If $f(x)$ is a **proper, convex**, and **lower-semicontinuous** function, then $f^{**} = f$ by the **Fenchel-Moreau theorem**, so

$$f(x) = \sup_z (\langle z, x \rangle - f^*(z)). \quad (29.9.16)$$

In particular we have the following inequality:

$$f(x) \geq \langle z, x \rangle - f^*(z), \quad \forall z \in \mathbb{R}^n. \quad (29.9.17)$$

29.9.3.6 Proximal mapping

If f is a convex function, then the **proximal operator** or **proximal mapping Moreau proximity operator** is defined [9–12] as

$$\text{prox}_f(z) \triangleq \arg \min_x \frac{1}{2} \|x - z\|_2^2 + f(x). \quad (29.9.18)$$

Because $\|x - z\|_2$ is strictly convex and $f(x)$ is convex, their sum is strictly convex so that sum has a unique minimizer for any z . This operator is particularly useful for the **proximal gradient method (PGM)** for non-smooth optimization and its relatives [13, 14].

If f is differentiable, then by differentiating (29.9.18) we have

$$0 = x - z + \nabla f(x) \Big|_{x=\text{prox}_f(z)} \implies \text{prox}_f(z) = (I + \nabla f)^{-1}(z).$$

This relationship holds even for nonsmooth convex functions using the **subdifferential**, so

$$\text{prox}_f = (\text{Id} + \partial f)^{-1},$$

where Id denotes the identity operator (not the identity matrix) and the inverse above denotes a function inverse, not a matrix inverse.

29.9.3.7 Moreau envelope

For a convex function f and $\lambda > 0$, the **Moreau envelope** of f with parameter λ is defined [9–12, 15] [16, Sect. 12.4] by

$$f_\lambda(z) \triangleq \inf_x f(x) + \frac{1}{2\lambda} \|x - z\|_2^2. \quad (29.9.19)$$

Comparing to the proximity operator (29.9.18), we see that

$$f_\lambda(z) = f(\text{prox}_{\lambda f}(z)) + \frac{1}{2\lambda} \|\text{prox}_{\lambda f}(z) - z\|_2^2.$$

The **Moreau envelope** is also the **infimal convolution** of $f(\cdot)$ with $\frac{1}{2\lambda} \|\cdot\|_2^2$.

29.9.3.8 Minimization of convex functions over convex sets

Many of the image reconstruction problems described in this book require finding the minimizer $\hat{\mathbf{x}}$ of a convex cost function $\Psi(\mathbf{x})$ over a convex set $\mathcal{C} \subset \mathbb{R}^{n_p}$, i.e.,

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathcal{C}} \Psi(\mathbf{x}).$$

In general the minimizer of $\Psi(\mathbf{x})$ may not be unique even over \mathcal{C} , so we define the solution set

$$\mathcal{X}^{(*)} \triangleq \{\mathbf{x}^{(*)} \in \mathcal{C} : \Psi(\mathbf{x}^{(*)}) \leq \Psi(\mathbf{x}), \forall \mathbf{x} \in \mathcal{C}\}.$$

Constrained minimization problems require more work than simply equating the gradient to zero. When Ψ is a convex cost function and \mathcal{C} is a convex set, the constrained minimizers are characterized by the following result.

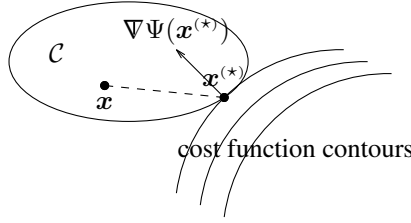
Theorem 29.9.26 [6, p. 203] *If Ψ is a differentiable convex function and \mathcal{C} is a convex set, then*

$$\mathbf{x}^{(*)} \in \mathcal{X}^{(*)} \iff \nabla \Psi(\mathbf{x}^{(*)}) (\mathbf{x} - \mathbf{x}^{(*)}) \geq 0, \quad \forall \mathbf{x} \in \mathcal{C}.$$

Proof of \Leftarrow [6, p. 200]:

