

IN550 Machine Learning

Decomposizione sparsa di segnali

Vincenzo Bonifaci

Minimizzazione della norma ℓ_1

$$\begin{aligned} &\text{minimize } |x_1| + |x_2| + \dots + |x_m| && \text{(L1)} \\ &\text{s.t. } Ax = b \\ & && x \in \mathbb{R}^m \end{aligned}$$

per $A \in \mathbb{R}^{n \times m}$, $b \in \mathbb{R}^n$

$$\begin{aligned} &\text{minimize } |x_1| + |x_2| + \dots + |x_m| && \text{(L1)} \\ &\text{s.t. } Ax = b \\ &\quad x \in \mathbb{R}^m \end{aligned}$$

per $A \in \mathbb{R}^{n \times m}$, $b \in \mathbb{R}^n$

Problema comune nell'elaborazione dei segnali e in statistica

Minimizzazione della norma ℓ_1

$$\begin{aligned} &\text{minimize } |x_1| + |x_2| + \dots + |x_m| && \text{(L1)} \\ &\text{s.t. } Ax = b \\ &\quad x \in \mathbb{R}^m \end{aligned}$$

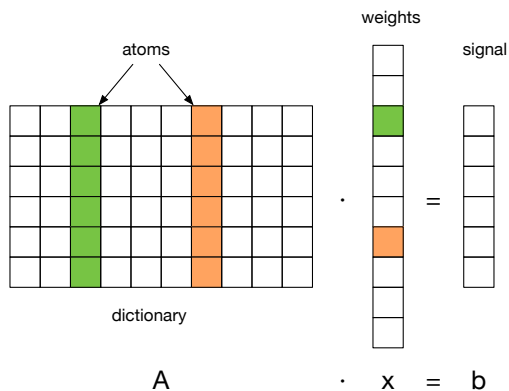
per $A \in \mathbb{R}^{n \times m}$, $b \in \mathbb{R}^n$

Problema comune nell'elaborazione dei segnali e in statistica

Assunzione standard: A ha **rango pieno**, pari ad n

Decomposizione di segnali

Dati $A \in \mathbb{R}^{n \times m}$, $b \in \mathbb{R}^n$: risolvere $Ax = b$ for $x \in \mathbb{R}^m$



Tipicamente $m > n$: il sistema è **sottodeterminato**

Regolarizzazione

Per selezionare una particolare soluzione: $\min_x J(x)$ s.t. $Ax = b$

$J(\cdot)$ è un criterio di **regolarizzazione**

Una scelta comune: $J(x) = \frac{1}{2} \|x\|_2^2$

Perché?

Regolarizzazione

Per selezionare una particolare soluzione: $\min_x J(x)$ s.t. $Ax = b$

$J(\cdot)$ è un criterio di **regolarizzazione**

Una scelta comune: $J(x) = \frac{1}{2} \|x\|_2^2$

Perché?

- Ha soluzione in **forma chiusa**: $\hat{x} = A^\top (AA^\top)^{-1} b$
- \hat{x} può essere calcolato risolvendo **1 sistema lineare** di equazioni

Ma non **sempre** è la scelta migliore

Regolarizzazione ℓ_p

$$\text{Funzione } \ell_p: \|x\|_p = \sqrt[p]{\sum_i |x_i|^p}$$

Regolarizzazione ℓ_p

Funzione ℓ_p : $\|x\|_p = \sqrt[p]{\sum_i |x_i|^p}$

Notazione:

Regolarizzazione ℓ_p

$$\text{Funzione } \ell_p: \|x\|_p = \sqrt[p]{\sum_i |x_i|^p}$$

Notazione:

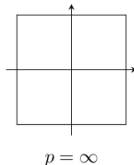
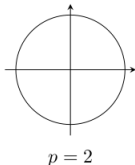
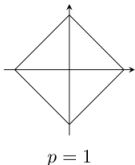
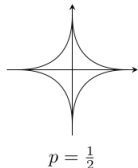
- $\|x\|_\infty \stackrel{\text{def}}{=} \max_i |x_i|$ (caso limite per $p \rightarrow \infty$)

