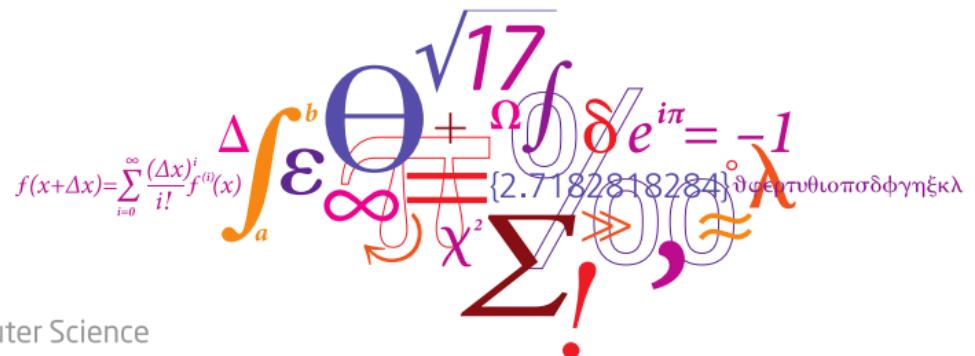


02471 Machine Learning for Signal Processing

## Sparsity-aware learning

Tommy Sonne Alstrøm

## Cognitive Systems section



# Outline

- Course admin
- Last week review
- Sparsity-aware learning
- Norms
- LASSO
- Next Week

Material: ML 8.2, 8.10.1–8.10.2, 9.1–9.5, 9.9.

- Feedback from last week:
  - List of symbols (todo, suggestion from week 4).
  - Latex files provided for problem set 2.
- Feedback from you is a critical component for improving both the course and my teaching.
- Type of feedback
  - Mention one thing that worked?
  - Mention one thing should be improved (both in current lecture and last weeks exercise)?
  - Mention one thing you would change if you gave the lecture.

## The story so far and what the future holds

- Problem set 1 submission corrections pending.
- Problem set 2 can be solved, except problem 2.6, which uses material for week 7.

What you have learned so far:

- Parameter estimation [L2 regularization, biased estimation, mean squared error minimization]. **Todo:** L1 regularization (this week), Bayesian parameter estimation (in 4 weeks)
- Filtering signals [Stochastic processes, correlation functions, Wiener filter, linear prediction, adaptive filtering using stochastic gradient decent (LMS, APA/NLMS), adaptive filtering using regularization (RLS)]

The topic the next three weeks will be sparse signal representations and dictionary learning

- Signal representations [Time frequency analysis, sparsity aware learning, factor models]
- Sparsity aware sensing (lasso, sparse priors), compressed sensing, dictionary learning [Independent component analysis, Non-negative matrix factorization,  $k$ -SVD]

# Learning objectives

## Learning objectives

A student who has met the objectives of the course will be able to:

- Explain, apply and analyze properties of discrete time signal processing systems
- Apply the short time Fourier transform to compute the spectrogram of a signal and analyze the signal content
- Explain compressed sensing and determine the relevant parameters in specific applications
- Deduce and determine how to apply factor models such as non-negative matrix factorization (NMF), independent component analysis (ICA) and sparse coding
- Deduce and apply correlation functions for various signal classes, in particular for stochastic signals
- Analyze filtering problems and demonstrate the application of least squares filter components such as the Wiener filter
- Describe, apply and derive non-linear signal processing methods based such as kernel methods and reproducing kernel Hilbert space for applications such as denoising
- Derive maximum likelihood estimates and apply the EM algorithm to learn model parameters
- Describe, apply and derive state-space models such as Kalman filters and Hidden Markov models
- Solve and interpret the result of signal processing systems by use of a programming language
- Design simple signal processing systems based on an analysis of involved signal characteristics, the objective of the processing system, and utility of methods presented in the course
- Describe a number of signal processing applications and interpret the results

## Last week review

Last week review

## The RLS algorithm

RLS minimized the exponentially weighted least-squares cost function

### Exponentially weighted least-squares

$$J(\boldsymbol{\theta}, \beta, \lambda) = \sum_{i=0}^n \beta^{n-i} (y_i - \boldsymbol{\theta}^T \mathbf{x}_i)^2 + \lambda \beta^{n+1} \|\boldsymbol{\theta}\|^2$$

datafit      L.S      Reg.      + g(θ)

Choose forget factor  $\beta$  and regularization  $\lambda$ .

We have an iterative update formula: update  $\boldsymbol{\theta}_n$  at each new observation  $(\mathbf{x}_n, y_n, )$

### Update formula

RLS  
algo

$$\begin{aligned} e_n &= y_n - \boldsymbol{\theta}_{n-1}^T \mathbf{x}_n \\ \mathbf{z}_n &= P_{n-1} \mathbf{x}_n \\ \mathbf{k}_n &= \mathbf{z}_n / (\beta + \mathbf{x}_n^T \mathbf{z}_n) \\ \boldsymbol{\theta}_n &= \boldsymbol{\theta}_{n-1} + \mathbf{k}_n e_n \\ P_n &= \beta^{-1} P_{n-1} - \beta^{-1} \mathbf{k}_n \mathbf{z}_n^T \end{aligned}$$

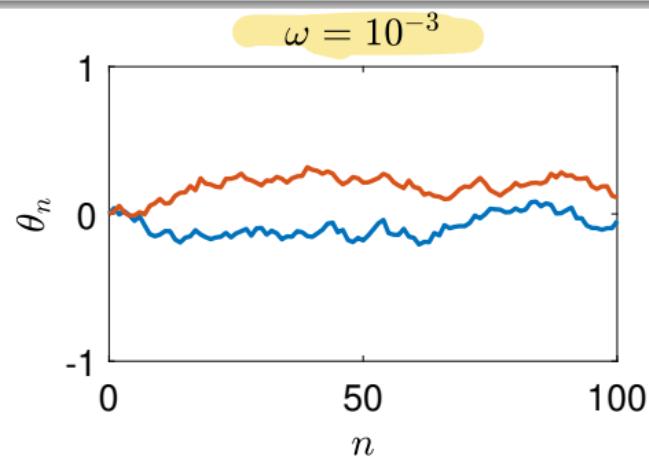
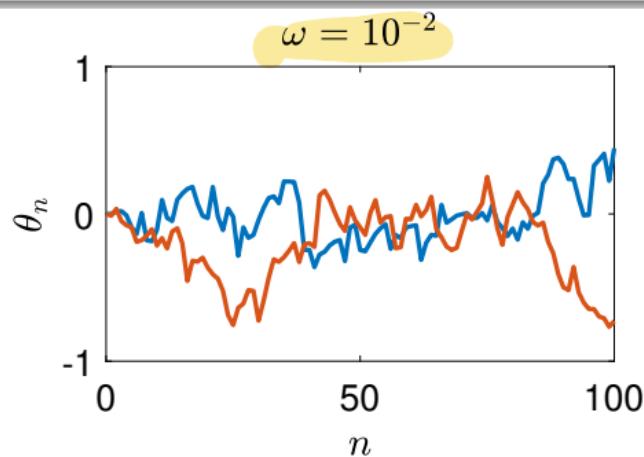
Last week review