Suppose $\nabla \Psi(\mathbf{x}^{(*)}) (\mathbf{x}^{(0)} - \mathbf{x}^{(*)}) < 0$ for some $\mathbf{x}^{(0)} \in \mathcal{C}$. Define $\mathbf{x}_\alpha = (1 - \alpha) \mathbf{x}^{(*)} + \alpha \mathbf{x}^{(0)} = \mathbf{x}^{(*)} + \alpha(\mathbf{x}^{(0)} - \mathbf{x}^{(*)})$ and note $\mathbf{x}_\alpha \in \mathcal{C}$ because \mathcal{C} is convex. Then by Taylor series: $\Psi(\mathbf{x}_\alpha) = \Psi(\mathbf{x}^{(*)}) + \alpha \nabla \Psi(\mathbf{x}^{(*)}) (\mathbf{x}_\alpha - \mathbf{x}^{(*)}) + o(\alpha)$ so $\Psi(\mathbf{x}_\alpha) < \Psi(\mathbf{x}^{(*)})$ for α sufficiently small, a contradiction. \square

The preceding characterization is easily understood geometrically: for $\mathbf{x} \in \mathcal{C}$, the vector $\mathbf{x} - \mathbf{x}^{(*)}$ has an angle less than $\pi/2$ away from the gradient $\nabla \Psi(\mathbf{x}^{(*)})$.



29.9.3.9 Exchanging order of minimization and maximization

In some optimization problems one would like to exchange the order of minimization and maximization as follows:

$$\min_{\mathbf{x} \in \mathcal{X}} \sup_{\mathbf{y} \in \mathcal{Y}} \Psi(\mathbf{x}, \mathbf{y}) \stackrel{?}{=} \sup_{\mathbf{y} \in \mathcal{Y}} \min_{\mathbf{x} \in \mathcal{X}} \Psi(\mathbf{x}, \mathbf{y}).$$

Sufficient conditions (including convexity in \mathbf{x} and concavity in \mathbf{y}) are given in [17, 18].

29.10 Minimizers with nonnegativity constraints (s,math,min,nonneg)

Theorem 29.10.1 *Let $f(\mathbf{x})$ be a differentiable function defined over the **nonnegative orthant** $\mathbb{R}_+^{n_p}$. Let $\mathbf{D} = \text{diag}\{d_j\}$ denote a diagonal $n_p \times n_p$ matrix with nonnegative diagonal elements. Let $\mathbf{x}^{(*)} \in \mathbb{R}_+^{n_p}$ denote a constrained minimizer of f , i.e., $f(\mathbf{x}^{(*)}) \leq f(\mathbf{x})$, $\forall \mathbf{x} \in \mathbb{R}_+^{n_p}$. Then*

$$\langle \nabla f(\mathbf{x}^{(*)}), \mathbf{D}(\mathbf{x} - \mathbf{x}^{(*)}) \rangle \geq 0, \quad \forall \mathbf{x} \in \mathbb{R}_+^{n_p}.$$

Proof:

By the **Karush-Kuhn-Tucker (KKT)** conditions,

$$\frac{\partial}{\partial x_j} f(\mathbf{x}^{(*)}) \begin{cases} = 0, & x_j^{(*)} > 0 \\ \geq 0, & x_j^{(*)} = 0, \end{cases} \quad (29.10.1)$$

so $x_j^{(*)} \frac{\partial}{\partial x_j} f(\mathbf{x}^{(*)}) = 0$ and hence $\sum_{j=1}^{n_p} d_j x_j^{(*)} \frac{\partial}{\partial x_j} f(\mathbf{x}^{(*)}) = 0$. Furthermore, $\mathbf{x} \in \mathbb{R}_+^{n_p}$ implies $x_j \geq 0$, so

$$\sum_{j=1}^{n_p} d_j x_j \frac{\partial}{\partial x_j} f(\mathbf{x}^{(*)}) = 0,$$

because $\frac{\partial}{\partial x_j} f(\mathbf{x}^{(*)}) \geq 0$. \square

29.11 Augmented Lagrangian methods (s,math,al)

This section summarizes the **augmented Lagrangian** method for constrained optimization [20, Ch. 3] [21], considering the case of complex vectors.

