

CSE 6740 A/ISyE 6740: Computational Data Analysis: Introductory lecture

Nisha Chandramoorthy

August 31, 2023

Last time

- ▶ Least squares regression

Last time

- ▶ Least squares regression
- ▶ Gauss-Markov theorem

Last time

- ▶ Least squares regression
- ▶ Gauss-Markov theorem
- ▶ Ridge regression, optimization and geometric perspectives

Last time

- ▶ Least squares regression
- ▶ Gauss-Markov theorem
- ▶ Ridge regression, optimization and geometric perspectives
- ▶ Shrinkage

Today

- ▶ Shrinkage by ridge regression

Today

- ▶ Shrinkage by ridge regression
- ▶ Geometric view of LASSO

Today

- ▶ Shrinkage by ridge regression
- ▶ Geometric view of LASSO
- ▶ Generalization of LASSO

Today

- ▶ Shrinkage by ridge regression
- ▶ Geometric view of LASSO
- ▶ Generalization of LASSO
- ▶ ℓ^0 regularization and compressed sensing

Today

- ▶ Shrinkage by ridge regression
- ▶ Geometric view of LASSO
- ▶ Generalization of LASSO
- ▶ ℓ^0 regularization and compressed sensing

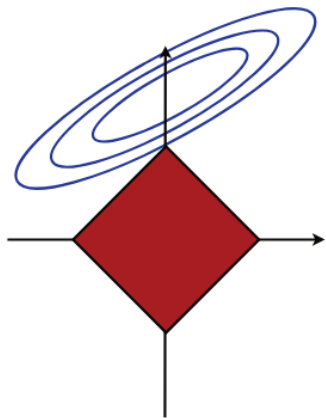
- ▶ Online linear regression algorithms

- ▶ Online linear regression algorithms
- ▶ Classification ERM

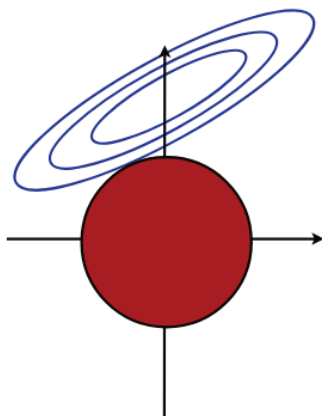
- ▶ Online linear regression algorithms
- ▶ Classification ERM
- ▶ Perceptron algorithm

- ▶ Online linear regression algorithms
- ▶ Classification ERM
- ▶ Perceptron algorithm
- ▶ Convergence proof of perceptron algorithm

Compression by LASSO



L1 regularization



L2 regularization

Shrinkage (Ridge)

Algorithms

$$\mathcal{H} := \{ \omega^T x + b : \omega \in \mathbb{R}^d, b \in \mathbb{R} \}$$

$$z = (x, y)$$

$$l(z, h) = (\omega^T x + b - y)^2$$

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m (\omega^T x_i + b - y_i)^2$$

$$(w, b) \quad (x_i, y_i) \in S$$

$$\omega^* = \underset{\omega \in \mathbb{R}^d}{\operatorname{argmin}} \hat{R}_S^{\text{ols}}(h) + \lambda \|\omega\|^2$$

$$(w^*, b^*)$$

$$\min_{\omega} \|X\omega - Y\|^2 + \lambda \|\omega\|^2$$

$$\omega^{\text{Ridge}} = (X^T X + \lambda I)^{-1} X^T Y$$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

$m \times d$

$$U = [u_1 | u_2 | \dots | u_d]$$

$$V = [v_1 | v_2 | \dots | v_d]$$

$$\Sigma = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_d \end{bmatrix}$$

$$Y^{\text{ridge}}_{\text{pred}} =$$

$$X \omega^{\text{ridge}} = X (X^T X + \lambda I)^{-1} X^T Y$$

$$= U \Sigma V^T (V \Sigma V^T + \lambda I)^{-1} V \Sigma U^T Y$$

$$V^T = V^{-1}$$

$$A^{-1} B^{-1} = (BA)^{-1}$$

$$= \sum_{i=1}^d u_i \frac{\sigma_i^2}{\lambda + \sigma_i^2} u_i^T Y$$

Shrinkage \uparrow when $\sigma_i \downarrow$.