## How do we simulate non-stationary environment?

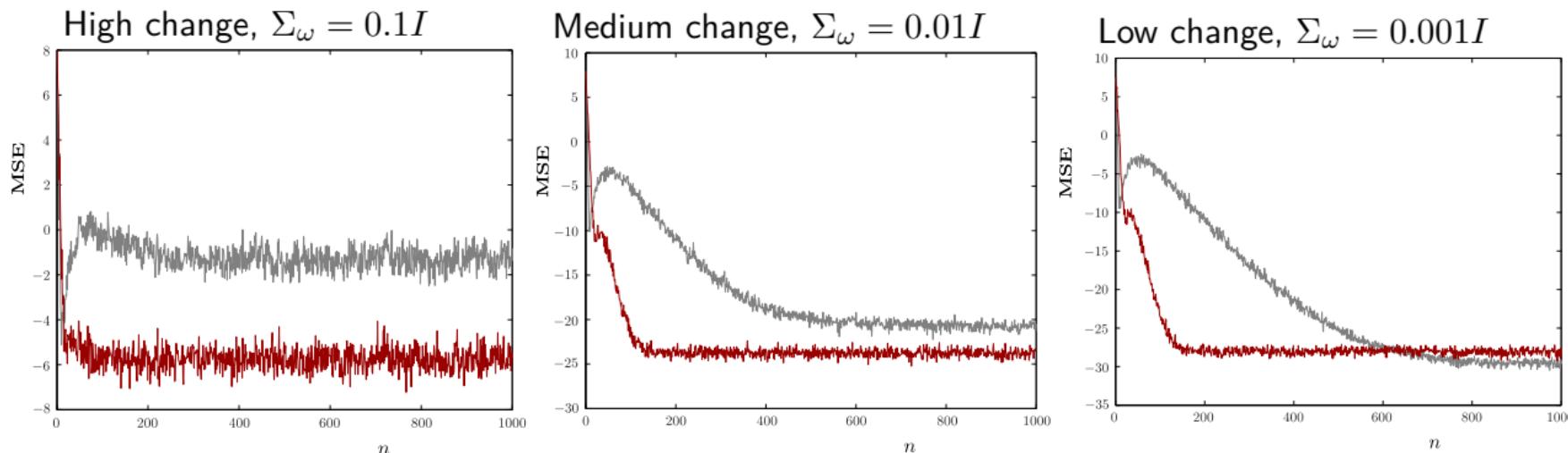
### Non-stationary environment

$$y_n = \theta_o^T x_n + \eta_n \quad \text{obs}$$
$$\theta_{o,n} = \alpha \theta_{o,n-1} + \omega_n \quad \text{state}$$

- $\alpha < 1$  is the autoregressive coefficient (memory).
- $\omega_n$  is zero mean Gaussian i.i.d with diagonal covariance (white noise).



# Convergence rates of RLS, time-tracking



Gray: RLS, Red: NLMS.  $y$ -axis is in dB.

Convergence results explained by  $J_{\text{exc}} = J_{\text{excess}} + J_{\text{lag}}$ .

For NLMS:  $J_{\text{exc}} = \frac{1}{2}\mu\sigma_\eta^2 \text{trace}\{\Sigma_x\} \frac{1}{\sigma_x^2(l-2)} + \frac{1}{2}\mu^{-1} \text{trace}\{\Sigma_x\} \text{trace}\{\Sigma_\omega\}$

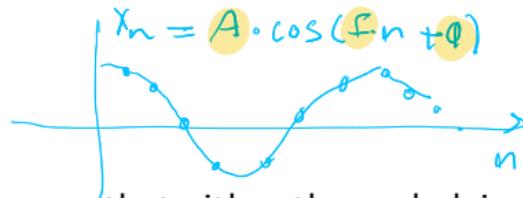
For RLS:  $J_{\text{exc}} = \frac{1}{2}(1-\beta)\sigma_\eta^2 l + \frac{1}{2}(1-\beta)^{-1} \text{trace}\{\Sigma_\omega \Sigma_x\}$ .

- RLS can under certain circumstances converge faster than NLMS
  - There exist more theory on convergence analysis of RLS vs LMS, on when to choose RLS over LMS, it depends on the signal properties
- We need to consider both time-tracking performance and stationary performance.
- We have models to simulate both scenarios.
- Even though the algorithms are "old", they are still very relevant.

## Sparsity-aware learning

## Sparsity-aware learning

### What is sparsity-aware learning



In a number of practical problems, it is known that either the underlying model is sparse or it is sparse in a “transform” domain, e.g., in the Fourier transform domain.

By **sparse models** is meant that most of the unknown coordinates in the model are zero.

Recall from Chapter 3, the use of a regularizer can shrink the norm of the obtained solution. In this vein, we will see that by adopting appropriate regularizers, one can help the optimization process to identify the coordinates of the zeros.

Leads to **compressed sensing**, where the goal is to directly acquire **as few samples as possible** that encode the minimum information, which is needed to obtain a **compressed signal representation**.

Algorithms covered so far will struggle with either **non-white noise**, or **presence of outliers**.

# Sparsity-aware learning

## Starting point

A linear system with multiple solutions (under-determined system):

$$y = \mathbf{a} \cdot \mathbf{x} \Rightarrow 2 = \begin{bmatrix} 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\hat{\mathbf{x}}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \ell_2$$

$$\hat{\mathbf{x}}_2 = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

$$\hat{\mathbf{x}}_3 = \begin{bmatrix} 0 \\ 2 \end{bmatrix}$$

Which solution would you pick?

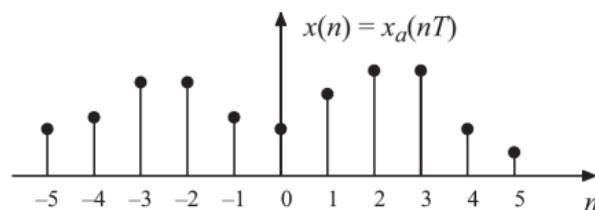
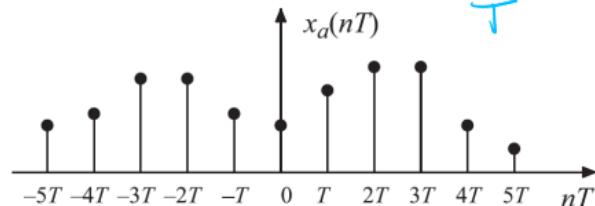
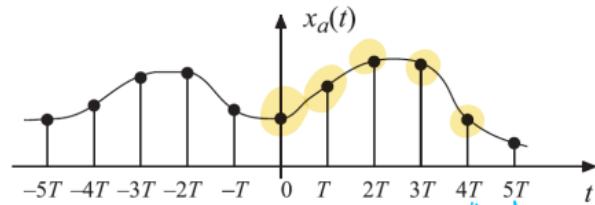
$$\|\hat{\mathbf{x}}_1\|_0 = 2 \quad \|\hat{\mathbf{x}}_1\|_1 = \|\hat{\mathbf{x}}_2\|_2$$

$$\|\hat{\mathbf{x}}_2\|_0 = 1 \quad \|\hat{\mathbf{x}}_1\|_2 < \|\hat{\mathbf{x}}_2\|_2$$