Let $\Psi : \mathbb{C}^{n_p} \mapsto \mathbb{R}$ denote a cost function and consider the constrained minimization problem

$$\arg \min_{\mathbf{x} \in \mathbb{C}^{n_p}} \Psi(\mathbf{x}) \quad \text{sub. to } \mathbf{B}\mathbf{x} = \mathbf{c}, \quad (29.11.1)$$

where \mathbf{B} is a $M \times n_p$ matrix (possibly complex valued) and $\mathbf{c} \in \mathbb{C}^M$.

It is easy to verify that (29.11.1) is equivalent to the following **saddle-point problem**:

$$\max_{\mathbf{q} \in \mathbb{C}^M} \arg \min_{\mathbf{x} \in \mathbb{C}^{n_p}} \Psi(\mathbf{x}) + \text{real}\{\mathbf{q}'(\mathbf{B}\mathbf{x} - \mathbf{c})\}, \quad (29.11.2)$$

because if we choose $\mathbf{q} = \alpha(\mathbf{B}\mathbf{x} - \mathbf{c})$ for $\alpha \in \mathbb{R}$ then the second term becomes $\alpha \|\mathbf{B}\mathbf{x} - \mathbf{c}\|^2$ which is unbounded as $\alpha \rightarrow \infty$ unless $\mathbf{B}\mathbf{x} = \mathbf{c}$. The vector \mathbf{q} is called a **Lagrange multiplier**.

The saddle-point problem (29.11.2) is also equivalent to the following problem:

$$\max_{\mathbf{q} \in \mathbb{C}^M} \arg \min_{\mathbf{x} \in \mathbb{C}^{n_p}} L_\mu(\mathbf{x}, \mathbf{q}), \quad (29.11.3)$$

where the **augmented Lagrangian** is defined by

$$L_\mu(\mathbf{x}, \mathbf{q}) \triangleq \Psi(\mathbf{x}) + \text{real}\{\mathbf{q}'(\mathbf{B}\mathbf{x} - \mathbf{c})\} + \frac{\mu}{2} \|\mathbf{B}\mathbf{x} - \mathbf{c}\|^2, \quad (29.11.4)$$

for any real $\mu > 0$. By **completing the square**, one can rewrite the **augmented Lagrangian** as²

$$L_\mu(\mathbf{x}, \mathbf{q}) = \Psi(\mathbf{x}) + \frac{\mu}{2} \left\| \mathbf{B}\mathbf{x} - \mathbf{c} + \frac{1}{\mu} \mathbf{q} \right\|^2 - \frac{\mu}{2} \left\| \frac{1}{\mu} \mathbf{q} \right\|^2. \quad (29.11.5)$$

Identifying $\boldsymbol{\eta} = \frac{-1}{\mu} \mathbf{q}$ we rewrite (29.11.3) as

$$\max_{\boldsymbol{\eta} \in \mathbb{C}^M} \arg \min_{\mathbf{x} \in \mathbb{C}^{n_p}} \tilde{L}_\mu(\mathbf{x}, \boldsymbol{\eta}) \quad \text{where } \tilde{L}_\mu(\mathbf{x}, \boldsymbol{\eta}) \triangleq L_\mu(\mathbf{x}, -\mu\boldsymbol{\eta}) = \Psi(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{B}\mathbf{x} - \mathbf{c} - \boldsymbol{\eta}\|^2 - \frac{\mu}{2} \|\boldsymbol{\eta}\|^2.$$

One way to solve this saddle-point problem is to alternate between updating \mathbf{x} and updating $\boldsymbol{\eta}$. We update \mathbf{x} using

$$\mathbf{x}^{(n+1)} = \arg \min_{\mathbf{x} \in \mathbb{C}^{n_p}} \tilde{L}_\mu(\mathbf{x}, \boldsymbol{\eta}^{(n)}) = \arg \min_{\mathbf{x} \in \mathbb{C}^{n_p}} \Psi(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{B}\mathbf{x} - \mathbf{c} - \boldsymbol{\eta}^{(n)}\|^2.$$