Convex quadratic program

\rightarrow QP

\rightarrow iterative methods

$$\hat{R}_S(\omega) = \frac{1}{m} \|X\omega - Y\|^2 + \lambda \|\omega\|^2$$

Gradient descent

$$\omega^{(t+1)} = \omega^{(t)} - \underset{\substack{\uparrow \\ \text{learning rate}}}{\eta} \nabla \hat{R}_S(\omega^{(t)})$$

Stochastic gradient descent (SGD)

$$\omega^{(t+1)} = \omega^{(t)} - \eta \nabla \tilde{R}_S(\omega^{(t)})$$

Batch $b < m$ samples

Best subset selection

$$\hat{R}_s(w) = \|Xw - Y\|^2 + \lambda \|w\|_0 \leftarrow$$

$\|w\|_0$: l^0 norm Tibshirani 2000s

"
of nonzero entries of w .

$$\begin{bmatrix} \overline{x_1^T \dots} \\ \vdots \\ \overline{x_m^T \dots} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

short wide
X
 $m \ll d$

$$\|Y - Xw\|$$

Optimization way of seeing compression in the LASSO

in the case of X orthonormal

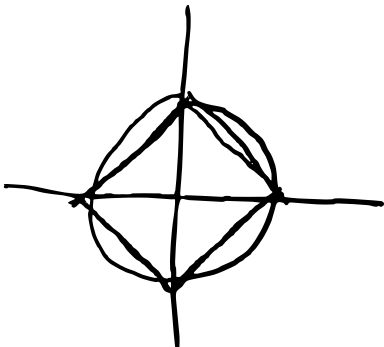
$$\hat{R}_S(\omega) = \|X\omega - Y\|^2 + \lambda \|\omega\|_1$$

$$\|\omega\|_1 = \sum_{j=1}^d |\omega_j|$$

LASSO estimator

ERM objective

$$\hat{R}_S^{\text{lasso}}, \quad \hat{R}_S^{\text{bss}}, \quad \hat{R}_S^{\text{ridge}}$$

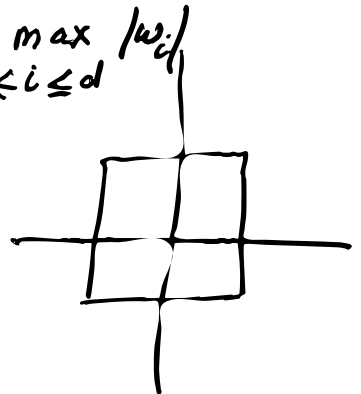


$$\|\omega\|_\infty = \gamma \quad \max_{1 \leq i \leq d} |\omega_i|$$

$$\|\omega\| = \gamma$$

$$\|\omega\|_1 = \gamma$$

$$|\omega_1| + |\omega_2|$$



$$\rightarrow \begin{array}{c} \downarrow \\ X \end{array} \omega \approx \begin{array}{c} \downarrow \\ Y \end{array}$$

$m \ll d$

Candes, Romberg, Tao 2005

$$\rightarrow X_{\omega} = Y + \varepsilon$$

Signal processing, image processing

$$S_j = \sum_{i=1}^d \underset{\uparrow}{x_i} \underset{\uparrow}{\omega_i}$$

Stable recovery of sparse w .

[Candes, Romberg, Tao 2005]

[Candes, Tao 2004]

Thm: $\arg \min_w \|w\|$, s.t. $Xw = Y$ \leftarrow

is the exact solution of $Xw = Y$
for any true sparse w if

$\text{nnz}(w) < S$, where X satisfies S -restricted isometry $\rightarrow \|w\|_0$

$$T \subseteq \{1, 2, 3, \dots, d\}$$

That is,

for all subsets of indices

with $|T| < S$, there is $\delta_S > 0$ s.t

$$(i) \quad (1 - \delta_S) \|v\|^2 < \underbrace{\|X_T v\|^2}_{T} \leq (1 + \delta_S) \|v\|^2$$

$$(ii) \quad \delta_S + \delta_{2S} + \delta_{3S} < 1 \leftarrow$$

$$X = [\dots | \dots]$$

X : Gaussian

Fourier basis

Stable recovery

Candes Romberg Tao 2005

Let $w^{\text{lasso}} :=$

$$\arg \min \|w\|_1 \quad \text{s.t.} \quad \|Xw - Y\| \leq \varepsilon$$

Let true w be sparse with $\text{nz}(w) \leq S$
with S s.t. $\underbrace{\delta_{3S} + 3\delta_{4S}}_{S\text{-isometry}} < 2$. $\|w\|_0$

