



Linear Inverse Problems (1/2)

Mathematical Foundations of Signal Processing

Dr. Matthieu Simeoni

November 9, 2020

EPFL

Table of contents

1 Introduction

- Problem Statement
- Common Sampling Functionals
- Common Noise Distributions
- Examples of Inverse Problems in Natural Sciences

2 Solving Inverse Problems

- Ill-Posedness
- Regularisation
- Existence and Unicity of Solutions
- Common Cost Functionals
- Common Regularisation Functionals and Representer Theorems
- Bayesian Interpretation

Introduction

Most real-life approximation problems can be formulated as **inverse problems**:

Inverse Problem

Consider an **unknown** signal $f \in \mathcal{L}^2(\mathbb{R}^d)$ and assume that the latter is *probed* by some **sensing device**, resulting in a **data vector** $\mathbf{y} = [y_1, \dots, y_L] \in \mathbb{R}^L$ of L **measurements**. Recovering f from the data vector \mathbf{y} is called an **inverse problem**.

We make the following assumptions:

1. To account for sensing *inaccuracies*, the data vector \mathbf{y} is assumed to be the outcome of a **random vector** $\mathbf{Y} = [Y_1, \dots, Y_L] : \Omega \rightarrow \mathbb{R}^L$, fluctuating according to some **noise distribution**. The entries of $\mathbb{E}[\mathbf{Y}] = \tilde{\mathbf{y}}$ are called the **ideal measurements** –these are the measurements that would be obtained in the absence of noise.
2. The measurements are assumed **unbiased** and **linear**, i.e. $\mathbb{E}[\mathbf{Y}] = \Phi^* f = [\langle f, \varphi_1 \rangle, \dots, \langle f, \varphi_L \rangle]$, for some **sampling functionals** $\{\varphi_1, \dots, \varphi_L\} \subset \mathcal{L}^2(\mathbb{R}^d)$, modelling the **acquisition system**.

Common Sampling Functionals

Common Sampling Functionals

- **Spatial Sampling:**

$$\tilde{y}_i = f(x_i) = \int_{\mathbb{R}^d} f(x) \delta(x - x_i) dx \quad \rightarrow \quad \varphi_i(x) = \delta(x - x_i), \quad x_i \in \mathbb{R}^d.$$

- **Fourier Sampling:**

$$\tilde{y}_{i1} = \int_{\mathbb{R}^d} f(x) \cos(\langle x, \omega_i \rangle) dx \quad \rightarrow \quad \varphi_{i1}(x) = \cos(\langle x, \omega_i \rangle), \quad \omega_i \in \mathbb{R}^d.$$

$$\tilde{y}_{i2} = \int_{\mathbb{R}^d} f(x) \sin(\langle x, \omega_i \rangle) dx \quad \rightarrow \quad \varphi_{i2}(x) = \sin(\langle x, \omega_i \rangle), \quad \omega_i \in \mathbb{R}^d.$$

- **Radon Sampling:**

$$\tilde{y}_i = \check{f}(p_i, \xi_i) = \int_{\mathbb{R}^d} f(x) \delta(p_i - \langle x, \xi_i \rangle) dx \quad \rightarrow \quad \varphi_i(x) = \delta(p_i - \langle x, \xi_i \rangle), \quad p_i > 0, \xi_i \in \mathbb{S}^{N-1}.$$

- **Filtering:**

$$\tilde{y}_i = \{f * \varphi\}(x_i) = \int_{\mathbb{R}^d} \varphi(x_i - x) f(x) dx \quad \rightarrow \quad \varphi_i(x) = \varphi(x_i - x), \quad x_i \in \mathbb{R}^d, \varphi: \mathbb{R}^d \rightarrow \mathbb{R}.$$

- **Mean-Pooling:**

$$\tilde{y}_i = \frac{1}{|\Omega_i|} \int_{\Omega_i} f(x) dx \quad \rightarrow \quad \varphi_i(x) = \frac{1}{|\Omega_i|} \chi_{\Omega_i}(x) := \begin{cases} |\Omega_i|^{-1} & \text{if } x \in \Omega_i \\ 0 & \text{otherwise} \end{cases}, \quad \Omega_i \subset \mathbb{R}^d.$$

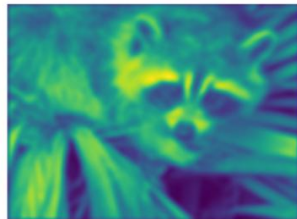
Example: Deblurring



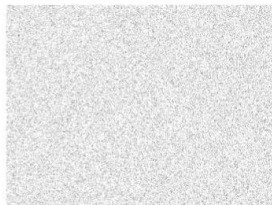
*



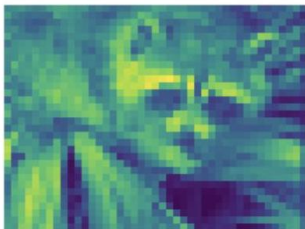
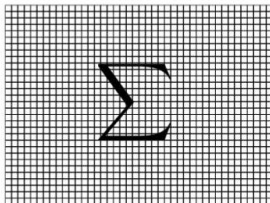
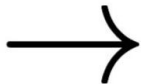
=



Example: Inpainting



Example: Unpooling

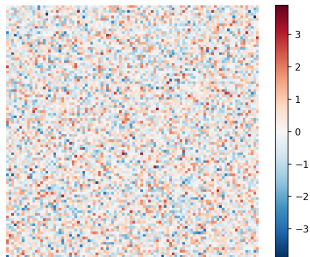
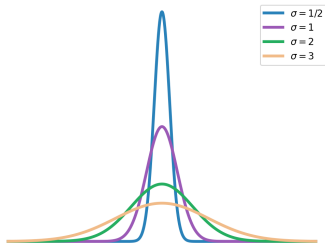


Gaussian White Noise

Assume that sensor inaccuracies are **independent** and result from the *sum of many independent perturbations*. Then, from the **central limit theorem**, sensor inaccuracies can be modelled as **independent realisations** of an **additive Gaussian white noise**:

$$Y_i = \tilde{y}_i + N_i, \quad \text{where } N_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2), \quad p_N(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{x^2}{2\sigma^2}\right), \quad x \in \mathbb{R},$$

where p_N is the noise **probability density function**. Notice that we have indeed $\mathbb{E}[Y_i] = \tilde{y}_i$ for each $i = 1, \dots, L$.

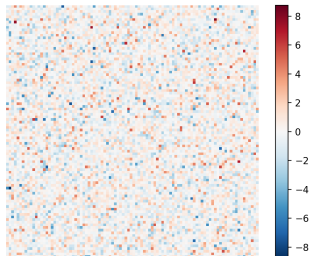
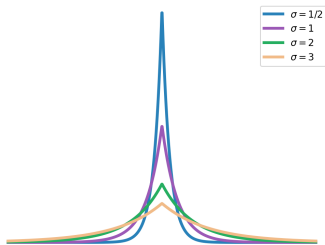


Laplacian/Salt-and-pepper White Noise

Assume that sensor inaccuracies are **independent** and present **strong outliers** (for example due to *malfunctioning* sensors). Then, sensor inaccuracies can be modelled as **independent realisations** of an *additive Laplacian white noise*, also called **salt-and-pepper noise**:

$$Y_i = \tilde{y}_i + N_i, \quad \text{where } N_i \stackrel{\text{i.i.d.}}{\sim} \text{Laplace}(0, \sigma), \quad p_N(x) = \frac{1}{2\sigma} \exp\left(-\frac{|x|}{\sigma}\right), \quad x \in \mathbb{R},$$

where p_N is the noise **probability density function**. Notice that we have indeed $\mathbb{E}[Y_i] = \tilde{y}_i$ for each $i = 1, \dots, L$.