At first it might seem natural to update $\boldsymbol{\eta}$ using

$$\boldsymbol{\eta}^{(n+1)} = \arg \max_{\boldsymbol{\eta} \in \mathbb{C}^M} \tilde{L}_\mu(\mathbf{x}^{(n+1)}, \boldsymbol{\eta}).$$

However, that would lead to $\boldsymbol{\eta}^{(n+1)}$ with infinite norm when $\mathbf{B}\mathbf{x}^{(n)} \neq \mathbf{c}$. Instead we apply a single **gradient ascent** update (using (29.2.16) in the complex case) for the Lagrange multiplier:

$$\boldsymbol{\eta}^{(n+1)} \triangleq \boldsymbol{\eta}^{(n)} + \alpha_n \nabla_{\boldsymbol{\eta}} \tilde{L}_\mu(\mathbf{x}^{(n+1)}, \boldsymbol{\eta}^{(n)}) = \boldsymbol{\eta}^{(n)} - \alpha_n \mu (\mathbf{B}\mathbf{x} - \mathbf{c}).$$

The parameter α_n is a step size that must be chosen appropriately to ensure convergence [20, Ch. 3]. Other methods for updating the Lagrange multiplier $\boldsymbol{\eta}$ have been proposed [23, 24].

29.12 Special functions (s,math,floor)

The **floor** and **ceiling** operations often are useful for imaging geometry calculations. They are defined as follows:

$$\begin{aligned} \lfloor x \rfloor &\triangleq \max \{n \in \mathbb{Z} : n \leq x\} \\ \lceil x \rceil &\triangleq \min \{n \in \mathbb{Z} : n \geq x\}. \end{aligned}$$

² Ramani et al. [22, after (13)] imply that the last term is an irrelevant constant. It is independent of \mathbf{x} , so irrelevant for updating \mathbf{x} , but it is essential for the update of $\boldsymbol{\eta}$.

Useful properties of $\lceil \cdot \rceil$ and $\lfloor \cdot \rfloor$ include the following:

$$\begin{aligned} x &\leq \lceil x \rceil < 1 + x \\ x - 1 &< \lfloor x \rfloor \leq x \\ \lceil x \rceil &\leq \lfloor x + 1 \rfloor = \lfloor x \rfloor + 1 \end{aligned} \tag{29.12.1}$$

$$\begin{aligned} \lfloor -x \rfloor &= -\lceil x \rceil \\ \lceil -x \rceil &= -\lfloor x \rfloor \\ \lceil x + b \rceil - \lfloor x + a \rfloor &\leq 1 + \lceil b - a \rceil. \end{aligned} \tag{29.12.2}$$

29.13 Convergence rates of iterations (s,math,rate)

One can evaluate the **convergence rate** of iterative methods using a variety of metrics.

- When a cost function Ψ has a minimum value $\Psi_* = \min_{\mathbf{z}} \Psi(\mathbf{z})$, we can examine how quickly $\Psi(\mathbf{x}^{(n)}) - \Psi_*$ approaches zero.
- If a cost function Ψ is differentiable, we can examine how quickly a norm of its gradient $\|\nabla \Psi(\mathbf{x}^{(n)})\|$ approaches zero.
- If a cost function Ψ has a unique minimizer $\mathbf{x}^{(*)}$, we can examine how quickly $\|\mathbf{x}^{(n)} - \mathbf{x}^{(*)}\|$ approaches zero.

The rest of this section focuses on the latter case because it is of the most interest.

We say a sequence $\{\mathbf{x}^{(n)}\}$ **converges** to a limit $\mathbf{x}^{(*)}$ with respect to a **norm** $\|\cdot\|$ if for all $\varepsilon > 0$ there exists a number N_ε such that

$$\|\mathbf{x}^{(n)} - \mathbf{x}^{(*)}\| < \varepsilon, \quad \forall n \geq N_\varepsilon.$$

There are various ways to quantify the **asymptotic convergence rate** of convergent sequences, and these are important for understanding the limiting behavior of optimization methods.