## Sparsity-aware learning

**Traditional sampling rate**

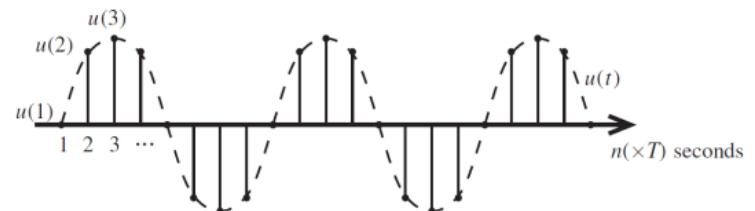
Traditional sampling



$$F_{\max} \leq 2 \cdot F_S$$

 $10 \text{ Hz}$  $20 \text{ Hz}$ 

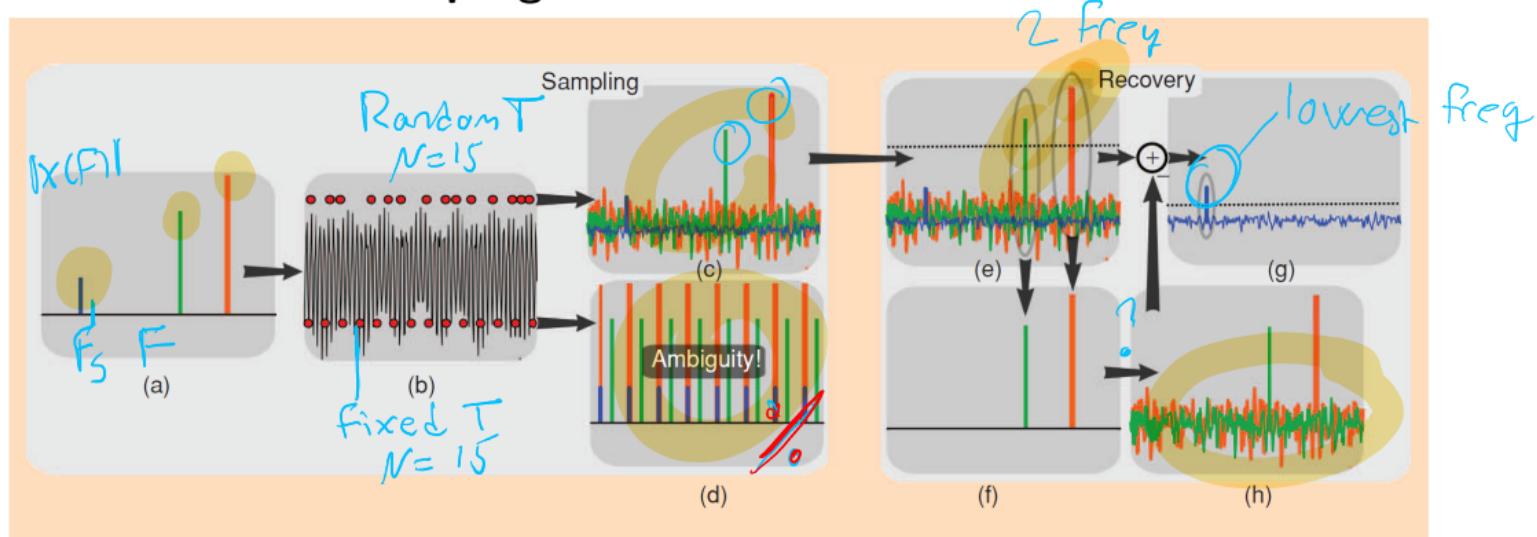
Suppose you were given this signal:



Can you sample more clever?

## Sparsity-aware learning

## Traditional under-sampling

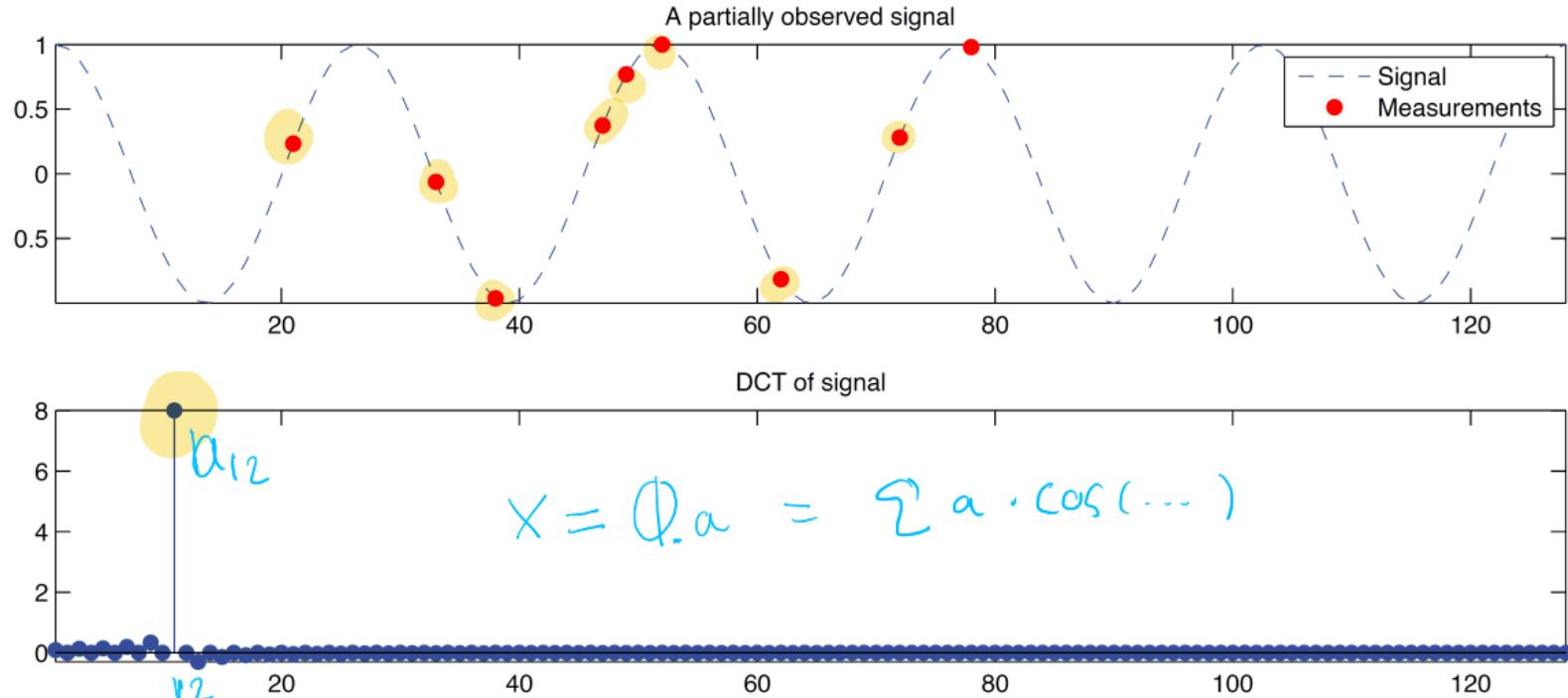


[FIG5] Heuristic procedure for reconstruction from undersampled data. A sparse signal (a) is 8-fold undersampled in its 1-D  $k$ -space domain (b). Equispaced undersampling results in signal aliasing (d) preventing recovery. Pseudo-random undersampling results in incoherent interference (c). Some strong signal components stick above the interference level, are detected and recovered by thresholding (e) and (f). The interference of these components is computed (g) and subtracted (h), thus lowering the total interference level and enabling recovery of weaker components.

M. Lustig, D. L. Donoho, J. M. Santos and J. M. Pauly, "Compressed Sensing MRI," in IEEE Signal Processing Magazine, vol. 25, no. 2, pp. 72-82, March 2008.

