

Lecture 23: Introduction to Sparse Recovery

Recovery

Up to this point in the course we have seen many problems & models, all of which roughly are of the form

$$\textcircled{*} \quad \min_{\underline{\beta}} \frac{1}{2} \|\underline{A}\underline{\beta} - \underline{Y}\|_2^2 + \underbrace{\frac{1}{2} \|\underline{\beta}\|_2^2}_{\text{Regularization}},$$

ex supervised learning, Ridge regression, graphical SSL, classification, kernel regression broadly.

We also saw, through Cross-Validation, that the choice of $\lambda \geq 0$ has a profound impact on the quality of our solutions. This is super important when A has fewer rows than columns! i.e., the system $\underline{A}\underline{\beta} = \underline{Y}$ is underdetermined.

Observe that if we differentiate $\textcircled{*}$ wrt $\underline{\beta}$ we get.

$$A^T(A\hat{\beta} - Y) + \lambda \hat{\beta}$$

$$\Rightarrow (A^T A + \lambda I) \hat{\beta} = A^T Y.$$

if $A \in \mathbb{R}^{N \times M}$ with $N < M$ then $A^T A \in \mathbb{R}^{M \times M}$ is rank deficient & so not invertible! Then adding λI to $A^T A$ results in an invertible matrix & stabilizes the solution $\hat{\beta}$.

Hence, the term "regularization" since λI makes the problem more regular.

But one might wonder, what is special about $\|\beta\|_2$? why not another norm?

This is an excellent question & as it turns out it has profound implications. Particular choices of the norm, ex $\|\beta\|_p$ with $p \leq 1$ turns out to be very helpful. This is an advanced topic at the intersection of signal processing, convex analysis, applied math & statistics!

This line of thought also gave rise to the field of Compressed Sensing!

An Introduction To Compressive Sampling

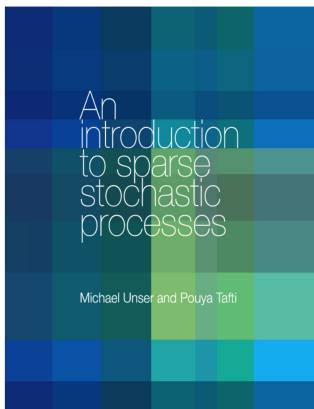
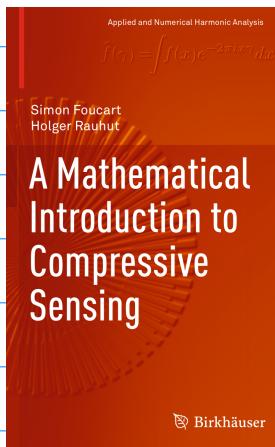
A sensing/sampling paradigm that goes against the common knowledge in data acquisition

Emmanuel J. Candès
and Michael B. Wakin

Conventional approaches to sampling signals or images follow Shannon's celebrated theorem: the sampling rate must be at least twice the maximum frequency present in the signal (the so-called Nyquist rate). In fact, this principle underlies nearly all signal acquisition protocols used in consumer audio and visual electronics, medical imaging devices, radio receivers, and so on. (For some signals, such as images, that are not naturally bandlimited, the sampling rate is often not set by the Shannon theorem but by the desired temporal or spatial resolution.) However, it is common in such systems to include a low-pass filter to bandlimit the signal before sampling, and so the Shannon theorem plays an implicit role. In the field of data compression, for example, standard analog-to-digital converter (ADC) technology implements the usual quantized Shannon representation: the signal is uniformly sampled at or above the Nyquist rate.

Digital Object Identifier 10.1109/SP.2007.414751

1053-5886/08/\$25.00/0008 IEEE
IEEE SIGNAL PROCESSING MAGAZINE | 21 | MARCH 2008

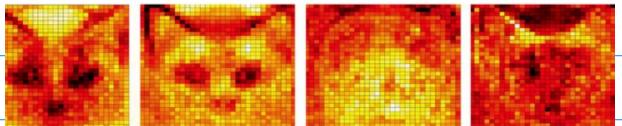


ch18

OXFORD

Data-Driven Modeling & Scientific Computation

Methods for Complex Systems & Big Data



J. NATHAN KUTZ



23.1 The Role of Sparsity

The notion of "sparsity" is going to be key for us.

Consider a linear system

$$A\beta = y, \quad \beta \in \mathbb{R}^N, \quad y \in \mathbb{R}^M, \quad A \in \mathbb{R}^{M \times N}$$

where $M < N$, so the system is underdetermined so it cannot be solved uniquely.