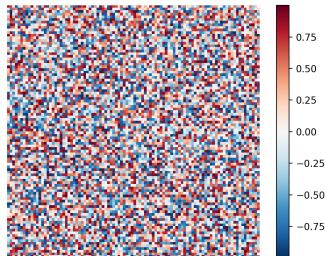
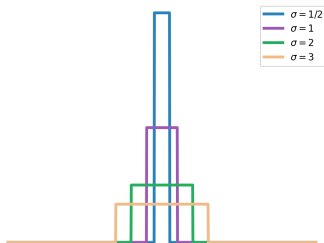


Uniform/Quantisation White Noise

Assume that sensor inaccuracies are **independent** and primarily caused by **quantisation artefacts** –i.e. round-off errors incurred by storing digits with finite precision. Then, sensor inaccuracies can be modelled as **independent realisations** of an *additive uniform white noise*, also called **quantisation noise**:

$$Y_i = \tilde{y}_i + N_i, \quad \text{where } N_i \stackrel{\text{i.i.d.}}{\sim} U\left(-\frac{\sigma}{2}, \frac{\sigma}{2}\right), \quad p_N(x) = \begin{cases} 1/\sigma & \text{if } x \in [-\sigma/2, \sigma/2] \\ 0 & \text{if } x \notin [-\sigma/2, \sigma/2]. \end{cases}$$

where p_N is the noise **probability density function**. Notice that we have indeed $\mathbb{E}[Y_i] = \tilde{y}_i$ for each $i = 1, \dots, L$.

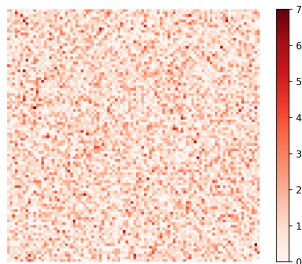
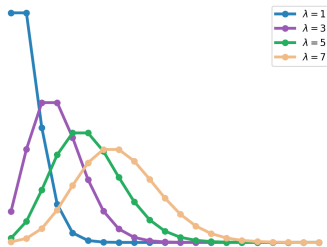


Poisson/Shot Noise

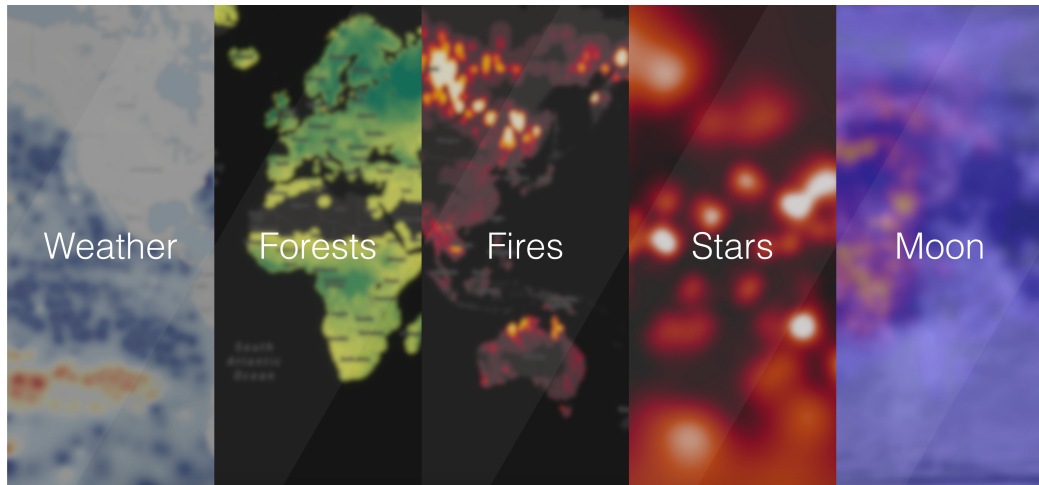
Assume that the measurements are **independent** and originate from a **counting process** –i.e. $Y: \Omega \rightarrow \mathbb{N}^L$. Then, sensor inaccuracies can be modelled as **independent realisations** of a **non additive Poisson noise**, also called **shot noise**:

$$Y_i \stackrel{\text{ind}}{\sim} \text{Poisson}(\tilde{y}_i), \quad p_{Y_i}(k) = \frac{\tilde{y}_i^k e^{-\tilde{y}_i}}{k!}, \quad \forall k \in \mathbb{N},$$

where p_{Y_i} is the **probability density function** for the i th measurement. Using properties from the Poisson distribution, we can indeed show that $\mathbb{E}[Y_i] = \tilde{y}_i$ for each $i = 1, \dots, L$.



Real-Life Examples: Meteorology, Forestry, Astronomy...



<https://matthieumeo.github.io/>

Pixelisation

Since the number of measurements is **finite**, it is reasonable to constrain the signal f to be **finite-dimensional**:¹

$$f = \sum_{n=1}^N \alpha_n \psi_n = \Psi \boldsymbol{\alpha}, \quad \boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_N] \in \mathbb{R}^N \quad (1)$$

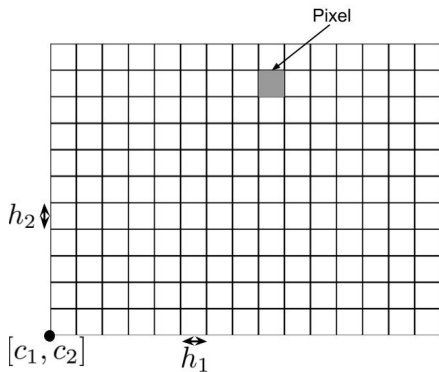
for some suitable basis functions $\{\psi_n, n = 1, \dots, N\} \subset \mathcal{L}^2(\mathbb{R}^d)$. Typically, the basis functions are chosen as **indicator functions** of **regular rectangular tiles** of \mathbb{R}^d called **pixels**. For example:

$$\psi_n(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in [c_1 + (n-1)h_1, c_1 + nh_1] \times \dots \times [c_d + (n-1)h_d, c_d + nh_d], \\ 0 & \text{otherwise,} \end{cases}$$

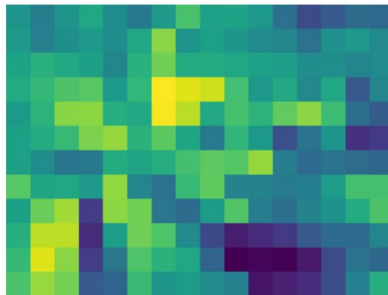
where $\mathbf{c} = [c_1, \dots, c_d]$ are the coordinates of the lower-left corner of the first pixel, and $\{h_1, \dots, h_d\}$ are the sizes of the pixels across each dimension. The parametric signal f in (1) is then a **piecewise constant signal** than can be **stored/manipulated/displayed efficiently** via **multi-dimensional array** (hence the popularity of pixel-based discretisation schemes).

¹Infinite-dimensional signals may indeed have an infinite number of degrees of freedom, which cannot hope to estimate from a finite number of measurements only.

Pixelisation



$$f(x, y) = \sum_{n=1}^N \alpha_n \psi_n(x, y)$$



Discrete Inverse Problems

Assuming the **parametric model (1)** induces a **discrete inverse problem**:

Find $\alpha \in \mathbb{R}^N$ from the noisy measurements $\mathbf{y} \leftarrow Y$ where $\mathbb{E}[Y] = \Phi^* \Psi \alpha = G\alpha$.

The operator $G: \mathbb{R}^N \rightarrow \mathbb{R}^L$ is a rectangular matrix given by:²

$$\mathbb{R}^{L \times N} \ni G = \begin{bmatrix} \langle \psi_1, \varphi_1 \rangle & \cdots & \langle \psi_N, \varphi_1 \rangle \\ \vdots & \ddots & \vdots \\ \langle \psi_1, \varphi_L \rangle & \cdots & \langle \psi_N, \varphi_L \rangle \end{bmatrix} = \begin{bmatrix} \int_{\Omega_1} \varphi_1(\mathbf{x}) d\mathbf{x} & \cdots & \int_{\Omega_N} \varphi_1(\mathbf{x}) d\mathbf{x} \\ \vdots & \ddots & \vdots \\ \int_{\Omega_1} \varphi_L(\mathbf{x}) d\mathbf{x} & \cdots & \int_{\Omega_N} \varphi_L(\mathbf{x}) d\mathbf{x} \end{bmatrix} \\ \simeq \eta \begin{bmatrix} \varphi_1(\xi_1) & \cdots & \varphi_1(\xi_N) \\ \vdots & \ddots & \vdots \\ \varphi_L(\xi_1) & \cdots & \varphi_L(\xi_N) \end{bmatrix},$$

where $\eta = \prod_{k=1}^d h_k$, and $\{\Omega_n\}_n \in \mathcal{P}(\mathbb{R}^d)$ and $\{\xi_n\}_n \subset \mathbb{R}^d$ are the *supports* and *centroids* of each pixel, respectively.