## Sparsity-aware learning

### Random sampling

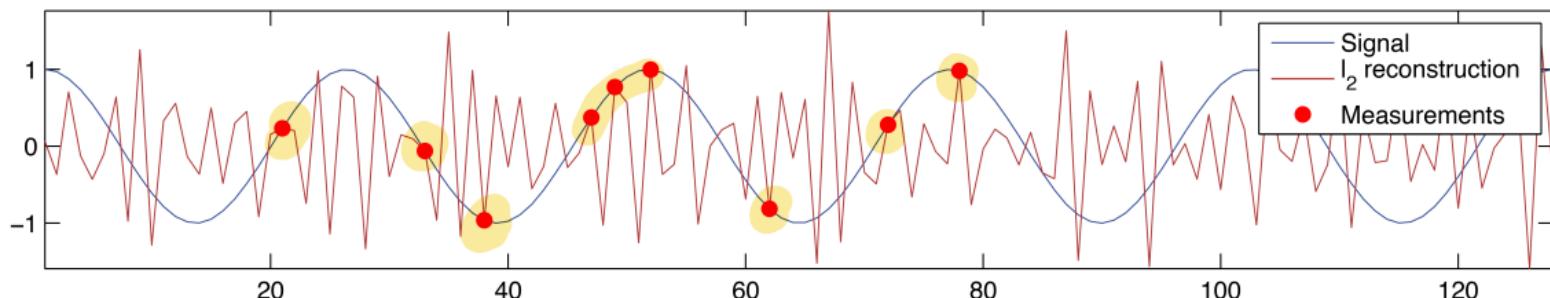


Example courtesy Prof. Paris Smaragdis

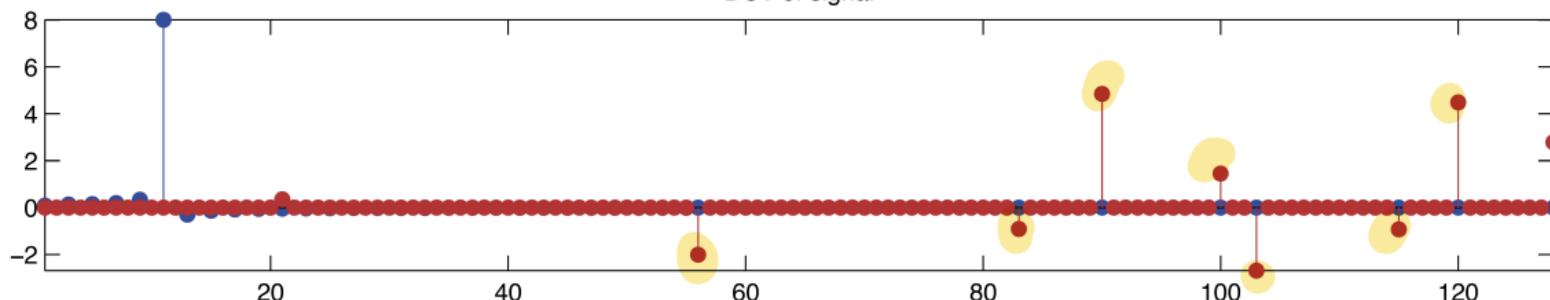
# Sparsity-aware learning Reconstruction using $\ell_2$

Ridge Reg

A partially observed signal

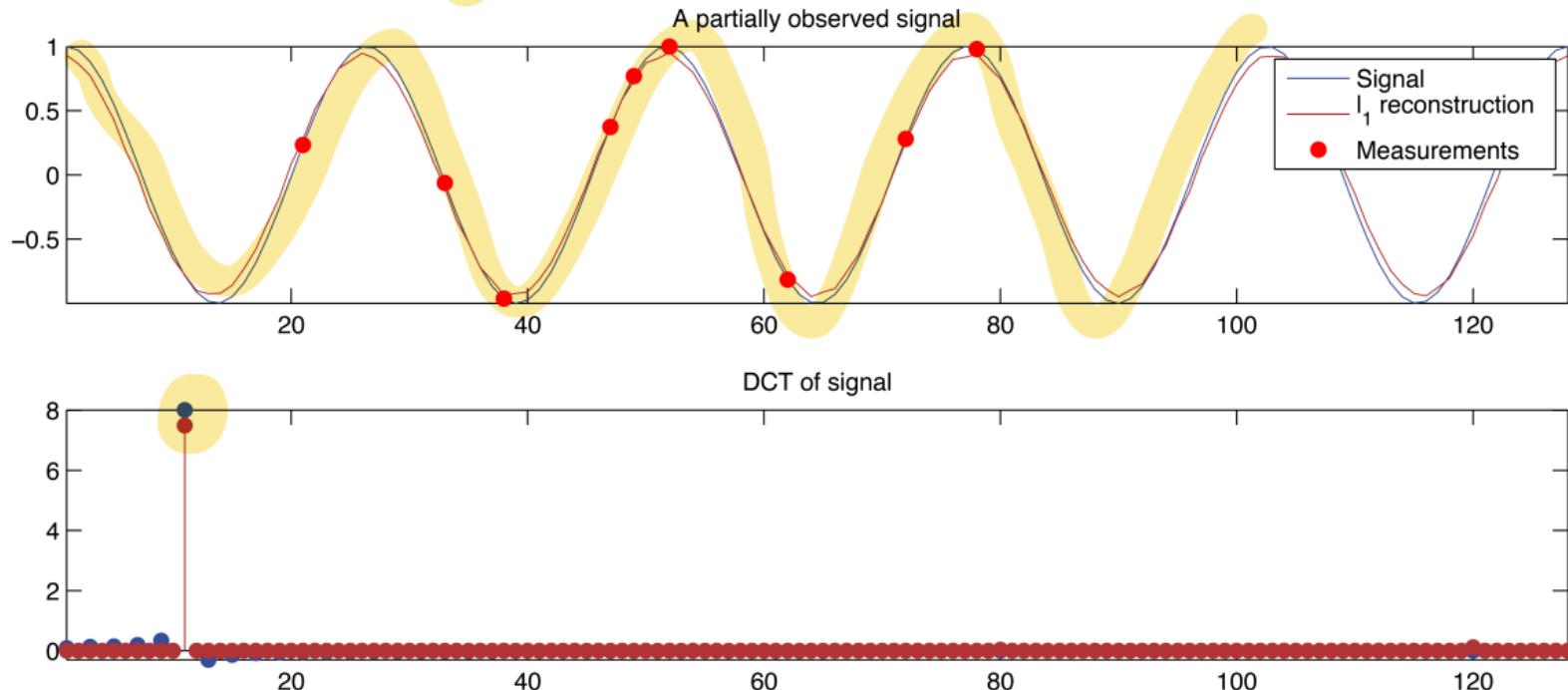


DCT of signal



Example courtesy Prof. Paris Smaragdis

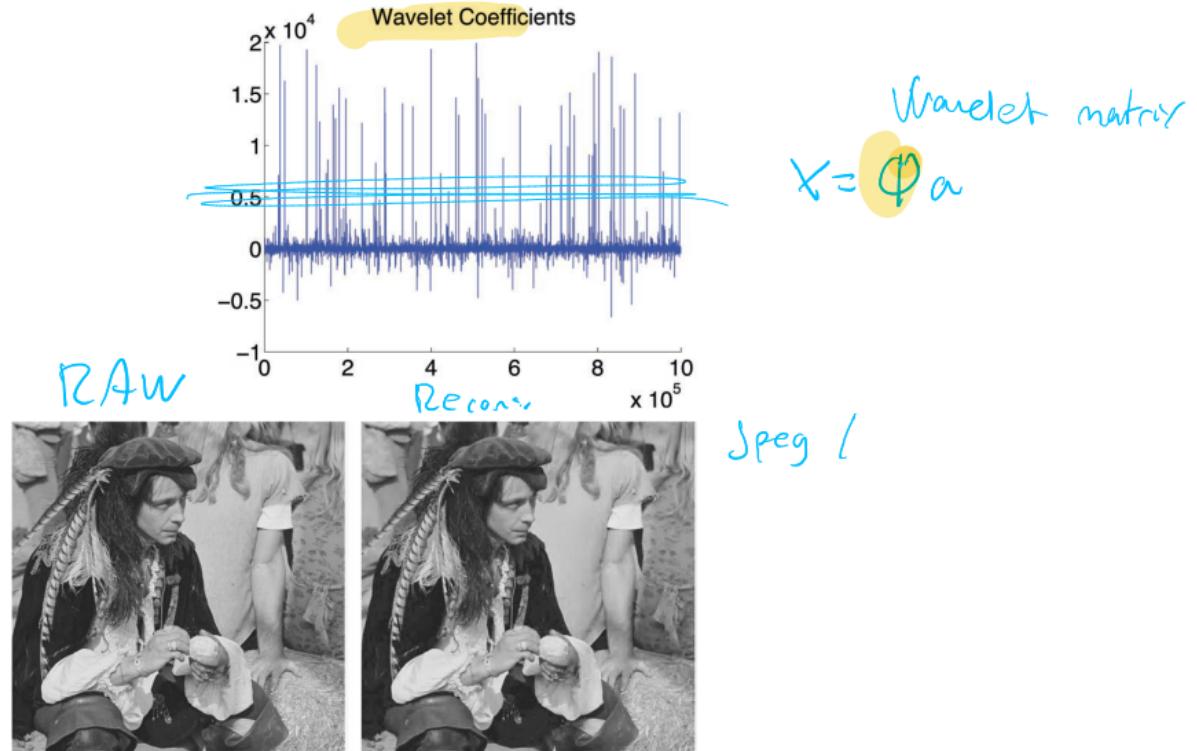
## Sparsity-aware learning

Reconstruction using  $\ell_1$ 

Example courtesy Prof. Paris Smaragdis

# Sparsity-aware learning

## Wavelet compression



Sparsity-aware learning have many cool applications, in statistics, signal processing, and machine learning.