The **quotient convergence factor** of a sequence $\{\mathbf{x}^{(n)}\}$ converging to a limit $\mathbf{x}^{(*)}$ is defined by [25, p. 281] [26]

$$Q_1(\{\mathbf{x}^{(n)}\}) \triangleq \limsup_{n \rightarrow \infty} \frac{\|\mathbf{x}^{(n+1)} - \mathbf{x}^{(*)}\|}{\|\mathbf{x}^{(n)} - \mathbf{x}^{(*)}\|}. \tag{29.13.1}$$

If this **limit superior** lies in the interval $(0, 1)$ then we say that the sequence $\{\mathbf{x}^{(n)}\}$ **converges linearly** to $\mathbf{x}^{(*)}$.

Another measure of convergence rate of a sequence $\{\mathbf{x}^{(n)}\}$ converging to $\mathbf{x}^{(*)}$ is the **root convergence factor** defined by [25, p. 288]

$$R_1(\{\mathbf{x}^{(n)}\}) \triangleq \limsup_{n \rightarrow \infty} \|\mathbf{x}^{(n)} - \mathbf{x}^{(*)}\|^{1/n}. \tag{29.13.2}$$

One can show that $0 \leq R_1 \leq 1$. Again, if this limit superior lies in the interval $(0, 1)$ then we say that the sequence **converges linearly** to $\mathbf{x}^{(*)}$.

Unlike the quotient convergence factor Q_1 , the root convergence factor R_1 is independent of the norm [25, p. 288]. Furthermore, $R_1 \leq Q_1$ for any norm [25, p. 296].

Often we are more interested in defining the convergence rate of an **iterative process** $\mathbf{x}^{(n+1)} = \mathcal{M}(\mathbf{x}^{(n)})$, where $\mathcal{M} : \mathbb{R}^{n_p} \rightarrow \mathbb{R}^{n_p}$, rather than that of a specific sequence generated by the process [27]. After all, if we happen to initialize the algorithm with $\mathbf{x}^{(0)} = \mathbf{x}^{(*)}$ where $\mathbf{x}^{(*)} = \mathcal{M}(\mathbf{x}^{(*)})$ is a **fixed point** of \mathcal{M} then we would converge immediately. So clearly the convergence rate of a particular sequence can depend on the initial condition. To avoid this dependence we examine the **worst case** over all possible initializers (or, more precisely, over all initializers that lead to convergence to $\mathbf{x}^{(*)}$):

$$Q_1(\mathcal{M}) \triangleq \sup_{\mathbf{x}^{(0)}} \{Q_1(\{\mathbf{x}^{(n)}\}) : \mathbf{x}^{(n+1)} = \mathcal{M}(\mathbf{x}^{(n)}) \rightarrow \mathbf{x}^{(*)}\}. \tag{29.13.3}$$

We call this the **quotient convergence factor** of the iterative process \mathcal{M} . Likewise we can define a **root convergence factor** of \mathcal{M} by

$$R_1(\mathcal{M}) \triangleq \sup_{\mathbf{x}^{(0)}} \{R_1(\{\mathbf{x}^{(n)}\}) : \mathbf{x}^{(n+1)} = \mathcal{M}(\mathbf{x}^{(n)}) \rightarrow \mathbf{x}^{(*)}\}. \tag{29.13.4}$$

Ostrowski's theorem [25, p. 300] states that if \mathcal{M} is continuous and differentiable on an open set containing $\mathbf{x}^{(*)}$ and

$$\rho(\nabla \mathcal{M}(\mathbf{x}^{(*)})) < 1,$$

where $\nabla \mathcal{M}(\mathbf{x})$ is the $n_p \times n_p$ matrix with elements

$$[\nabla \mathcal{M}(\mathbf{x})]_{k,j} = \frac{\partial}{\partial x_j} [\mathcal{M}(\mathbf{x})]_k,$$

and where $\rho(\cdot)$ in (27.1.2) denotes the **spectral radius** (largest eigenvalue magnitude) of a square matrix, then the root convergence factor of \mathcal{M} is

$$R_1(\mathcal{M}) = \rho(\nabla \mathcal{M}(\mathbf{x}^{(*)})). \quad (29.13.5)$$