Regolarizzazione ℓ_p

$$\text{Funzione } \ell_p: \|x\|_p = \sqrt[p]{\sum_i |x_i|^p}$$

Notazione:

- $\|x\|_\infty \stackrel{\text{def}}{=} \max_i |x_i|$ (caso limite per $p \rightarrow \infty$)
- $\|x\|_0 \stackrel{\text{def}}{=} \text{numero di componenti non-nulle di } x$ (**attenzione**: non è una norma, nonostante la notazione!)

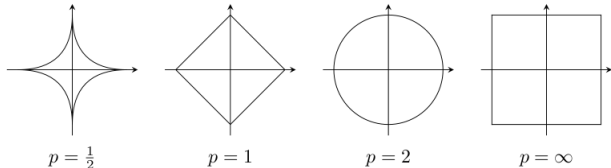


Regolarizzazione ℓ_p

$$\text{Funzione } \ell_p: \|x\|_p = \sqrt[p]{\sum_i |x_i|^p}$$

Notazione:

- $\|x\|_\infty \stackrel{\text{def}}{=} \max_i |x_i|$ (caso limite per $p \rightarrow \infty$)
- $\|x\|_0 \stackrel{\text{def}}{=} \text{numero di componenti non-nulle di } x$ (**attenzione**: non è una norma, nonostante la notazione!)

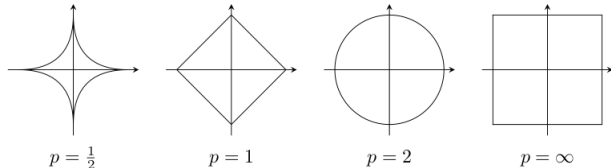


Regolarizzazione ℓ_p

Funzione ℓ_p : $\|x\|_p = \sqrt[p]{\sum_i |x_i|^p}$

Notazione:

- $\|x\|_\infty \stackrel{\text{def}}{=} \max_i |x_i|$ (caso limite per $p \rightarrow \infty$)
- $\|x\|_0 \stackrel{\text{def}}{=} \text{numero di componenti non-nulle di } x$ (**attenzione**: non è una norma, nonostante la notazione!)



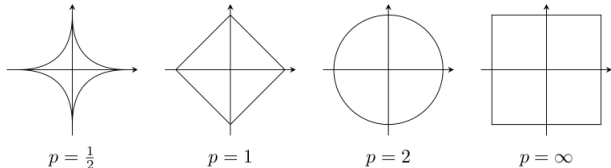
- ℓ_p è una **norma** sse $p \geq 1$

Regolarizzazione ℓ_p

$$\text{Funzione } \ell_p: \|x\|_p = \sqrt[p]{\sum_i |x_i|^p}$$

Notazione:

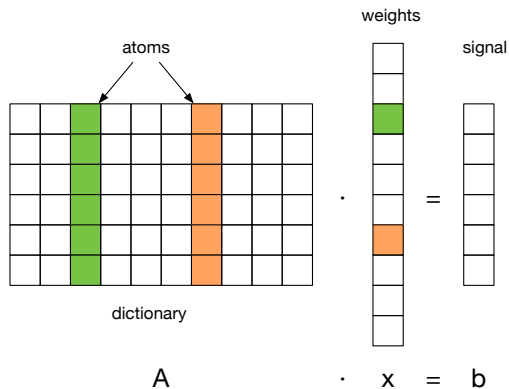
- $\|x\|_\infty \stackrel{\text{def}}{=} \max_i |x_i|$ (caso limite per $p \rightarrow \infty$)
- $\|x\|_0 \stackrel{\text{def}}{=} \text{numero di componenti non-nulle di } x$ (**attenzione**: non è una norma, nonostante la notazione!)