²The last approximate equality results from the **midpoint rule**.

Inverse Problems are Ill-Posed

To solve the inverse problem one can approximate the mean $\mathbb{E}[Y]$ by its *one-sample empirical estimate* y and solve the linear problem:

$$y = G\alpha. \quad (2)$$

Unfortunately, (2) is in general *ill-posed*:

1. **There may exist no solutions to (2).** If $N \ll L$ indeed (or more generally if G is *not surjective*), $\mathcal{R}(G) \subsetneq \mathbb{R}^N$. Therefore the noisy data vector y is not guaranteed to belong to $\mathcal{R}(G)$.
2. **There may exist more than one solution to (2).** If $N \gg L$ indeed (or more generally if G is *not injective*), $\mathcal{N}(G) \neq \{0\}$. Therefore, if α^\star is a solution to (2), then $\alpha^\star + \beta$ is also a solution $\forall \beta \in \mathcal{N}(G)$:

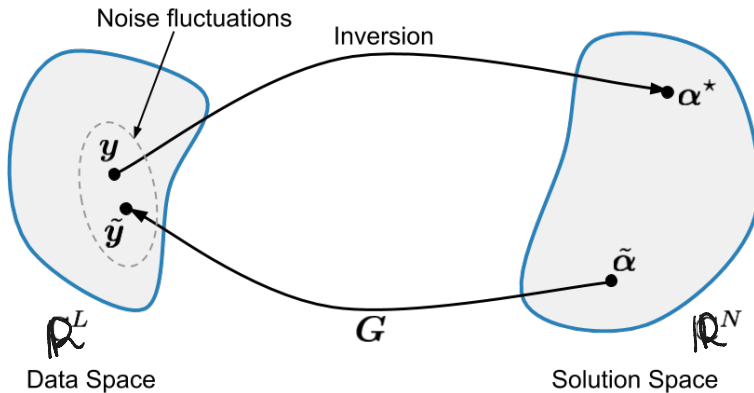
$$G(\alpha^\star + \beta) = G\alpha^\star + G\beta = G\alpha^\star = y.$$

3. **Solutions to (2) may be numerically unstable.** If G is *surjective* for example, then $G^\dagger = G^T(GG^T)^{-1}$ is a *right-inverse* of G and $\alpha^\star(y) = G^\dagger y = G^T(GG^T)^{-1}y$ is a solution to (2). We have then

$$\|\alpha^\star(y)\|_2 \leq \|G\|_2 \|(G^T G)^{-1}\|_2 \|y\|_2 = \underbrace{\frac{\sqrt{\lambda_{\max}(G^T G)}}{\lambda_{\min}(G^T G)}}_{\text{Can be very large!}} \|y\|_2, \quad \forall y \in \mathbb{R}^L.$$

The reconstruction linear map $y \mapsto \alpha^\star(y)$ can hence be *virtually unbounded* making it *unstable*.

Inverse Problems are Unstable



**Small perturbations on the data
affect greatly the solution!**

Regularising Inverse Problems

The linear system (2) is not only **ill-posed** but also **non sensible**: matching exactly the measurements is not desirable since the latter are in practice **corrupted** by instrumental noise.

A more sensible approach consists instead in solving the inverse problem by means of a **penalised optimisation problem**, confronting the physical evidence to the analyst's a priori beliefs about the solution (e.g. **smoothness**, **sparsity**) via a **data-fidelity** and **regularisation** term, respectively:

$$\min_{\alpha \in \mathbb{R}^N} F(\mathbf{y}, \mathbf{G}\alpha) + \lambda \mathcal{R}(\alpha). \quad (3)$$

The various quantities involved in (3) can be interpreted as follows:

- $F: \mathbb{R}^L \times \mathbb{R}^L \rightarrow \mathbb{R}_+ \cup \{+\infty\}$ is a **cost/data-fidelity functional**, measuring the discrepancy between the **observed** and **predicted** measurements \mathbf{y} and $\mathbf{G}\alpha$ respectively.
- $\mathcal{R}: \mathbb{R}^N \rightarrow \mathbb{R}_+ \cup \{+\infty\}$ is a **regularisation/penalty functional** favouring **simple** and **well-behaved** solutions (typically with a finite number of degrees of freedom).
- $\lambda > 0$ is a **regularisation/penalty parameter** which controls the **amount of regularisation** by putting the regularisation functional and the cost functional on a similar scale.

Existence of Solutions

Theorem: (Existence of Solutions to (3))

Consider the following set of assumptions:

1. For all $\mathbf{y} \in \mathbb{R}^L$, the univariate **cost trace functionals**

$$F(\mathbf{y}, \cdot) : \begin{cases} \mathbb{R}^L \rightarrow \mathbb{R}_+ \cup \{+\infty\} \\ \mathbf{z} \mapsto F(\mathbf{y}, \mathbf{z}) \end{cases}$$

and the **regularisation functional** $\mathcal{R} : \mathbb{R}^N \rightarrow \mathbb{R}_+ \cup \{+\infty\}$ are **proper**, **convex** and **lower semi-continuous** (see Slide 20 for a definition).

2. The objective functional of (3) is **coercive**, i.e. $\lim_{\|\boldsymbol{\alpha}\|_2 \rightarrow +\infty} F(\mathbf{y}, \mathbf{G}\boldsymbol{\alpha}) + \lambda \mathcal{R}(\boldsymbol{\alpha}) = +\infty$.

Then, the solution set $\mathcal{V} = \operatorname{argmin}_{\boldsymbol{\alpha} \in \mathbb{R}^N} F(\mathbf{y}, \mathbf{G}\boldsymbol{\alpha}) + \lambda \mathcal{R}(\boldsymbol{\alpha})$, is **non empty**, **convex** and **compact**.³

The proof of this theorem can be deduced from [1, Proposition 8] (for reference only do not check it!).

³In finite dimension, a **compact** set is a **closed** and **bounded** set.

Proper, Convex, Lower Semi-Continuous Functional

Definition: (Proper Convex Functional)

A function $F: \mathbb{R}^N \rightarrow \mathbb{R} \cup \{-\infty, +\infty\}$ is called **convex** if

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^N, \forall \theta \in [0, 1]: \quad F(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) \leq \theta F(\mathbf{x}) + (1 - \theta) F(\mathbf{y}), \quad (4)$$

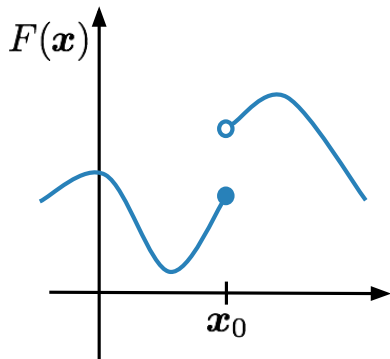
and **strictly convex** if the inequality in (4) is **strict**. If moreover, $F(\mathbf{x}) > -\infty$ for all $\mathbf{x} \in \mathbb{R}^N$ and $D = \{\mathbf{x} \in \mathbb{R}^N : F(\mathbf{x}) < +\infty\} \neq \emptyset$, then F is called a **proper (strictly) convex** function.⁴

Definition: (Lower Semi-Continuity)

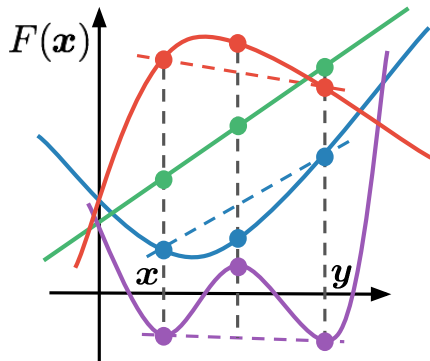
A function $F: \mathbb{R}^N \rightarrow \mathbb{R} \cup \{-\infty, +\infty\}$ is said **lower semi-continuous (lwsc)** at $\mathbf{x}_0 \in \mathbb{R}^N$ if for every $y < F(\mathbf{x}_0)$ there exists a **neighborhood** $U \subset \mathbb{R}^N$ of \mathbf{x}_0 such that $F(\mathbf{x}) \geq y, \forall \mathbf{x} \in U$.