As a sketch of why this equality holds, for $\mathbf{x}^{(n)} \approx \mathbf{x}^{(*)}$:

$$\mathbf{x}^{(n+1)} \approx \mathcal{M}(\mathbf{x}^{(*)}) + \nabla \mathcal{M}(\mathbf{x}^{(*)})(\mathbf{x}^{(n)} - \mathbf{x}^{(*)}) = \mathbf{x}^{(*)} + \nabla \mathcal{M}(\mathbf{x}^{(*)})(\mathbf{x}^{(n)} - \mathbf{x}^{(*)})$$

so

$$\mathbf{x}^{(n+1)} - \mathbf{x}^{(*)} \approx \nabla \mathcal{M}(\mathbf{x}^{(*)})(\mathbf{x}^{(n)} - \mathbf{x}^{(*)}).$$

Thus the asymptotic rate of convergence of $\mathbf{x}^{(n)}$ to $\mathbf{x}^{(*)}$ is governed by the eigenvalues of $\nabla \mathcal{M}(\mathbf{x}^{(*)})$.

Example 29.13.1 Consider a sequence generated by the affine recursion

$$\mathbf{x}^{(n+1)} = \mathcal{M}(\mathbf{x}^{(n)}) \triangleq \mathbf{B}\mathbf{x}^{(n)} + \mathbf{u} \quad (29.13.6)$$

where $\mathbf{I} - \mathbf{B}$ is an invertible matrix and we define $\mathbf{x}^{(*)} = [\mathbf{I} - \mathbf{B}]^{-1} \mathbf{u}$ so that $\mathbf{x}^{(n+1)} - \mathbf{x}^{(*)} = \mathbf{B}(\mathbf{x}^{(n)} - \mathbf{x}^{(*)})$ and hence $\mathbf{x}^{(n)} - \mathbf{x}^{(*)} = \mathbf{B}^n(\mathbf{x}^{(0)} - \mathbf{x}^{(*)})$. If $\rho(\mathbf{B}) < 1$ then it is easy to show [25, p. 303] using the **Jordan form** of \mathbf{B} that

$$R_1(\mathcal{M}) = \rho(\mathbf{B}).$$

As a more concrete example, the **PGD** iteration $\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \mathbf{P} \nabla \Psi(\mathbf{x}^{(n)})$ for the quadratic cost function $\Psi(\mathbf{x}) = \frac{1}{2} \mathbf{x}' \mathbf{H} \mathbf{x} - \mathbf{b}' \mathbf{x}$ has the form (29.13.6) with $\mathbf{B} = \mathbf{I} - \mathbf{P} \mathbf{H}$ and $\mathbf{u} = \mathbf{P} \mathbf{b}$. So the root convergence factor for PGD in this case is $R_1 = \rho(\mathbf{I} - \mathbf{P} \mathbf{H})$.

29.14 Problems (s,math,prob)

Problem 29.1 Prove affine scaling property (29.2.5) of the Fourier transform.

(Need typed.)

Problem 29.2 Prove the Laplace transform pair for $\cos(bt) e^{-at} \text{step}(t)$ in §29.4.1.

Problem 29.3 Prove the Laplace transform pair for $\sin(bt) e^{-at} \text{step}(t)$ in §29.4.1.

Problem 29.4 Prove the Laplace transform pair for $\cos(bt) e^{-a|t|}$ in §29.4.1

Problem 29.5 Prove the Laplace transform pair for $\sin(b|t|) e^{-a|t|}$ in §29.4.1

Problem 29.6 Prove that if f is convex on \mathbb{R}^n then $g(\mathbf{z}) = f(\mathbf{M}\mathbf{z} + \mathbf{b})$ is convex.

(Need typed.)

Problem 29.7 Prove that $f(x) = |x|^p$ is strictly convex for $p > 1$, for Example 29.9.10.