- Signal recovery in under-determined system.
- Leads to dictionary learning. *Week 8*
- ... and to compressed sensing, where we sample and compress at the same time.

## Norms

$$\lambda \|\theta\|_1^1$$

## Ridge regression

$$J(\theta, \lambda) = \sum_{i=1}^N (y_i - \theta^T x_i)^2 + \lambda \|\theta\|_2^2 \quad \text{eq.(3.39)}$$

DESIGN

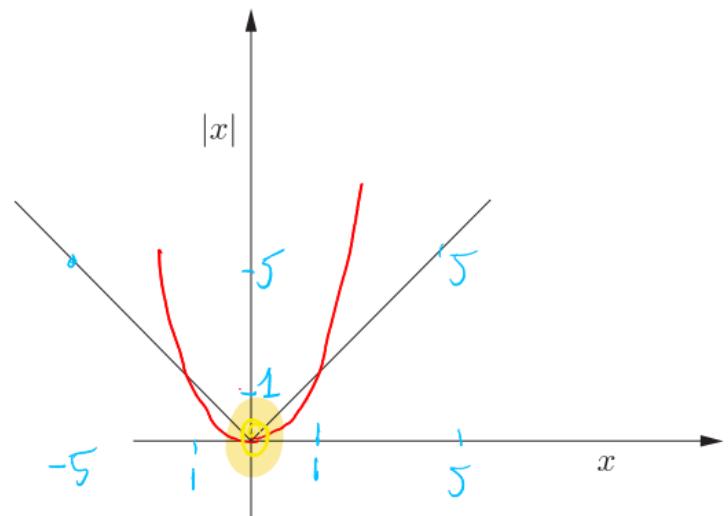
Can we think of another way to construct our regularization to promote sparsity

$$\|\begin{bmatrix} 1 \\ 0 \end{bmatrix}\| > \|\begin{bmatrix} 2 \\ 0 \end{bmatrix}\| \quad \checkmark$$

## How about the absolute value?

$$\frac{\partial J(\theta)}{\partial \theta} = 0$$

$$f(x) = |x| \quad \ell_1 \text{ in } 1d$$



Is there any problems in using this function?

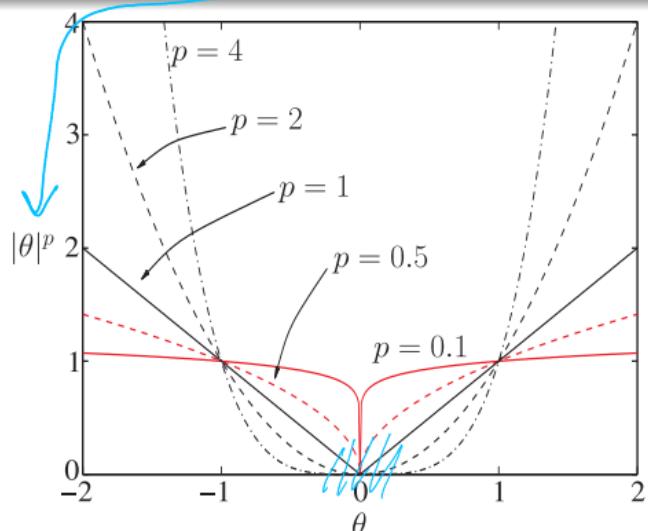
# The $\ell_p$ norm

## The $\ell_p$ norm

$$\|\theta\|_p := \left( \sum_{i=1}^l |\theta_i|^p \right)^{1/p}$$

$$\|\theta\|_p^p = \sum_{i=1}^l |\theta_i|^p$$

$l$  is the size of vector  $\theta$ .



## Norms

## Norms

## The $\ell_p$ norm

$$\|\boldsymbol{\theta}\|_p := \left( \sum_{i=1}^l |\theta_i|^p \right)^{1/p}, \quad p \geq 1$$

ex 6-1

$l$  is the size of vector  $\theta$ .

## Definition of norm

Let  $V$  be a vector space. A norm on  $V$  is a function  $\|\cdot\| : V \rightarrow \mathbb{R}$  that satisfies the following three conditions,  $\forall \theta \in V$ :

- positive

  - $\|\theta\|_p \geq 0$ ,  $\|\theta\|_p = 0 \Leftrightarrow \theta = 0$
  - \*  $\|\alpha\theta\|_p = |\alpha| \|\theta\|_p$ ,  $\forall \alpha \in \mathbb{R}$
  - \*  $\|\theta_1 + \theta_2\|_p \leq \|\theta_1\|_p + \|\theta_2\|_p$  (triangle inequality)

$$\begin{array}{ccc}
 \text{LHS} & & \text{RHS} \\
 \theta_1 = 1 \quad \theta_2 = 1 \quad p = 1/2 & & \left(\frac{1^{1/2}}{2}\right)^{1/2} = 1 \\
 \left(1^{1/2} + 1^{1/2}\right)^{1/(1/2)} = 2^2 = 4 & \leq & 2
 \end{array}$$

Why is  $\ell_p$  not a norm for  $p < 1$ ? can you create an example that breaks the triangle inequality?

## Norms

How does the  $\ell_p$  norm behave

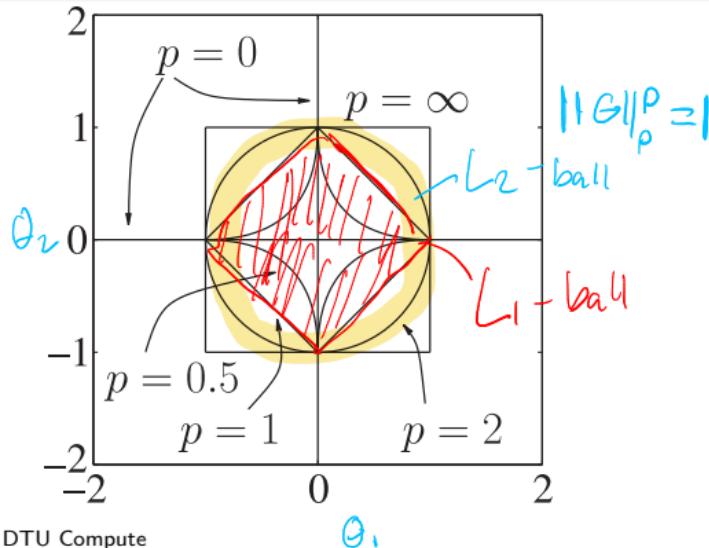
## Special cases

$$\|\theta\|_\infty := \arg \max_i |\theta_i| \quad (\text{max element})$$

$$\sum |\theta_i|^p$$

$$\sum |\theta_i|^{0.0001}$$

$$\|\theta\|_0 := |\{i \mid \theta_i \neq 0, i = 1, \dots, l\}| \quad (\text{number of nonzeros})$$



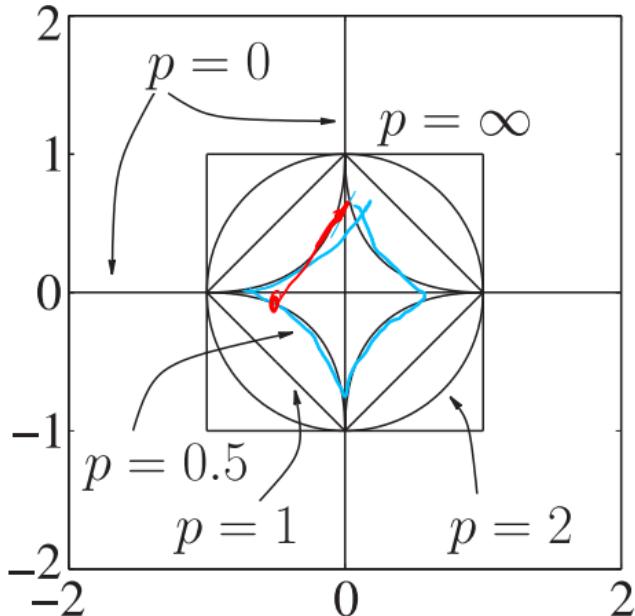
Notation remarks:

$|x|$  is the numerical value when  $x \in \mathbb{R}$ .