⁴In short, a convex function is proper if its domain is nonempty and it never attains $-\infty$.

Proper, Convex, Lower Semi-Continuous Functional



Example of Lower
Semi-Continuous Function



Example of **Strictly Convex**,
Convex, **Concave** and **Non**
Convex Functions

Unicity of Solutions

Theorem: (Unicity of Solutions)

Assume that F and \mathcal{R} are as in Slide 19 and that the objective functional $\mathcal{J}(\alpha) := F(y, G\alpha) + \lambda \mathcal{R}(\alpha)$ is strictly convex. Then (3) admits a unique solution.

Proof: Assume that there exists at least two distinct solutions $\alpha_1, \alpha_2 \in \mathcal{V}$. Then, by the strict convexity of \mathcal{J} , we have $\forall \theta \in [0, 1]: \mathcal{J}(\theta \alpha_1 + (1 - \theta) \alpha_2) < \theta \mathcal{J}(\alpha_1) + (1 - \theta) \mathcal{J}(\alpha_2)$, and hence α_1, α_2 do not minimise \mathcal{J} which is a contradiction. \square

Sufficient conditions for the strict convexity of \mathcal{J} are: $F(y, \cdot)$ is strictly convex and G is injective, or \mathcal{R} is strictly convex. When \mathcal{J} is not strictly convex we can still retain a weaker form of unicity:

Theorem: (Unicity of Predicted Measurements)

Assume that F and \mathcal{R} are as in Slide 19 and that $F(y, \cdot)$ is strictly convex. Then there exists a unique $y^* \in \mathbb{R}^L$ such that $G\alpha^* = y^*$, $\forall \alpha^* \in \mathcal{V} = \operatorname{argmin}_{\alpha \in \mathbb{R}^N} \mathcal{J}(\alpha)$, i.e. every solution yield the same predicted measurements.

Proof (Unicity of Predicted Measurements)

$$\exists \alpha_1, \alpha_2 \in V \quad G\alpha_1 = y_1 \quad G\alpha_2 = y_2 \quad y_1 \neq y_2$$

$$J(\alpha) = F(y, G\alpha) + \lambda R(\alpha) = F_y(G\alpha) + \lambda R(\alpha)$$

$F_y(\cdot)$ is strictly convex. Consider $\alpha_3 = \theta\alpha_1 + (1-\theta)\alpha_2$
 $\in V$ convex

$$\begin{aligned} J(\alpha_3) &= F_y(\theta\alpha_1 + (1-\theta)\alpha_2) + \lambda R(\theta\alpha_1 + (1-\theta)\alpha_2) \\ &\leq \theta F_y(\alpha_1) + (1-\theta) F_y(\alpha_2) + \lambda (\theta R(\alpha_1) + (1-\theta) R(\alpha_2)) \\ &\stackrel{\text{strictly convex } F}{=} \theta J(\alpha_1) + (1-\theta) J(\alpha_2) = J^* \end{aligned}$$

J^* ← min value of J

Choosing the Cost Functional (Noiseless Case)

In a **noiseless setup**, one has *full trust* in the measurements. It is therefore natural to require that any solution of (3) be **consistent** with the data at hand, i.e. $\mathbf{y} = \mathbf{G}\boldsymbol{\alpha} \forall \boldsymbol{\alpha} \in \mathcal{V}$. This can be achieved by choosing the cost functional as $F(\mathbf{y}, \mathbf{G}\boldsymbol{\alpha}) = \iota(\mathbf{y} - \mathbf{G}\boldsymbol{\alpha})$, where $\iota: \mathbb{R}^L \rightarrow \{0, +\infty\}$ is the **indicator function**

$$\iota(\mathbf{z}) = \begin{cases} 0 & \text{if } \mathbf{z} = \mathbf{0}, \\ +\infty & \text{otherwise.} \end{cases}$$

Problem (3) becomes then a **generalised interpolation problem**:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \iota(\mathbf{y} - \mathbf{G}\boldsymbol{\alpha}) + \lambda \mathcal{R}(\boldsymbol{\alpha}) = \min_{\boldsymbol{\alpha} \in \mathbb{R}^N, \mathbf{y} = \mathbf{G}\boldsymbol{\alpha}} \mathcal{R}(\boldsymbol{\alpha}).$$

24

Penalised Problems with Strictly Convex Cost Functional are Interpolation Problems

Under the assumptions of the Theorem “**Unicity of Predicted Measurements**” on Slide 22 we have that:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^N} F(\mathbf{y}, \mathbf{G}\boldsymbol{\alpha}) + \lambda \mathcal{R}(\boldsymbol{\alpha}) = \min_{\boldsymbol{\alpha} \in \mathbb{R}^N, \mathbf{y}^* = \mathbf{G}\boldsymbol{\alpha}} \mathcal{R}(\boldsymbol{\alpha}),$$

for some (unknown) $\mathbf{y}^* \in \mathbb{R}^L$. Hence, every penalised optimisation problem with **strictly convex cost functional** is **equivalent to a generalised interpolation problem**.

Choosing the Cost Functional (Noisy Case)

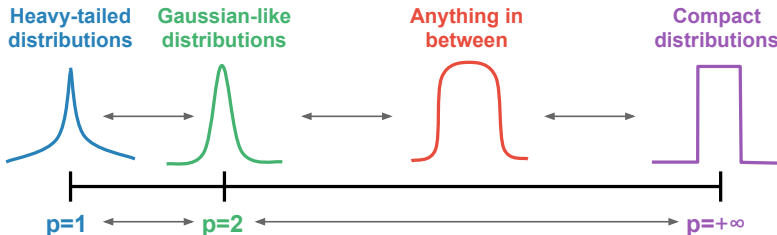
In a **noisy setup**, consistency is not desired anymore, as it almost always leads to **overfitting** the noisy data. One approach consists then in using the **negative log-likelihood** of the data \mathbf{y} as a measure of discrepancy:

$$F(\mathbf{y}, \mathbf{G}\boldsymbol{\alpha}) = -\ell(\boldsymbol{\alpha}|\mathbf{y}) = -\log p_{Y_1, \dots, Y_L}(y_1, \dots, y_L|\boldsymbol{\alpha}).$$

When the noise distribution is **not fully known** or the likelihood **too complex**, one can also use general ℓ_p cost functionals

$$F(\mathbf{y}, \mathbf{G}\boldsymbol{\alpha}) = \|\mathbf{y} - \mathbf{G}\boldsymbol{\alpha}\|_p^p = \sum_{i=1}^L \left| y_i - \sum_{n=1}^N G_{in} \alpha_n \right|^p,$$

where $p \in [1, +\infty]$ is typically chosen according to the **tail behaviour** of the noise distribution [2].



Example: Cost Functional for Gaussian Noise

Assume the following [multivariate Gaussian noise model](#):

$$\mathbf{Y} = \mathbf{G}\boldsymbol{\alpha} + \mathbf{N}, \quad \text{where} \quad \mathbf{N} \stackrel{d}{\sim} \mathcal{N}_L(\mathbf{0}, \boldsymbol{\Sigma}), \quad p_{\mathbf{N}}(\mathbf{y}) = \frac{1}{|\boldsymbol{\Sigma}|^{1/2} (2\pi)^{L/2}} \exp\left(-\frac{1}{2} \mathbf{y}^T \boldsymbol{\Sigma}^{-1} \mathbf{y}\right), \quad \mathbf{y} \in \mathbb{R}^L.$$

Then we have:

$$\begin{aligned} F(\mathbf{y}, \mathbf{G}\boldsymbol{\alpha}) &= -\ell(\boldsymbol{\alpha}|\mathbf{y}) = -\log p_{\mathbf{Y}}(\mathbf{y}|\boldsymbol{\alpha}) \\ &= -\log\left(\frac{1}{|\boldsymbol{\Sigma}|^{1/2} (2\pi)^{L/2}} \exp\left(-\frac{1}{2} (\mathbf{y} - \mathbf{G}\boldsymbol{\alpha})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{G}\boldsymbol{\alpha})\right)\right) \\ &= \frac{1}{2} \left\| \boldsymbol{\Sigma}^{-1/2} (\mathbf{y} - \mathbf{G}\boldsymbol{\alpha}) \right\|_2^2 + \underbrace{\frac{1}{2} \log |\boldsymbol{\Sigma}| + \frac{L}{2} \log(2\pi)}_{\text{Independent of } \boldsymbol{\alpha}} \\ &\propto \left\| \boldsymbol{\Sigma}^{-1/2} (\mathbf{y} - \mathbf{G}\boldsymbol{\alpha}) \right\|_2^2. \end{aligned}$$

This is the [weighted least-squares functional](#). For white noise, we have $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}_L$ and the cost functional becomes proportional to $\|\mathbf{y} - \mathbf{G}\boldsymbol{\alpha}\|_2^2$, which is the regular [least-squares functional](#).