Problem 29.8 Prove Corollary 29.9.21, i.e., (29.9.13).

Problem 29.9 Prove Corollary 29.9.20, i.e., for a twice differentiable, convex cost function Ψ , the Lipschitz condition (29.9.12) is equivalent to $\nabla^2 \Psi(\mathbf{x}) \preceq \mathbf{S} \mathbf{S}'$. Hint: use Theorem 27.5.3.

29.15 Bibliography

- [1] D. G. Luenberger. *Optimization by vector space methods*. New York: Wiley, 1969. URL: <http://books.google.com/books?id=1ZU0CAH4RccC> (cit. on p. 29.9).
- [2] S. Boyd and L. Vandenberghe. *Convex optimization*. UK: Cambridge, 2004. URL: <http://web.stanford.edu/~boyd/cvxbook.html> (cit. on pp. 29.10, 29.11, 29.12).

- [3] A. Wiesel. “Geodesic convexity and covariance estimation.” In: *IEEE Trans. Sig. Proc.* 60.12 (Dec. 2012), 6182–9. DOI: [10.1109/TSP.2012.2218241](https://doi.org/10.1109/TSP.2012.2218241) (cit. on p. 29.10).
- [4] D. P. Bertsekas. *Constrained optimization and Lagrange multiplier methods*. New York: Academic-Press, 1982. URL: <http://www.athenasc.com/> (cit. on p. 29.10).
- [5] R. A. Horn and C. R. Johnson. *Matrix analysis*. Cambridge: Cambridge Univ. Press, 1985 (cit. on p. 29.11).
- [6] B. T. Polyak. *Introduction to optimization*. New York: Optimization Software Inc, 1987 (cit. on pp. 29.12, 29.13, 29.15).
- [7] Y. Nesterov. *Introductory lectures on convex optimization: A basic course*. Kluwer, 2004. DOI: [10.1007/978-1-4419-8853-9](https://doi.org/10.1007/978-1-4419-8853-9) (cit. on pp. 29.12, 29.13).
- [8] P. Giselsson and S. Boyd. “Monotonicity and restart in fast gradient methods.” In: *Proc. Conf. Decision and Control*. 2014, 5058–63. DOI: [10.1109/CDC.2014.7040179](https://doi.org/10.1109/CDC.2014.7040179) (cit. on p. 29.13).
- [9] J.-J. Moreau. “Fonctions convexes duales et points proximaux dans un espace hilbertien.” In: *C. R. Acad. Sci. Paris* 255 (1962), 2897–289 (cit. on p. 29.14).
- [10] J. J. Moreau. “Proximité et dualité dans un espace hilbertien.” In: *Bulletin de la Société Mathématique de France* 93 (1965), 273–99. URL: http://www.numdam.org/item?id=BSMF_1965__93__273_0 (cit. on p. 29.14).
- [11] P. Combettes and V. Wajs. “Signal recovery by proximal forward-backward splitting.” In: *siam-jmms* 4.4 (2005), 1168–200. DOI: [10.1137/050626090](https://doi.org/10.1137/050626090) (cit. on p. 29.14).
- [12] N. Parikh and S. Boyd. “Proximal algorithms.” In: *Found. Trends in Optimization* 1.3 (2013), 123–231. DOI: [10.1561/2400000003](https://doi.org/10.1561/2400000003) (cit. on p. 29.14).
- [13] A. Beck and M. Teboulle. “A fast iterative shrinkage-thresholding algorithm for linear inverse problems.” In: *SIAM J. Imaging Sci.* 2.1 (2009), 183–202. DOI: [10.1137/080716542](https://doi.org/10.1137/080716542) (cit. on p. 29.14).
- [14] A. Beck and M. Teboulle. “Fast gradient-based algorithms for constrained total variation image denoising and deblurring problems.” In: *IEEE Trans. Im. Proc.* 18.11 (Nov. 2009), 2419–34. DOI: [10.1109/TIP.2009.2028250](https://doi.org/10.1109/TIP.2009.2028250) (cit. on p. 29.14).