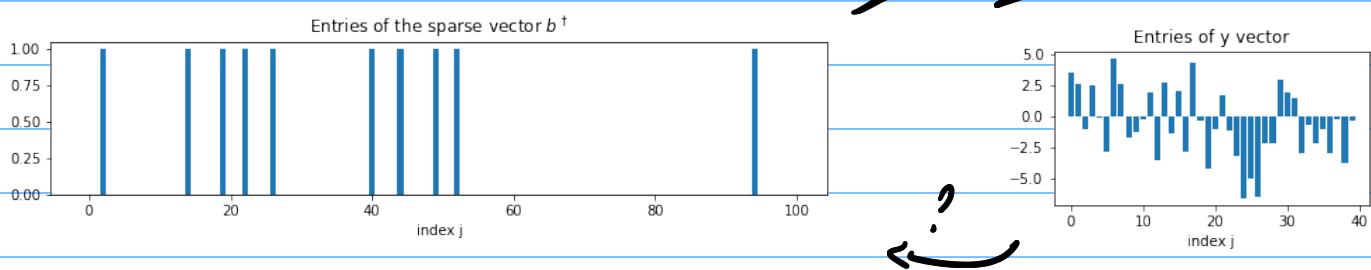
However, suppose that we have additional information that the vector β is sparse, more precisely it only has s -non-zero entries! (we say β is s -sparse)

Then, you might wonder whether you could find β if $s \ll m < n$?

Let us look at a quick example.

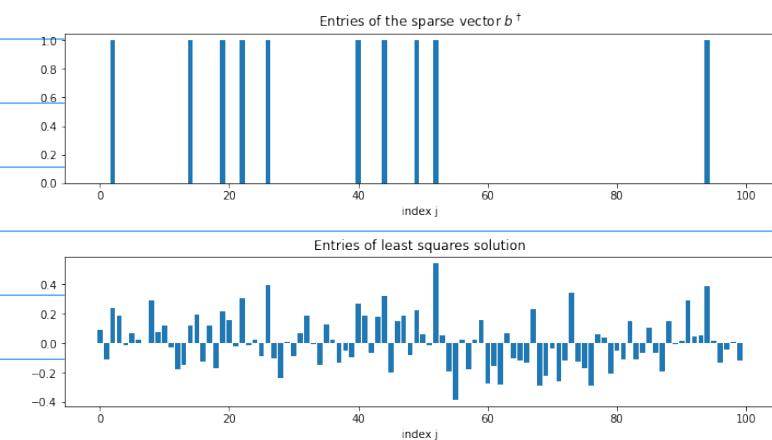
A vector $\beta^* \in \mathbb{R}^{100}$ with only $s=10$ non-zero entries & a matrix $A \in \mathbb{R}^{40 \times 100}$ with random entries $A_{ij} \stackrel{iid}{\sim} N(0, 1)$. Then set $y = A\beta^* \in \mathbb{R}^{40}$.

Our goal is then to infer β^* from y , i.e., "solve" $A\beta = y$.



First thing we might try is least squares but this does not work well!

$$\beta_{\text{lsq}} = \underset{\beta \in \mathbb{R}^{100}}{\arg \min} \|A\beta - y\|_2^2$$