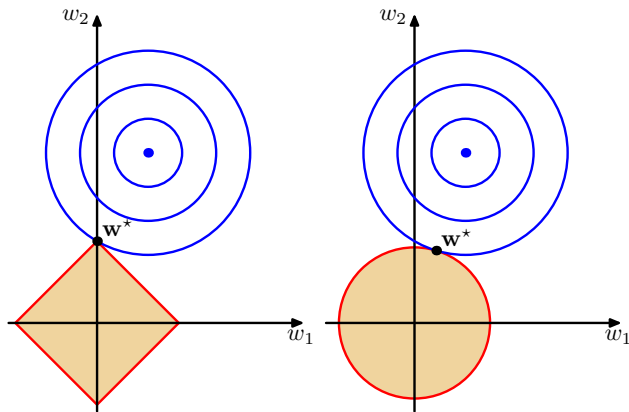
- ℓ_p è una **norma** sse $p \geq 1$
- ℓ_p è **convessa** sse $p \geq 1$

Decomposizione di segnali in norma ℓ_p

$$\min \|x\|_p \text{ s.t. } Ax = b$$



Effetto di usare ℓ_1 vs. ℓ_2



La funzione obiettivo $\|x\|_1$ promuove decomposizioni **sparse**

ℓ_1 promuove decomposizioni sparse

$$\begin{aligned} & \text{minimize} && \|x\|_1 && \text{(L1)} \\ & \text{s.t.} && Ax = b \\ & && x \in \mathbb{R}^m \end{aligned}$$

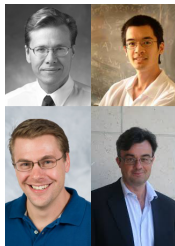
Theorem

(L1) ammette un minimizzante con al più n componenti non-nulle.

Cio è interessante perché spesso $n \ll m$

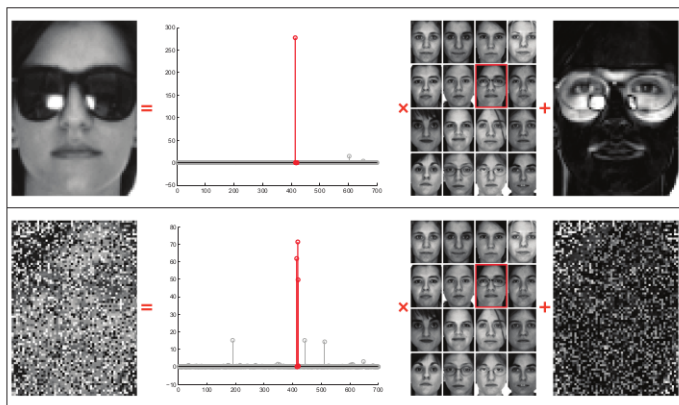
Tipicamente, invece, i minimizzanti in norma ℓ_2 sono **densi**

D. Donoho, T. Tao, J.K. Romberg, E. Candes, ed altri hanno dimostrato che un segnale K -sparso può essere ricostruito (in modo robusto) con $O(K)$ osservazioni



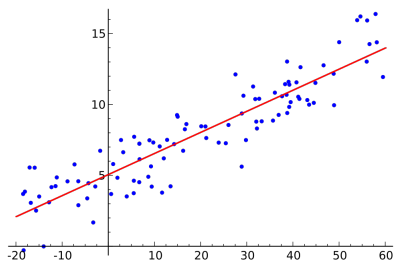
Face Recognition Breakthrough

By using sparse representation and compressed sensing, researchers have been able to demonstrate significant improvements in accuracy over traditional face-recognition techniques.



Problemi di regressione lineare

Dati di input $D \in \mathbb{R}^{m \times n}$, dati di output $y \in \mathbb{R}^m$ ($m \geq n$)



i	D_i	y_i
1	-14.0	0.01
2	2.5	7.52
...
m	58.2	11.9

Fit di una funzione lineare: $\min_{w \in \mathbb{R}^n} \|y - Dw\|$

Cerchiamo soluzioni **approssimate** del sistema **sovradeterminato** $Dw = y$

Quale norma per i residui?