- [15] A. Beck and M. Teboulle. “Smoothing and first order methods: A unified framework.” In: *SIAM J. Optim.* 22.2 (2012), 557–80. DOI: [10.1137/100818327](https://doi.org/10.1137/100818327) (cit. on p. 29.14).
- [16] H. H. Bauschke and P. L. Combettes. *Convex analysis and monotone operator theory in Hilbert spaces*. Springer, 2011. DOI: [10.1007/978-1-4419-9467-7](https://doi.org/10.1007/978-1-4419-9467-7) (cit. on p. 29.14).
- [17] M. Sion. “On general minimax theorems.” In: *Pacific J. Math.* 8.1 (1958), 171–6. DOI: [10.2140/pjm.1958.8.171](https://doi.org/10.2140/pjm.1958.8.171) (cit. on p. 29.15).
- [18] H. Komiyama. “Elementary proof for Sion’s minimax theorem.” In: *Kodai Mathematical J.* 11.1 (1988), 5–7. DOI: [10.2996/kmj/1138038812](https://doi.org/10.2996/kmj/1138038812) (cit. on p. 29.15).
- [19] R. Mourya et al. “Augmented Lagrangian without alternating directions: practical algorithms for inverse problems in imaging.” In: *Proc. IEEE Intl. Conf. on Image Processing*. 2015, 1205–9. DOI: [10.1109/ICIP.2015.7350991](https://doi.org/10.1109/ICIP.2015.7350991).
- [20] R. Glowinski and P. L. Tallec. *Augmented Lagrangian and operator-splitting methods in nonlinear mechanics*. Soc. Indust. Appl. Math., 1989. DOI: [10.1137/1.9781611970838.ch3](https://doi.org/10.1137/1.9781611970838.ch3) (cit. on p. 29.16).
- [21] E. G. Birgin and J. M. Martinez. *Practical augmented Lagrangian methods for constrained optimization*. ISBN 978-1-611973-35-8. Soc. Indust. Appl. Math., 2014. URL: <http://www.ec-securehost.com/SIAM/FA10.html> (cit. on p. 29.16).
- [22] S. Ramani and J. A. Fessler. “Parallel MR image reconstruction using augmented Lagrangian methods.” In: *IEEE Trans. Med. Imag.* 30.3 (Mar. 2011), 694–706. DOI: [10.1109/TMI.2010.2093536](https://doi.org/10.1109/TMI.2010.2093536) (cit. on p. 29.16).
- [23] J. Eckstein. “Parallel alternating direction multiplier decomposition of convex programs.” In: *J. Optim. Theory Appl.* 80.1 (Jan. 1994), 39–62. DOI: [10.1007/BF02196592](https://doi.org/10.1007/BF02196592) (cit. on p. 29.16).
- [24] S. Boyd et al. “Distributed optimization and statistical learning via the alternating direction method of multipliers.” In: *Found. & Trends in Machine Learning* 3.1 (2010), 1–122. DOI: [10.1561/2200000016](https://doi.org/10.1561/2200000016) (cit. on p. 29.16).

ortega:70

- [25] J. M. Ortega and W. C. Rheinboldt. *Iterative solution of nonlinear equations in several variables*. New York: Academic, 1970. DOI: [10.1137/1.9780898719468](https://doi.org/10.1137/1.9780898719468) (cit. on pp. [29.17](#), [29.18](#)).

potra:89:oqo

- [26] F. A. Potra. “On Q-order and R-order of convergence.” In: *J. Optim. Theory Appl.* 63.3 (1989), 415–31. DOI: [10.1007/BF00939805](https://doi.org/10.1007/BF00939805) (cit. on p. [29.17](#)).

dubeau:14:fpa

- [27] F. Dubeau and C. Gnan. “Fixed point and Newton’s methods for solving a nonlinear equation: from linear to high-order convergence.” In: *SIAM Review* 56.4 (Dec. 2014), 691–708. DOI: [10.1137/130934799](https://doi.org/10.1137/130934799) (cit. on p. [29.17](#)).