Example: Cost Functional for Laplacian Noise

Assume the following **Laplacian white noise model**:

$$Y_i = (\mathbf{G}\boldsymbol{\alpha})_i + N_i, \quad \text{where } N_i \stackrel{\text{i.i.d.}}{\sim} \text{Laplace}(0, \sigma), \quad p_N(x) = \frac{1}{2\sigma} \exp\left(-\frac{|x|}{\sigma}\right), \quad x \in \mathbb{R}.$$

Then we have:

$$\begin{aligned} F(\mathbf{y}, \mathbf{G}\boldsymbol{\alpha}) &= -\ell(\boldsymbol{\alpha}|\mathbf{y}) = -\log p_{Y_1, \dots, Y_L}(y_1, \dots, y_L|\boldsymbol{\alpha}) \\ &= -\log\left(\frac{1}{(2\sigma)^L} \prod_{i=1}^L \exp\left(-\frac{|y_i - (\mathbf{G}\boldsymbol{\alpha})_i|}{\sigma}\right)\right) \\ &= \frac{1}{\sigma} \sum_{i=1}^L |y_i - (\mathbf{G}\boldsymbol{\alpha})_i| + \underbrace{L \log(2\sigma)}_{\text{Independent of } \boldsymbol{\alpha}} \\ &\propto \|\mathbf{y} - \mathbf{G}\boldsymbol{\alpha}\|_1. \end{aligned}$$

This is the **least absolute deviations** functional. It is **less affected by outliers** than the least-squares functional. The **weighted least absolute deviations** functional can also be defined but cannot be interpreted as the negative log-likelihood of a **multivariate Laplacian distribution**.

Example: Cost Functional for Poisson Noise

Assume **positive measurements** $Y: \Omega \rightarrow \mathbb{R}_+^L$ and the following **Poisson noise model**:

$$Y_i \stackrel{\text{ind}}{\sim} \text{Poisson}((\mathbf{G}\boldsymbol{\alpha})_i), \quad p_{Y_i}(k) = \frac{(\mathbf{G}\boldsymbol{\alpha})_i^k e^{-(\mathbf{G}\boldsymbol{\alpha})_i}}{k!}, \quad \forall k \in \mathbb{N}.$$

Then we have:

$$\begin{aligned} F(\mathbf{y}, \mathbf{G}\boldsymbol{\alpha}) &= -\ell(\boldsymbol{\alpha}|\mathbf{y}) = -\log p_{Y_1, \dots, Y_L}(y_1, \dots, y_L|\boldsymbol{\alpha}) \\ &= -\log \left(\prod_{i=1}^L \frac{(\mathbf{G}\boldsymbol{\alpha})_i^{y_i} e^{-(\mathbf{G}\boldsymbol{\alpha})_i}}{y_i!} \right) \\ &= \sum_{i=1}^L (\mathbf{G}\boldsymbol{\alpha})_i - y_i \log((\mathbf{G}\boldsymbol{\alpha})_i) + \underbrace{\log(y_i!)}_{\text{Independent of } \boldsymbol{\alpha}} \\ &\propto \sum_{i=1}^L (\mathbf{G}\boldsymbol{\alpha})_i - y_i \log((\mathbf{G}\boldsymbol{\alpha})_i) \end{aligned}$$

Example: Cost Functional for Poisson Noise (Continued)

$$\begin{aligned} &\propto \sum_{i=1}^L (\mathbf{G}\boldsymbol{\alpha})_i - y_i \log((\mathbf{G}\boldsymbol{\alpha})_i) + \underbrace{y_i \log(y_i) - y_i}_{\text{Can add anything independent of } \boldsymbol{\alpha}} \\ &= \sum_{i=1}^L y_i \log\left(\frac{y_i}{(\mathbf{G}\boldsymbol{\alpha})_i}\right) + (\mathbf{G}\boldsymbol{\alpha})_i - y_i \\ &= D_{KL}(\mathbf{y} || \mathbf{G}\boldsymbol{\alpha}), \end{aligned}$$

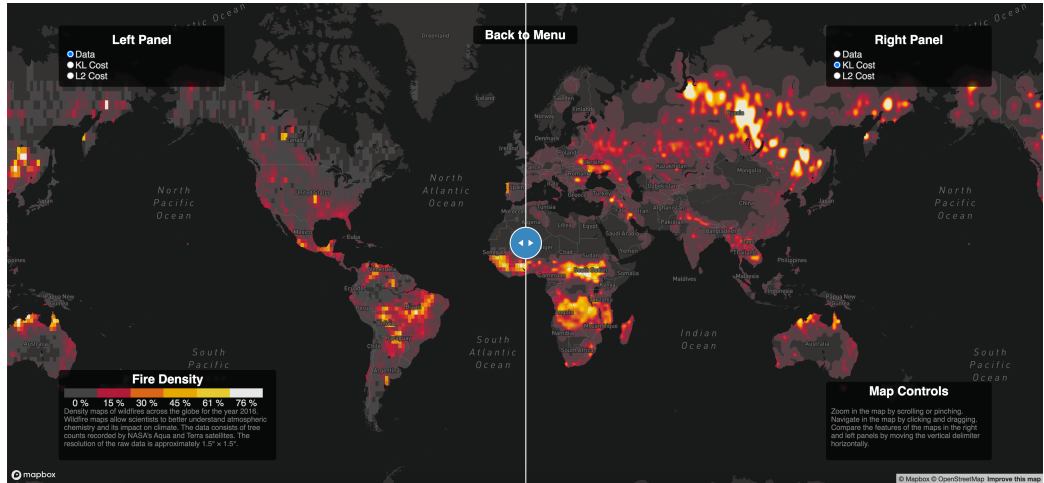
where

$$D_{KL}(\mathbf{y} || \mathbf{z}) = \sum_{i=1}^L y_i \log\left(\frac{y_i}{z_i}\right) - y_i + z_i, \quad \forall \mathbf{y}, \mathbf{z} \in \mathbb{R}_+^L, \quad (5)$$

is the **generalised Kullback-Leibler (KL) divergence** [3] for discrete **positive vectors** which **do not necessarily sum to one**. In information theory, and in the case where $\mathbf{1}^T \mathbf{z} = \mathbf{1}^T \mathbf{y} = 1$,⁵ the KL-divergence (5) can be interpreted as the **relative entropy** of \mathbf{y} with respect to \mathbf{z} , i.e. the **amount of information lost when** using \mathbf{z} to approximate \mathbf{y} . Note that the KL-divergence is **not a distance** (no symmetry/subadditivity).

⁵so that \mathbf{z} and \mathbf{y} can be interpreted as **discrete probability distributions**

Real-Life Example: Wild Fires



https://matthieumeo.github.io/fire_density.html

Choosing the Regularisation Functional

The regularisation functional is used to favour **physically-admissible** solutions with **simple behaviours**. It can be interpreted as implementing **Occam's razor principle**:

Occam's Razor Principle (*Lex parsimoniae*)

Occam's razor principle is a philosophical principle also known as the "**law of briefness**" or in Latin "**lex parsimoniae**". It was supposedly formulated by William of Ockham in the 14th century, who wrote in Latin "*Entia non sunt multiplicanda praeter necessitatem*". In English, this translates to "**More things should not be used than are necessary**".

In essence, this principle states that when two equally good explanations for a given phenomenon are available, one should always favour the simplest, i.e. the one that introduces **the least explanatory variables**.

What exactly is meant by "simple" solutions will depend on the specific application at hand.