I **residui** sono le differenze $y - Dw$

- **Metodo dei minimi quadrati**: norma ℓ_2
- **Deviazioni assolute minime** (Least Absolute Deviations o LAD):
norma ℓ_1

Quale norma per i residui?

I **residui** sono le differenze $y - Dw$

- **Metodo dei minimi quadrati**: norma ℓ_2
- **Deviazioni assolute minime** (Least Absolute Deviations o LAD):
norma ℓ_1

minimize $(z - z_1)^2 + (z - z_2)^2 + (z - z_3)^2$
 $\Rightarrow z =$ **media** degli z_i

Quale norma per i residui?

I **residui** sono le differenze $y - Dw$

- **Metodo dei minimi quadrati**: norma ℓ_2
- **Deviazioni assolute minime** (Least Absolute Deviations o LAD): norma ℓ_1

$$\text{minimize } (z - z_1)^2 + (z - z_2)^2 + (z - z_3)^2$$

$\Rightarrow z = \text{media}$ degli z_i

$$\text{minimize } |z - z_1| + |z - z_2| + |z - z_3|$$

$\Rightarrow z = \text{mediana}$ degli z_i

$\Rightarrow \ell_1$ può essere più robusta rispetto ad esempi anomali (**outlier**)

Breve digressione storica

- Metodo dei minimi quadrati:

- Metodo dei minimi quadrati: Adrien-Marie Legendre, 1805

Breve digressione storica

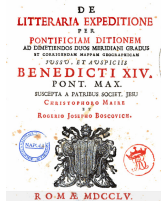
- Metodo dei minimi quadrati: Adrien-Marie Legendre, 1805
- Deviazioni assolute minime:

Breve digressione storica

- Metodo dei minimi quadrati: Adrien-Marie Legendre, 1805
- Deviazioni assolute minime: Ruggero Boscovich, 1760

Breve digressione storica

- Metodo dei minimi quadrati: Adrien-Marie Legendre, 1805
- Deviazioni assolute minime: Ruggero Boscovich, 1760



Boscovich, R. J. (1760). *De recentissimis graduum dimensionibus, et figura, ac magnitudine terrae inde derivanda.*

Regressione LAD è equivalente al problema (L1)

Il problema di deviazioni assolute minime

$$\min_{w \in \mathbb{R}^n} \|y - Dw\|_1$$

è **equivalente** a

$$\min_{x \in \mathbb{R}^m} \|x\|_1 \quad \text{s.t. } Ax = b$$

per opportune A, b

Dimostrazione

Si prendano A, b tali che $\ker A = \text{im } D$ e $b = Ay$, allora

Regressione LAD è equivalente al problema (L1)

Il problema di deviazioni assolute minime

$$\min_{w \in \mathbb{R}^n} \|y - Dw\|_1$$

è **equivalente** a

$$\min_{x \in \mathbb{R}^m} \|x\|_1 \quad \text{s.t.} \quad Ax = b$$

per opportune A, b

Dimostrazione

Si prendano A, b tali che $\ker A = \text{im } D$ e $b = Ay$, allora

$$Ax = b$$

Regressione LAD è equivalente al problema (L1)

Il problema di deviazioni assolute minime

$$\min_{w \in \mathbb{R}^n} \|y - Dw\|_1$$

è **equivalente** a

$$\min_{x \in \mathbb{R}^m} \|x\|_1 \quad \text{s.t.} \quad Ax = b$$

per opportune A, b

Dimostrazione

Si prendano A, b tali che $\ker A = \text{im } D$ e $b = Ay$, allora

$$Ax = Ay$$

Regressione LAD è equivalente al problema (L1)