We now try a different approach call Lasso which solves

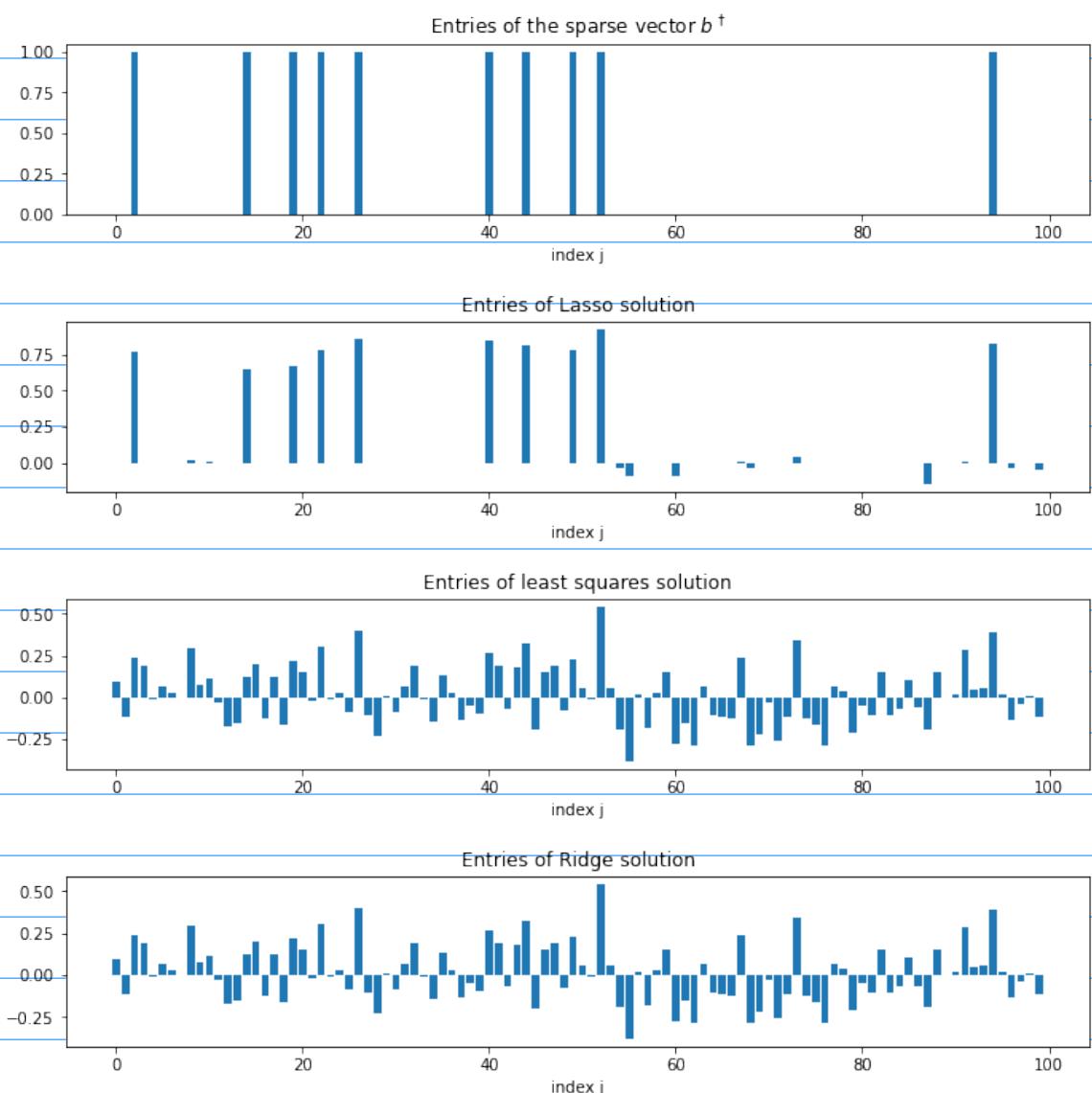
$$\beta_{\text{Lasso}} = \underset{\beta \in \mathbb{R}^{100}}{\arg \min} \frac{1}{2} \|A\beta - y\|_2^2 + \lambda \|\beta\|_1,$$

Compare to Ridge

$$\beta_{\text{Ridge}} = \underset{\beta \in \mathbb{R}^{100}}{\arg \min} \frac{1}{2} \|A\beta - y\|_2^2 + \lambda \|\beta\|_2^2.$$

essentially we use a different regularization term.

But this ends up having a huge effect!



At an intuitive level the reason for the success of Lasso is that the $\|\cdot\|_1$ -norm "promotes/prefers" solutions that are sparse while the $\|\cdot\|_2$ -norm does not. Hence, Lasso / 1-norm regularization are preferable when we expect \underline{x} to be sparse!

23.2 Intuition for Lasso

Solving a problem of the form

$$\textcircled{I} \quad \underline{\beta}^* = \underset{\underline{\beta} \in \mathbb{R}^N}{\operatorname{argmin}} \frac{1}{2} \|\underline{A}\underline{\beta} - \underline{y}\|_2^2 + \lambda \|\underline{\beta}\|_p^p$$

is equivalent to solving

$$\textcircled{II} \quad \underline{\beta}^* = \underset{\underline{\beta} \in \mathbb{R}^N}{\operatorname{argmin}} \frac{1}{2} \|\underline{A}\underline{\beta} - \underline{y}\|_2^2$$

$$\text{subject to } \|\underline{\beta}\|_p^p \leq t,$$

The values of t & λ are related to each other, indeed, \textcircled{I} is the Lagrangian form of \textcircled{II} .

Here we can see the effect of the choice of the norm $(\|\cdot\|_p)$ in the regularization term.