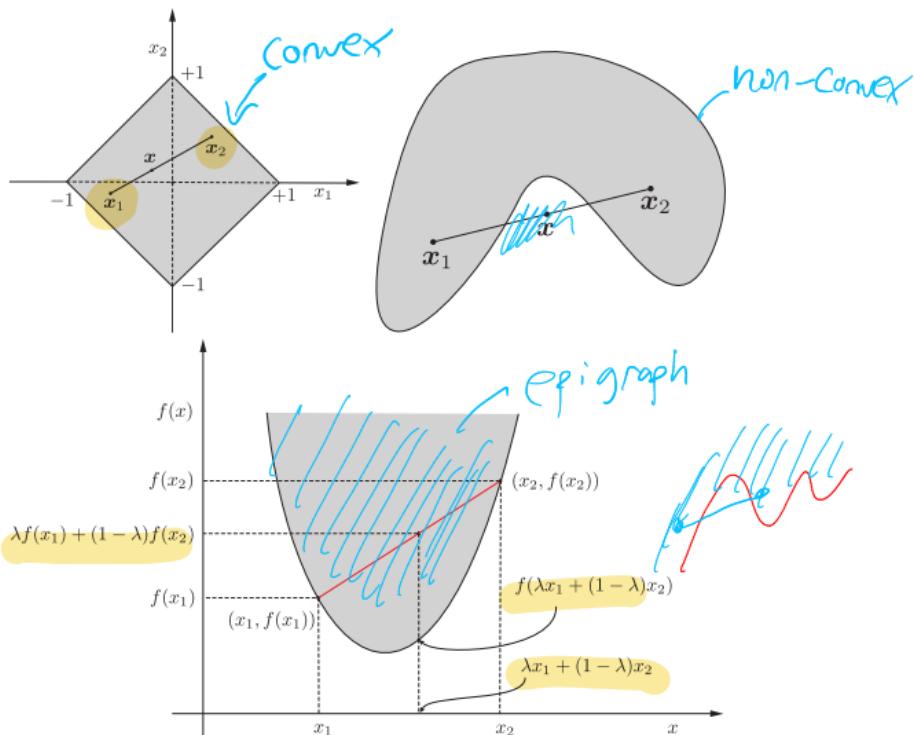
$|\mathcal{X}|$  is the cardinality (size) of the set  $\mathcal{X}$ .

How can these different norms promote sparsity?

**The  $\ell_p$  norm is convex ( $p \geq 1$ )**



IYL 8.2 3 pages



Can be proven by showing triangle inequality implies convexity.

- The  $\ell_p$  norm can in different configurations regularize our learning problem.
- $\ell_p$  is only a true norm for  $p \geq 1$ .
- $\ell_0$  is the most sparse "norm", and counts the number of nonzero elements.
- The  $\ell_1$  norm is the most sparse true norm, and is convex, hence we can optimize it.
- The  $\ell_1$  norm is not differentiable though. *need fixing*

For a rigorous treatment of norms: 01325 – Function spaces and mathematical analysis.

For a rigorous treatment of optimization: 02612 – Constrained Optimization.

# LASSO

## Exponentially weighted least-squares

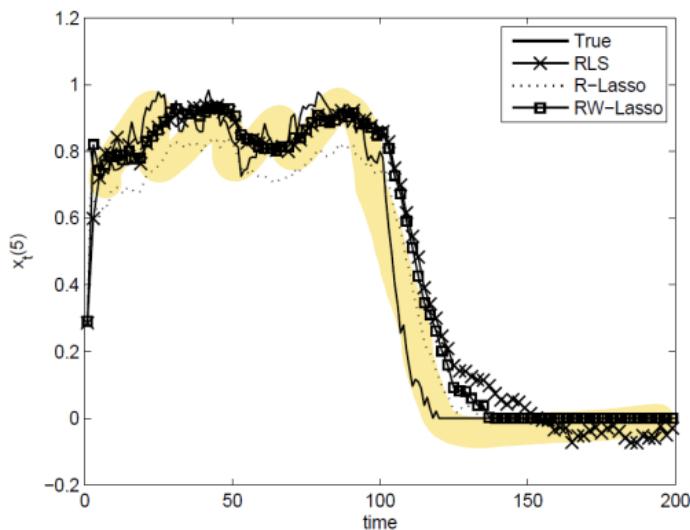
$$J(\boldsymbol{\theta}, \beta, \lambda) = \sum_{i=0}^n \beta^{n-i} (y_i - \boldsymbol{\theta}^T \mathbf{x}_i)^2 + \lambda \beta^{n+1} \|\boldsymbol{\theta}\|_2^2 \quad \text{RLS}$$

## Recursive LASSO cost function

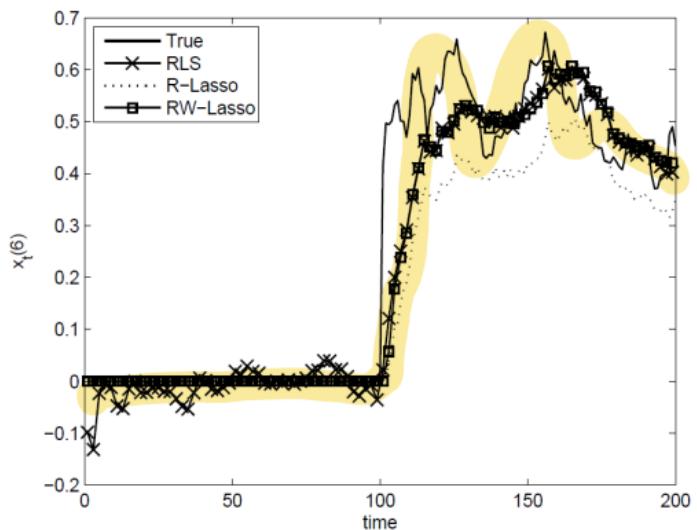
$$J(\boldsymbol{\theta}, \beta, \lambda) = \sum_{i=0}^n \beta^{n-i} (y_i - \boldsymbol{\theta}^T \mathbf{x}_i)^2 + \lambda \beta^{n+1} \|\boldsymbol{\theta}\|_1$$

Ref: D. Angelosante and G. B. Giannakis, "RLS-weighted Lasso for adaptive estimation of sparse signals," 2009 IEEE International Conference on Acoustics, Speech and Signal Processing, Taipei, 2009, pp. 3245-3248.

## RLS-weighted Lasso convergence



**Fig. 3.** Time evolution of the  $x_5$  signal entry.



**Fig. 4.** Time evolution of the  $x_6$  signal entry.

Ref: D. Angelosante and G. B. Giannakis, "RLS-weighted Lasso for adaptive estimation of sparse signals," 2009 IEEE International Conference on Acoustics, Speech and Signal Processing, Taipei, 2009, pp. 3245-3248.