(generalised) Tikhonov Regularisation

A common regularisation strategy consists in penalising the **squared ℓ_2 -norm** of the solutions, i.e.

$$\mathcal{R}(\boldsymbol{\alpha}) = \|\boldsymbol{\alpha}\|_2^2, \quad \boldsymbol{\alpha} \in \mathbb{R}^N. \quad (6)$$

This strategy is called **Tikhonov regularisation** and tends to favour **smooth solutions**. Different notions of smoothness can be achieved by introducing a **positive semi-definite finite-difference differential operator** $\mathbf{D} \in \mathbb{R}^{N \times N}$ in (6), yielding a **generalised Tikhonov regularisation**:

$$\mathcal{R}(\boldsymbol{\alpha}) = \|\mathbf{D}\boldsymbol{\alpha}\|_2^2, \quad \boldsymbol{\alpha} \in \mathbb{R}^N. \quad (7)$$

The Tikhonov functional (6) is **strictly convex**, hence yielding unique solutions when used in conjunction with a **convex cost functional**. The generalised Tikhonov functional (7) is **strictly convex** if \mathbf{D} is **injective** and simply **convex** otherwise. In the latter case, solutions to (3) **exist** if $\mathcal{N}(\mathbf{G}) \cap \mathcal{N}(\mathbf{D}) = \{\mathbf{0}\}$ and F is **coercive** but are in general **non unique**.⁶

⁶A sufficient condition for uniqueness is that F is **proper strictly convex**.

Form of Solutions with generalised Tikhonov Regularisation

Representer Theorem: (generalised Tikhonov Regularisation)

Assume that:

1. $\mathbf{G} \in \mathbb{R}^{L \times N}$ is **surjective** (i.e. full row rank), \mathbf{D} is **positive semi-definite** and $\mathcal{N}(\mathbf{G}) \cap \mathcal{N}(\mathbf{D}) = \{\mathbf{0}\}$.
2. $F(\mathbf{y}, \cdot) : \mathbb{R}^L \rightarrow \mathbb{R}_+$ is **proper strictly convex**, **coercive** and **lower semi-continuous** for every $\mathbf{y} \in \mathbb{R}^L$.

Then the optimisation problem:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^N} F(\mathbf{y}, \mathbf{G}\boldsymbol{\alpha}) + \lambda \|\mathbf{D}\boldsymbol{\alpha}\|_2^2$$

admits a **unique solution** which can be written as

$$\boldsymbol{\alpha}^* = (\mathbf{D}^T \mathbf{D})^\dagger \mathbf{G}^T \boldsymbol{\beta}^* + \boldsymbol{\gamma}^*,$$

for some $\boldsymbol{\beta}^* \in \mathbb{R}^L$ and $\boldsymbol{\gamma}^* \in \mathcal{N}(\mathbf{D})$.

When $\mathbf{D} = \mathbf{I}_N$ (standard Tikhonov regularisation) of \mathbf{D} is **invertible** then the theorem holds for F **proper convex** and **lwsc** and we get $\boldsymbol{\alpha}^* = (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{G}^T \boldsymbol{\beta}^*$ this case is discussed in [4, Corollary 7]).

Ridge Estimate

strict cvx + proper + lws c (coercive and strict)

$$\min_{\alpha \in \mathbb{R}^N} \underbrace{\frac{1}{2} \|y - G\alpha\|_2^2}_{J} + \underbrace{\left(\frac{\lambda}{2} \|\alpha\|_2^2\right)}_{D=I} \leftarrow \frac{\partial J}{\partial \alpha} = 0$$

Theorem: $\alpha^* = \underline{G^T \beta^*}$? $\beta^* \in \mathbb{R}^L$

$$\frac{\partial J}{\partial \alpha}(\alpha) = G^T G \alpha^* - G^T y + \lambda \alpha^* = 0$$

$$\Leftrightarrow (G^T G + \lambda I) \alpha^* = G^T y$$

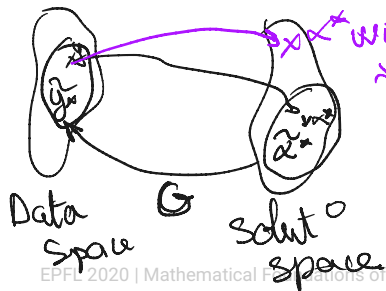
$$\Rightarrow \alpha^* = \underbrace{(G^T G + \lambda I)^{-1} G^T y}_{\in \mathcal{R}(G^T)}$$

Stability of Ridge Estimate

$$y \mapsto \alpha^*(y) = (G^T G + \lambda I)^{-1} G^T y \quad \text{bounded?}$$

$$\|\alpha^*(y)\|_2 \leq \|(G^T G + \lambda I)^{-1}\|_2 \|G^T\|_2 \|y\|_2$$

$$= \frac{\sqrt{\lambda_{\max}(G^T G)}}{\lambda_{\min}(G^T G) + \lambda} \|y\|_2 \quad \lambda > 0$$



$$\leq C \|y\|_2$$

$$C < +\infty$$

ℓ_1 /TV Regularisation

A common regularisation strategy consists in penalising the ℓ_1 -norm of the solutions, i.e.

$$\mathcal{R}(\alpha) = \|\alpha\|_1, \quad \alpha \in \mathbb{R}^N. \quad (8)$$

This strategy tends to favour **sparse solutions** with only a **few non zero coefficients**. Different notions of sparsity can be achieved by introducing a **positive semi-definite finite-difference differential operator** $D \in \mathbb{R}^{N \times N}$ in (8), yielding a **total variation (TV) regularisation**:

$$\mathcal{R}(\alpha) = \|D\alpha\|_1, \quad \alpha \in \mathbb{R}^N. \quad (9)$$

The ℓ_1 and TV functionals are **convex**. Solutions to (9) **exist** if $\mathcal{N}(G) \cap \mathcal{N}(D) = \{0\}$ and F is **coercive** but are in general **non unique**.⁷

Examples:

- **LASSO/Penalised Basis Pursuit:** $\min_{\alpha \in \mathbb{R}^N} \frac{1}{2} \|y - G\alpha\|_2^2 + \lambda \|\alpha\|_1.$
- **Generalised LASSO:** $\min_{\alpha \in \mathbb{R}^N} \frac{1}{2} \|y - G\alpha\|_2^2 + \lambda \|D\alpha\|_1.$

⁷Sufficient conditions for uniqueness are: G is **injective** and F is **strictly convex**.

Form of Solutions with TV Regularisation

Representer Theorem I: (TV Regularisation)

Assume that:

1. $G \in \mathbb{R}^{L \times N}$ is **invertible**, D is **positive semi-definite**.
2. $F(\mathbf{y}, \cdot) : \mathbb{R}^L \rightarrow \mathbb{R}_+$ is **proper strictly convex**, **coercive** and **lower semi-continuous** for every $\mathbf{y} \in \mathbb{R}^L$.

Then the optimisation problem:

$$\mathcal{V} = \underset{\boldsymbol{\alpha} \in \mathbb{R}^N}{\operatorname{argmin}} F(\mathbf{y}, G\boldsymbol{\alpha}) + \lambda \|D\boldsymbol{\alpha}\|_1$$

admits a **unique solution** of the form:

$$\boldsymbol{\alpha}^* = D^\dagger \boldsymbol{\beta}_K^* + \boldsymbol{\gamma}^*,$$

for some **K -sparse vector** $\boldsymbol{\beta}_K^* \in \mathbb{R}^N$, $K \leq L$ and $\boldsymbol{\gamma}^* \in \mathcal{N}(\mathbf{0})$.

When $D = I_N$ (ℓ_1 regularisation) or D is **invertible** then the theorem holds for F **proper strictly convex** and **lwsc** (no coercivity needed) and we have $\boldsymbol{\alpha}^* = D^{-1} \boldsymbol{\beta}_K^*$.

Form of Solutions with TV Regularisation

Representer Theorem II: (TV Regularisation)

Assume that:

1. $G \in \mathbb{R}^{L \times N}$ is **surjective** (i.e. full row rank), D is **positive semi-definite** and $\mathcal{N}(G) \cap \mathcal{N}(D) = \{\mathbf{0}\}$.
2. $F(\mathbf{y}, \cdot) : \mathbb{R}^L \rightarrow \mathbb{R}_+$ is **proper convex**, **coercive** and **lower semi-continuous** for every $\mathbf{y} \in \mathbb{R}^L$.