Il problema di deviazioni assolute minime

$$\min_{w \in \mathbb{R}^n} \|y - Dw\|_1$$

è **equivalente** a

$$\min_{x \in \mathbb{R}^m} \|x\|_1 \quad \text{s.t.} \quad Ax = b$$

per opportune A, b

Dimostrazione

Si prendano A, b tali che $\ker A = \text{im } D$ e $b = Ay$, allora

$$Ax = Ay \Leftrightarrow y - x \in \ker A$$

Regressione LAD è equivalente al problema (L1)

Il problema di deviazioni assolute minime

$$\min_{w \in \mathbb{R}^n} \|y - Dw\|_1$$

è **equivalente** a

$$\min_{x \in \mathbb{R}^m} \|x\|_1 \quad \text{s.t.} \quad Ax = b$$

per opportune A, b

Dimostrazione

Si prendano A, b tali che $\ker A = \text{im } D$ e $b = Ay$, allora

$$Ax = Ay \Leftrightarrow y - x \in \text{im } D$$

Regressione LAD è equivalente al problema (L1)

Il problema di deviazioni assolute minime

$$\min_{w \in \mathbb{R}^n} \|y - Dw\|_1$$

è **equivalente** a

$$\min_{x \in \mathbb{R}^m} \|x\|_1 \text{ s.t. } Ax = b$$

per opportune A, b

Dimostrazione

Si prendano A, b tali che $\ker A = \text{im } D$ e $b = Ay$, allora

$$Ax = Ay \Leftrightarrow y - x \in \text{im } D \Leftrightarrow x = y - Dw \text{ per qualche } w \in \mathbb{R}^n$$

(L1) è formulabile come problema di ottimizzazione lineare

(L1)

$$\begin{aligned} \min \quad & \|x\|_1 \\ \text{s.t.} \quad & Ax = b \\ & x \in \mathbb{R}^m \end{aligned}$$

(LP)

$$\begin{aligned} \min \quad & \mathbf{1}^\top z \\ \text{s.t.} \quad & z + x \geq 0 \\ & z - x \geq 0 \\ & Ax = b \\ & x, z \in \mathbb{R}^m \end{aligned}$$

m variabili, n vincoli $\Rightarrow 2m$ variabili, $n + 2m$ vincoli

(L1) è formulabile come problema di ottimizzazione lineare

(L1)

$$\begin{aligned} \min \quad & \|x\|_1 \\ \text{s.t.} \quad & Ax = b \\ & x \in \mathbb{R}^m \end{aligned}$$

(LP)

$$\begin{aligned} \min \quad & \mathbf{1}^\top z \\ \text{s.t.} \quad & z + x \geq 0 \\ & z - x \geq 0 \\ & Ax = b \\ & x, z \in \mathbb{R}^m \end{aligned}$$

m variabili, n vincoli $\Rightarrow 2m$ variabili, $n + 2m$ vincoli

Metodi di tipo Iteratively Reweighted Least Squares (IRLS)

Riscriviamo $\|x\|_1 = \sum_{i=1}^m \frac{1}{|x_i|} |x_i|^2 = x^\top Wx$

Schema tipo IRLS

- 1 Trova $\min x^\top Wx$ tale che $Ax = b$
- 2 Aggiorna la matrice dei pesi W (in qualche modo opportuno)
- 3 Ripeti fino a convergere ad un punto fisso

Il passo 1 è un problema ai **minimi quadrati** (\Rightarrow soluzione in forma chiusa)

I vari metodi differiscono nel passo 2

Vantaggi dei metodi di tipo IRLS

- Semplici da programmare
- **Ottimizzatore lineare non richiesto**, solo solutore di equazioni lineari
- In genere numericamente stabili

R. Chartrand, W. Yin. *Iteratively reweighted algorithms for compressive sensing*. IEEE Conf. on Acoustics, Speech and Signal Processing, 2008.

1000+ citazioni

Ma la complessità dei metodi di tipo IRLS non è ben compresa
Pochi risultati teorici