# The Least Absolute Shrinkage and Selection Operator (LASSO)

We have the familiar regression task (on matrix form)

$$\underset{\text{Input}}{\mathbf{y}} = \underset{\text{Target}}{X\boldsymbol{\theta}} + \underset{\text{noise}}{\boldsymbol{\eta}}, \quad \mathbf{y} \in \mathbb{R}^N, X \in \mathbb{R}^{N \times l}, \boldsymbol{\theta} \in \mathbb{R}^l, \boldsymbol{\eta} \in \mathbb{R}^N,$$

## LASSO cost function

$$J(\boldsymbol{\theta}, \lambda) = \sum_{i=1}^N (y_i - \boldsymbol{\theta}^T \mathbf{x}_i)^2 + \lambda \|\boldsymbol{\theta}\|_1$$

Note:  $\sum_{i=1}^N a_i^2 = \mathbf{a} \cdot \mathbf{a} = \mathbf{a}^T \mathbf{a} = \|\mathbf{a}\|_2^2$

The LASSO cost function can equivalently be written as

$$J(\boldsymbol{\theta}, \lambda) = (\mathbf{y} - X\boldsymbol{\theta})^T (\mathbf{y} - X\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_1$$

**LASSO minimization,  $p = 1$** 

$$\hat{\boldsymbol{\theta}}_1 = \arg \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}, \lambda) = \sum_{i=1}^N (y_i - \boldsymbol{\theta}^T \mathbf{x}_i)^2 + \lambda \|\boldsymbol{\theta}\|_p^p$$

## Nomenclature

- $\hat{\boldsymbol{\theta}}_{LS}$  denotes the least squares solution  $\lambda = 0$
- $\hat{\boldsymbol{\theta}}_R$  denotes the least squares solution with  $\ell_2$  regularization (Ridge regression)  $p = 2$
- $\hat{\boldsymbol{\theta}}_1$  denotes the least squares solution with  $\ell_1$  regularization (LASSO)  $p = 1$

## Input–output of the weights when $X^T X = I$

### The weight estimates when $X^T X = I$

$$\hat{\theta}_{LS} = X^T y$$

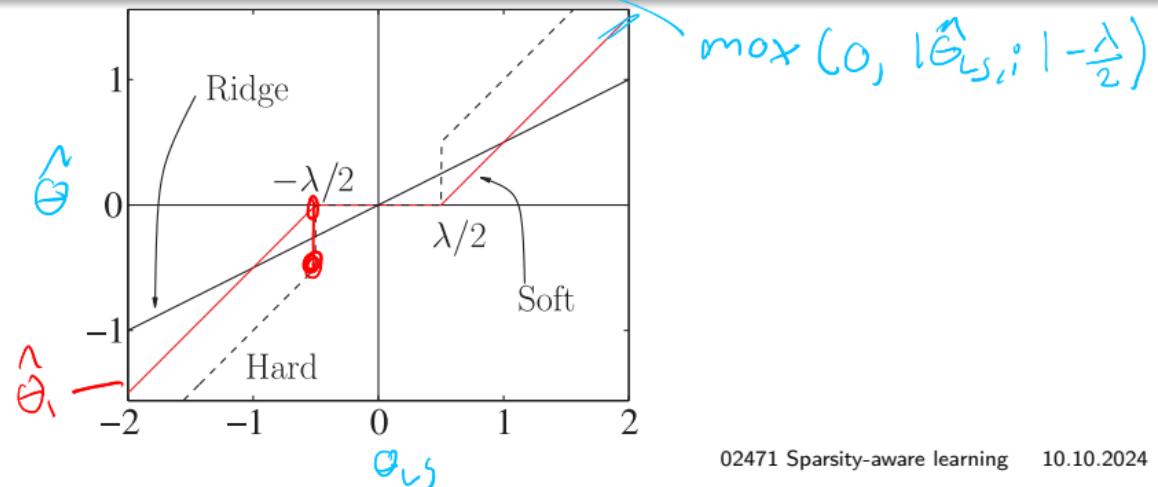
$$\hat{\theta}_R = \frac{1}{1+\lambda} \hat{\theta}_{LS}$$

$$\hat{\theta}_{1,i} = \text{sgn}(\hat{\theta}_{LS,i}) \left( |\hat{\theta}_{LS,i}| - \frac{\lambda}{2} \right)_+ \quad i = 1, 2, \dots, l$$

$$(X^T X)^{-1} X^T y = X^T y \cdot \frac{1}{1+\lambda}$$

$$(X^T X + \lambda I)^{-1} X^T y = (I + \lambda I)^{-1} \hat{\theta}_{LS}$$

ex6.2



## An example

Assume  $\lambda = 1$  and  $X^T X = I$ , then given a solution for  $\hat{\theta}_{LS}$ , compute the weights for the Ridge, Soft, and Hard thresholds. We have

$$\hat{\theta}_{LS} = [0.4, 0.6]^T$$

$$\hat{\theta}_R = \frac{1}{1 + \lambda} \hat{\theta}_{LS} = \frac{1}{2} [0.4 \ 0.6]^T = [0.2 \ 0.3]^T$$

$$\hat{\theta}_{1,i} = \text{sgn}(\hat{\theta}_{LS,i}) \left( |\hat{\theta}_{LS,i}| - \frac{\lambda}{2} \right)_+ \quad i = 1, 2, \dots, l$$

We get ...

Compute on your own

$$\hat{\theta}_1 = \begin{bmatrix} 0.4 - 0.5 \\ 0.6 - 0.5 \end{bmatrix}_+ = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

## LASSO minimization

$$\hat{\theta}_1 = \arg \min_{\theta} J(\theta, \lambda) = (\mathbf{y} - X\theta)^T(\mathbf{y} - X\theta) + \lambda \|\theta\|_1$$

The minimization problem can equivalently be written as

$$\hat{\theta}_1 = \arg \min_{\theta} (\mathbf{y} - X\theta)^T(\mathbf{y} - X\theta)$$

Subject to :  $\|\theta\|_1 \leq \rho$

$\left| \begin{array}{l} \text{LS} \\ 0 \leq \rho - \|\theta\|_1 \end{array} \right.$

For a specific choice of  $\lambda$  and  $\rho$ .

To see this, use Lagrangian multipliers (appendix C):

$$\text{minimize : } J(\theta)$$

$$\text{subject to : } f(\theta) \geq 0$$

$$\text{Lagrangian : } L(\theta, \lambda) = J(\theta) - \lambda f(\theta)$$

$$\begin{aligned}
 &= \text{LS} - \lambda(-\|\theta\|_1 + \rho) \\
 &= +\lambda\|\theta\|_1 - \lambda\rho
 \end{aligned}$$

## The LASSO solution

Take the derivative of the cost function, and put to zero, solve for  $\theta$ :  $\frac{d}{d\theta} J(\theta, \lambda) = 0$

Problem: the  $\ell_1$  norm is not differentiable.

Solution: use the subgradient method (sec 8.10).

### Subgradient Algorithm

$$\theta^{(i)} = \theta^{(i-1)} - \mu_i J'(\theta^{(i-1)}) \quad \text{GD}$$

$J'(\cdot)$  denotes any subgradient of  $J(\cdot)$ .

Converges if the stepsize  $\mu_i$  diminishes over time. Additionally

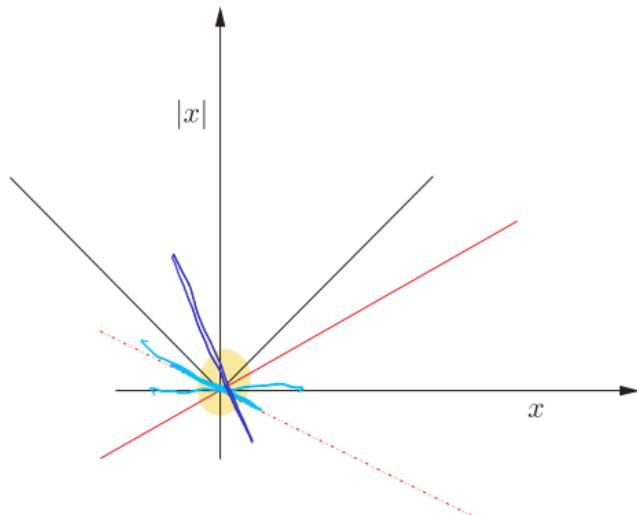
### Condition for a minimizer

Given a convex function  $f$ , the zero vector must belong to the subgradient at  $x_*$ , then  $x_*$  is a minimizer, i.e.

$$0 \in \partial f(x_*) \quad \checkmark$$

$\partial f(x_*)$  is the subdifferential of  $f(x)$  at  $x_*$  (the set of all subgradients).

$$\partial f(x) = \begin{cases} \text{sgn}(x), & \text{if } x \neq 0 \\ [-1, 1], & \text{if } x = 0 \end{cases}$$



Can formally be derived using definition of subdifferentials (example 8.4).