Then the solution set:

$$\mathcal{V} = \arg \min_{\alpha \in \mathbb{R}^N} F(\mathbf{y}, G\alpha) + \lambda \|D\alpha\|_1$$

is **non empty**, **compact** and the **convex-hull** of **extreme point solutions** of the form:

$$\alpha^* = D^\dagger \beta_K^* + \gamma^*,$$

for some **K -sparse vector** $\beta_K^* \in \mathbb{R}^N$, $K \leq L$ and $\gamma^* \in \mathcal{N}(D)$.

When $D = I_N$ (ℓ_1 regularisation) or D is **invertible** then the theorem holds for F **proper convex** and **lwsc** (no coercivity needed) and we have $\alpha^* = D^{-1} \beta_K^*$ (this case is discussed in [4, Corollary 8]).

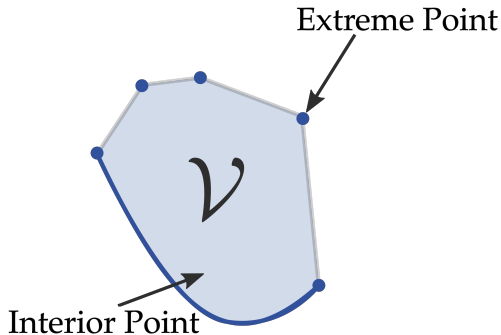
Solution Set is the Convex-Hull of Sparse Extreme Points

Definiton: (Extreme Point)

Let \mathcal{V} be a **convex set**. An **extreme point** $v \in \mathcal{V}$ is a point such that

$$\nexists (w, v) \in \mathcal{V}^2, \theta \in]0, 1[: \quad v = \theta w + (1 - \theta)v.$$

In plain words, v is a point in \mathcal{V} which **does not lie in any open line segment** joining two points of \mathcal{V} .



Example: Finite Difference Operator in $\mathbb{R}^{3 \times 3}$

$$D = \begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & -1 & 1 \end{bmatrix}$$

$$D \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 - x_1 \\ x_3 - x_2 \end{pmatrix} \rightarrow \text{finite-difference operator}$$

$$V = \operatorname{argmin}_{x \in \mathbb{R}^3} F(y, Gx) + \lambda \|Dx\|_1$$

$y \in \mathbb{R}^2, \quad G \in \mathbb{R}^{2 \times 3}$

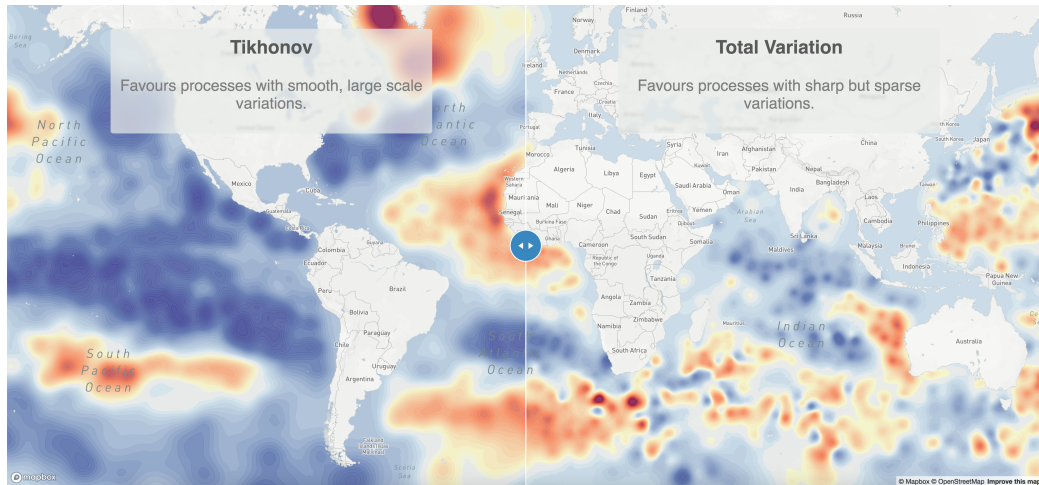
V convex-hull of extreme points of the form:

$$x^* = D^{-1} \beta_k^* = \sum_{k=1}^K \beta_k D_{:,n_k}^{-1} \quad k \leq 2$$

$k \leq 2$

$$D^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}$$

Real-Life Example: Sea Surface Temperatures



https://matthieumeo.github.io/tikhonov_vs_tv_en.html

Maximum Entropy Regularisation

The **maximum entropy** regularisation strategy considers the following regularisation functional:

$$\mathcal{R}(\boldsymbol{\alpha}) = \sum_{n=1}^N \alpha_i \log(\alpha_i), \quad \boldsymbol{\alpha} \in \mathbb{R}_+^N. \quad (10)$$

When $\mathbf{1}^T \boldsymbol{\alpha} = 1$ this is the negative **Shannon entropy** [5, 6] of $\boldsymbol{\alpha}$, a mathematical generalisation of entropy as introduced by Boltzmann in thermodynamics. It favours **positive, featureless** solutions which: smooth vectors indeed carry much less information than vectors with sharp, localised features, and hence have **higher entropy**.

This regularisation can be generalised by considering the **negative relative entropy** w.r.t. a **reference discrete distribution** $\boldsymbol{\eta} \in \mathbb{R}_+^N$:

$$\mathcal{R}(\boldsymbol{\alpha}) = \sum_{n=1}^N \alpha_i \log\left(\frac{\alpha_i}{\eta_i}\right), \quad \boldsymbol{\alpha} \in \mathbb{R}_+^N. \quad (11)$$

The functional (11) favours solutions with similar features as the reference distribution $\boldsymbol{\eta}$. Both functionals (10) and (11) are **strictly convex** and **coercive**, hence yielding unique solutions when they exist.

Form of Solutions with Maximum Entropy Interpolation

Representer Theorem: (Maximum Entropy Interpolation)

Consider the generalised interpolation problem:

$$\min_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^N \\ \mathbf{y} = \mathbf{G}\boldsymbol{\alpha}}} \sum_{i=1}^N \alpha_i \log(\alpha_i), \quad (12)$$

for some $\mathbf{y} \in \mathbb{R}^L$. If (12) admits a solution, the latter is unique and can be written as [7]:

$$\boldsymbol{\alpha}^\star = \gamma \exp(\mathbf{G}^T \boldsymbol{\beta}^\star),$$

for some $\gamma > 0$ and $\boldsymbol{\beta}^\star \in \mathbb{R}^L$.

The nonlinear exponential map kills low intensity features and boosts prominent ones.

Nonnegativity Regularisation

The **nonnegativity** regularisation strategy considers the following regularisation functional:

$$\mathcal{R}(\boldsymbol{\alpha}) = \begin{cases} 0 & \text{if } \boldsymbol{\alpha} \in \mathbb{R}_+^N, \\ +\infty & \text{otherwise.} \end{cases} \quad (13)$$

It constrains the solutions to be **positive**. This functional is **convex** and **non coercive**.

The functional (13) is sometimes replaced by the **log-barrier** functional (**convex** and **non coercive**):

$$\mathcal{R}(\boldsymbol{\alpha}) = - \sum_{n=1}^N \log(\alpha_i), \quad \boldsymbol{\alpha} \in \mathbb{R}_+^N,$$

which also **promotes positive solutions**.

Nonnegative Least-Squares (NNLS):

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}_+^N} \frac{1}{2} \|\mathbf{y} - \mathbf{G}\boldsymbol{\alpha}\|_2^2.$$

Form of Solutions with Nonnegativity Constraints

Representer Theorem: (Nonnegativity Constraints)

Assume that:

1. $G \in \mathbb{R}^{L \times N}$ is **injective**.
2. $F(y, \cdot) : \mathbb{R}^L \rightarrow \mathbb{R}_+$ is **proper strictly convex, coercive** and **lower semi-continuous** for every $y \in \mathbb{R}^L$.

Then the optimisation problem:

$$\min_{\alpha \in \mathbb{R}_+^N} F(y, G\alpha)$$

admits a **unique L -sparse solution**.

The proof to this Theorem follows from [8, Proposition 4.1].