Then,

$$\|w^{\text{lasso}} - w\| \leq C_S \varepsilon$$

$$w^{\text{ols}} = (X^T X)^{-1} X^T Y$$

$$\begin{aligned} w^{\text{ols}} - w &= (X^T X)^{-1} X^T Y - (X^T X)^{-1} X^T X w \\ &= (X^T X)^{-1} X^T (Y - X w) \end{aligned}$$

$$\|w^{\text{ols}} - w\| \leq C \|X^T \varepsilon\| \approx C \varepsilon$$

Compressed sensing

$\xrightarrow{m \leq d} (X^T X)^{\dagger}$ is computationally expensive
 $O(d^3)$

\rightarrow Interpretability

Generalization of regression

→ Hoeffding's inequality: $S_n = X_1 + X_2 + \dots + X_n$
 $X_i \perp\!\!\!\perp X_j$ $0 \leq X_i \leq L$

$$\mathbb{P}(S_n - \mathbb{E}S_n \geq t) \leq e^{-\frac{2t^2}{nL^2}}$$

Thm: $\sup_{z, h} l(z, h) = L$. Let \mathcal{H} be finite.

Then, for every $\delta > 0$, with probability at least $\underline{1-\delta}$, $\forall h \in \mathcal{H}$,

$$R(h) \leq \hat{R}_S(h) + L \sqrt{\frac{\log|\mathcal{H}| + \log 1/\delta}{2m}}$$

Proof: Use Hoeffding's inequality

Generalization for Ridge & Lasso

\mathcal{H} Complexity

Ridge

$$\underline{R(h)} \leq \hat{R}_S(h) + 4L \sqrt{\frac{\underline{\kappa}^2 \Lambda^2}{m}} + L^2 \sqrt{\frac{\log 1/\delta}{2m}}$$

where $\underline{\kappa}^2 \geq \frac{\Phi^T(x) \Phi(x)}{\|x\|^2} \forall x$
 $\|w\| < \Lambda$

Lasso $\mathcal{H} := \{h_{w,b}(x) : w \in \mathbb{R}^d, b \in \mathbb{R}, \|w\|_1 < \Lambda\}$

$$R(h) \leq \hat{R}_S(h) + 2\underline{\kappa}_\infty \Lambda_1 L \sqrt{\frac{2 \log(2d)}{m}} + L^2 \sqrt{\frac{\log 1/\delta}{2m}}$$

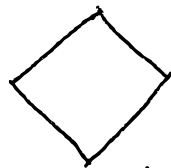
where $\|X_i\|_\infty \leq \kappa_\infty$,

$$\|w\|_1 \leq \Lambda_1$$

and $|h(x) - y| \leq L$

$$\|XW - Y\|^2$$

$$X^T X = Id$$



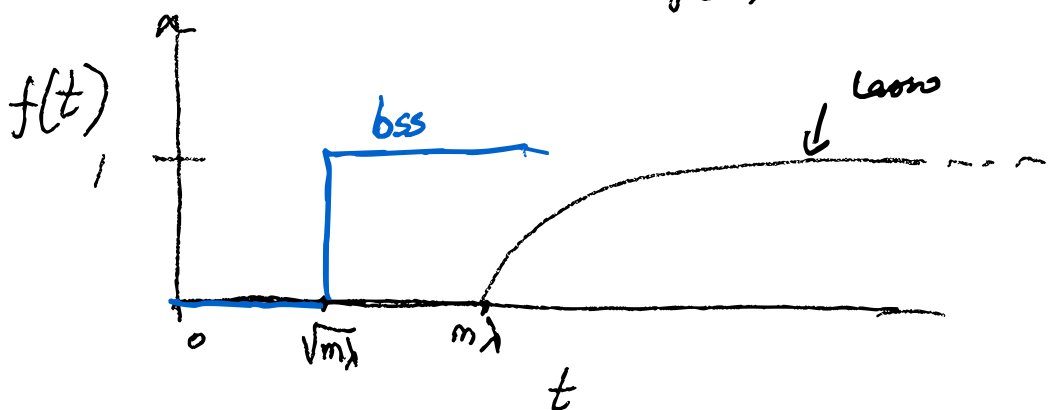
Subgradient
Proximal-descent

$$1) \quad w^{OLS} = X^T Y$$

$$2) \quad w^{ridge} = (X^T X + \lambda)^{-1} X^T Y \\ = (I + \lambda)^{-1} X^T Y$$

$$3) \quad w_j^{lasso} = w_j^{ols} \max(0, 1 - \frac{m\lambda}{w_j^{ols}})$$

$f(t)$



$$4) \quad w_j^{bss} = w_j^{ols} \frac{1}{\{w_j^{ols} > \sqrt{m\lambda}\}}$$