## LASSO

## Simplified analysis of weights

Assume  $X^T X = I$ , then

$$\hat{\theta}_{LS} = (X^T X)^{-1} X^T \mathbf{y} = X^T \mathbf{y}$$

$$\hat{\theta}_R = (X^T X + \lambda I)^{-1} X^T \mathbf{y} = \frac{1}{1 + \lambda} \hat{\theta}_{LS}$$

For  $\ell_1$ , we get

$$J(\boldsymbol{\theta}, \lambda) = (\mathbf{y} - X\boldsymbol{\theta})^T (\mathbf{y} - X\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_1 \quad \text{Lasso}$$

$$\partial_{\boldsymbol{\theta}} J(\boldsymbol{\theta}, \lambda) = -2X^T \mathbf{y} + 2X^T X \boldsymbol{\theta} + \lambda \partial \|\boldsymbol{\theta}\|_1$$

Which lead to

$$0 \in -\theta_{LS,i} + \hat{\theta}_{i,1} + \frac{\lambda}{2} \partial |\hat{\theta}_{i,1}|$$

Which again leads to (which you will derive in the exercise)

$$\hat{\theta}_{1,i} = \operatorname{sgn}(\hat{\theta}_{LS,i}) \left( |\hat{\theta}_{LS,i}| - \frac{\lambda}{2} \right)_+ \quad i = 1, 2, \dots, l$$

$\boldsymbol{\theta} > 0$   
 $\boldsymbol{\theta} < 0$   
 $\boldsymbol{\theta} = 0$

When is  $\ell_1$  unique? Example of recovery I

- Consider an unknown system.  $f(x)$
- The system only has two parameters, and you would like to recover those parameters.
- You know the system is on the form  $f(x_1, x_2) = x_1\theta_1 + x_2\theta_2$ .  $\theta_1$  and  $\theta_2$
- You can only measure once, that is, you can only select one pair of  $(x_1, x_2)$ .
- To simulate this environment, we create the  $\theta$ 's for the unknown system, and select those as  $(\theta_1, \theta_2) = (0, 1)$  (arbitrary choice).
- We will now see what happens if we select 3 possible measurements,  $x_a = (1/2, 1)$ ,  $x_b = (1, 1)$ ,  $x_c = (2, 1)$ .
- These three measurements will all give a response of  $f(x_1, x_2) = 1$  for the unknown system.

When is  $\ell_1$  unique? Example of recovery II

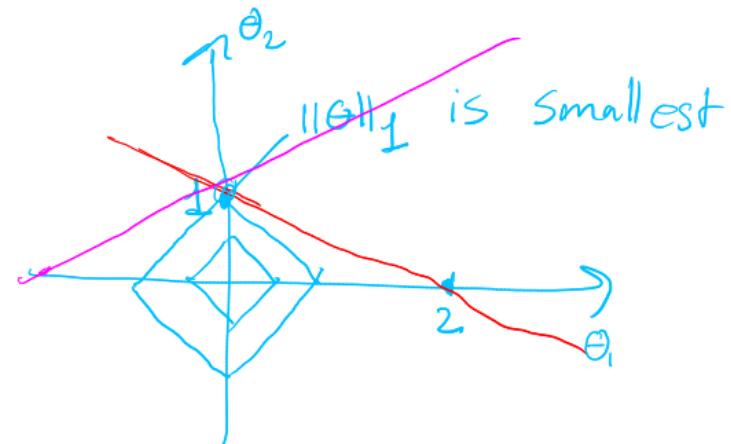
True  $\theta = (0, 1)$ . Find  $\theta$

Test measurement:  $x_a = (0.5, 1) \Rightarrow f(x) = 1$

System:  $f(x) = x_1\theta_1 + x_2\theta_2$  "box"

$$1 = 0.5 \cdot \theta_1 + 1 \cdot \theta_2$$

$$1 = 0.5 \theta_1 + 1 \cdot \theta_2$$

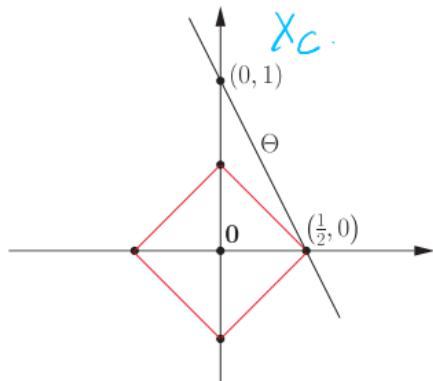
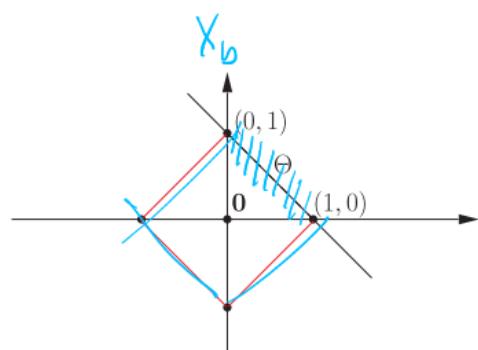
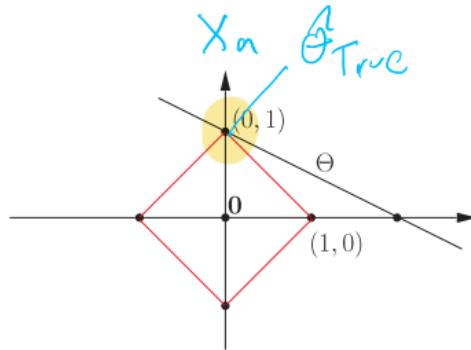


When is  $\ell_1$  unique? Example of recovery II

True  $\theta = (0, 1)$ .

Test measurements:  $x_a = (1/2, 1)$ ,  $x_b = (1, 1)$ ,  $x_c = (2, 1)$  (all result in  $f(x) = 1$ ).

Let us draw possible solutions to the linear system,  $f(x) = x_1\theta_1 + x_2\theta_2 = 1$ , and select the solution with the lowest  $\ell_1$  norm.



Compressed sensing explores the construction of  $X$ , and it turns out that a good  $X$  is a random  $X$  (sec 9.6–9.8).

- In the pursuit of sparse solution, we arrive at using the  $\ell_1$  norm as regularizer (LASSO) as the computationally most efficient norm.
  - The norm is convex.
  - No closed form solution, but solved with subgradients.
  - Has a soft thresholding operation, in the special case  $X^T X = I$ , sets weights to zero once they are numerically smaller than  $\lambda/2$ .
  - Almost always have a unique solution.
- The  $\ell_0$  "norm" leads to the sparsest solution, but is not convex.
- Under special circumstances, the  $\ell_1$  regularization will find the sparsest solution.
- We only have ONE solution (so far) for the LASSO when  $X^T X = I$ . Algorithms will be looked at next week.

Material: ML 10.1, 10.2–10.2.1 (until p.476), 10.2.2 (until p.482), 10.5–10.6.

- Estimating the LASSO solution:
  - LARS and Matching Pursuit algorithms
- Towards dictionary learning
  - Sparse analysis models
  - Time-frequency analysis

2022. pl.3