Bayesian Interpretation

In certain cases, the penalised optimisation problem (3) can be interpreted as a **maximum a posteriori (MAP)** problem. Adopting a **Bayesian view**, assume for example a **Gaussian a priori distribution** for α and a **Gaussian likelihood function** (i.e. a **Gaussian white noise model**):

$$p(\alpha) \propto \exp\left(-\frac{1}{2\xi^2} \|\alpha\|_2^2\right), \quad p(y|\alpha) \propto \exp\left(-\frac{1}{2\sigma^2} \|y - G\alpha\|_2^2\right). \quad (14)$$

From **Baye's theorem**, the **posterior distribution** of α knowing the data y is then given by

$$p(\alpha|y) = \frac{p(y|\alpha)p(\alpha)}{\int_{\mathbb{R}^N} p(y|\alpha)p(\alpha) d\alpha}. \quad (15)$$

A **maximum a posteriori (MAP)** estimate is then defined as

$$\alpha_{MAP}^* \in \arg\max_{\alpha \in \mathbb{R}^N} p(\alpha|y) = \arg\max_{\alpha \in \mathbb{R}^N} p(y|\alpha)p(\alpha) = \arg\max_{\alpha \in \mathbb{R}^N} L(\alpha|y)p(\alpha) = \arg\min_{\alpha \in \mathbb{R}^N} -\ell(\alpha|y) - \log(p(\alpha)).$$

For the prior distribution and likelihood assumed in (14) this yields:

$$\alpha_{MAP}^* \in \arg\min \frac{1}{2\sigma^2} \|y - G\alpha\|_2^2 + \frac{1}{2\xi^2} \|\alpha\|_2^2 = \arg\min \frac{1}{2} \|y - G\alpha\|_2^2 + \frac{\lambda}{2} \|\alpha\|_2^2, \quad \text{with } \lambda = \frac{\sigma^2}{\xi^2}. \quad (16)$$

Bayesian Interpretation (continued)

We recognise in (16) the **Ridge estimate** (see Slide 34), obtained when choosing a **least-squares cost functional** and a **Tikhonov** penalty in (3). Notice moreover that the regularisation parameter λ is equal to the **ratio** of the likelihood and prior variances.

Note that the prior distribution can be **improper** (i.e. *not summable*) as long as $\int_{\mathbb{R}^N} p(\mathbf{y}|\boldsymbol{\alpha})p(\boldsymbol{\alpha})d\boldsymbol{\alpha} < +\infty$ so that (15) is still **well-defined**.

This allows us to extend the previous analysis to many classical optimisation problems:

Example: (Penalised Optimisation Pbs as MAP)

- **Weighted Least-Squares with Generalised Tikhonov** ($\min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \frac{1}{2} \|\boldsymbol{\Sigma}^{-1/2}(\mathbf{y} - \mathbf{G}\boldsymbol{\alpha})\|_2^2 + \frac{1}{2} \|\mathbf{D}\boldsymbol{\alpha}\|_2^2$):
 $p(\boldsymbol{\alpha}) \propto \exp\left(-\frac{1}{2} \boldsymbol{\alpha}^T \mathbf{D}^T \mathbf{D} \boldsymbol{\alpha}\right), \quad p(\mathbf{y}|\boldsymbol{\alpha}) \propto \exp\left(-\frac{1}{2} (\mathbf{y} - \mathbf{G}\boldsymbol{\alpha})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{G}\boldsymbol{\alpha})\right), \quad \lambda = 1.$
- **Least Absolute Deviations with Tikhonov** ($\min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{G}\boldsymbol{\alpha}\|_1 + \frac{\lambda}{2} \|\boldsymbol{\alpha}\|_2^2$):
 $p(\boldsymbol{\alpha}) \propto \exp\left(-\frac{1}{2\xi^2} \|\boldsymbol{\alpha}\|_2^2\right), \quad p(\mathbf{y}|\boldsymbol{\alpha}) \propto \exp\left(-\frac{1}{\sigma} \|\mathbf{y} - \mathbf{G}\boldsymbol{\alpha}\|_1\right), \quad \lambda = \frac{\sigma}{\xi^2}.$

Bayesian Interpretation (continued)

Example: (Penalised Optimisation Pbs as MAP)

- **LASSO** ($\min_{\alpha \in \mathbb{R}^N} \frac{1}{2} \|\mathbf{y} - \mathbf{G}\alpha\|_2^2 + \lambda \|\alpha\|_1$):
 $p(\alpha) \propto \exp\left(-\frac{1}{\xi} \|\alpha\|_1\right), \quad p(\mathbf{y}|\alpha) \propto \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{G}\alpha\|_2^2\right), \quad \lambda = \frac{\sigma^2}{\xi}.$
- **Generalised LASSO** ($\min_{\alpha \in \mathbb{R}^N} \frac{1}{2} \|\mathbf{y} - \mathbf{G}\alpha\|_2^2 + \lambda \|\mathbf{D}\alpha\|_1$):
 $p(\alpha) \propto \exp(-\|\mathbf{D}\alpha\|_1), \quad p(\mathbf{y}|\alpha) \propto \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{G}\alpha\|_2^2\right), \quad \lambda = \sigma^2.$
- **Least-squares with Maximum Entropy** ($\min_{\alpha \in \mathbb{R}^N} \frac{1}{2} \|\mathbf{y} - \mathbf{G}\alpha\|_2^2 + \lambda \sum_{n=1}^N \alpha_n \log(\alpha_n)$):
 $p(\alpha) \propto \exp\left(-\sum_{n=1}^N \alpha_n \log(\alpha_n)\right), \quad p(\mathbf{y}|\alpha) \propto \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{G}\alpha\|_2^2\right), \quad \lambda = \sigma^2.$
- **Nonnegative Least-Squares** ($\min_{\alpha \in \mathbb{R}_+^N} \|\mathbf{y} - \mathbf{G}\alpha\|_2$):
 $p(\alpha) \propto \exp\left(-\iota_{\mathbb{R}_+^N}(\alpha)\right), \quad p(\mathbf{y}|\alpha) \propto \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{G}\alpha\|_2^2\right).$
- **KL-Divergence with Tikhonov** ($\min_{\alpha \in \mathbb{R}^N} D_{KL}(\mathbf{y}||\mathbf{G}\alpha) + \frac{\lambda}{2} \|\alpha\|_2^2$):
 $p(\alpha) \propto \exp\left(-\frac{1}{2\xi^2} \|\alpha\|_2^2\right), \quad p(\mathbf{y}|\alpha) \propto \exp\left(-D_{KL}(\mathbf{y}||\mathbf{G}\alpha)\right), \quad \lambda = \frac{1}{\xi^2}.$

References I

- [1] Harshit Gupta, Julien Fageot, and Michael Unser.
Continuous-domain solutions of linear inverse problems with tikhonov versus generalized tv regularization.
IEEE Transactions on Signal Processing, 66(17):4670–4684, 2018.
- [2] John R Rice and John S White.
Norms for smoothing and estimation.
SIAM review, 6(3):243–256, 1964.
- [3] Mario Bertero, Patrizia Boccacci, Giorgio Talenti, Riccardo Zanella, and Luca Zanni.
A discrepancy principle for poisson data.
Inverse problems, 26(10):105004, 2010.
- [4] Michael Unser, Julien Fageot, and Harshit Gupta.
Representer theorems for sparsity-promoting l1 regularization.
IEEE Transactions on Information Theory, 62(9):5167–5180, 2016.
- [5] Ramesh Narayan and Rajaram Nityananda.
Maximum entropy image restoration in astronomy.
Annual review of astronomy and astrophysics, 24(1):127–170, 1986.

References II

- [6] David L Donoho, Iain M Johnstone, Jeffrey C Hoch, and Alan S Stern.
Maximum entropy and the nearly black object.
Journal of the Royal Statistical Society: Series B (Methodological), 54(1):41–67, 1992.
- [7] Lawrence D Brown.
Fundamentals of statistical exponential families: with applications in statistical decision theory.
Ims, 1986.
- [8] Claire Boyer, Antonin Chambolle, Yohann De Castro, Vincent Duval, Frédéric De Gournay, and Pierre Weiss.
On representer theorems and convex regularization.
SIAM Journal on Optimization, 29(2):1260–1281, 2019.