

Lasso (L₁-Regularized Regression)

Note Title

3/28/2016

Input $P \in \mathbb{R}^{d \times 1}$ n points

P_1, P_2, \dots, P_n

$$P_i = (P_{i1}, P_{i2}, \dots, P_{id-1}, \underbrace{P_{id}}_{y_i})$$



only important unit in d -dim-coord = y

$$A = [a_0, a_1, a_2, \dots, a_d] \quad P_x \quad a_0 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

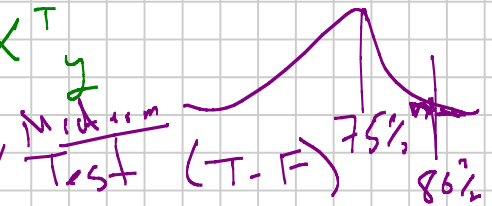
$$\sum_{P_i \in P} (f_A(P_i) - y_i)^2$$

$$X = \begin{bmatrix} 1 & P_x \end{bmatrix}$$

$$\min_A \|XA - y\|_2 \Rightarrow A = (X^T X)^{-1} X^T y$$

Regularization

- Bayesian Prior
- "Regression to the mean"



- Know 50 Qs
- Guess 50 Qs

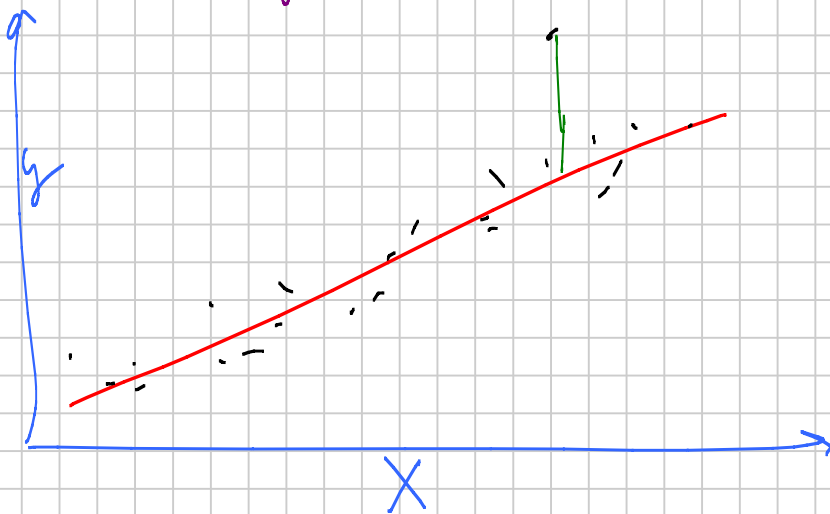
$$E[\text{score}] = 75\%$$

Take top 10 student
 $Ave[\text{score}] = 80\%$

Final

Expected of top students?

$$\hookrightarrow E[\text{top}] = 75\%$$



Tikhonov Regularization (Ridge Regression)

$$\hat{A} = \arg \min_A \left(\|XA - g\|_2 + s \|A\|_2 \right)$$

$$\hat{A} = (X^T X + sI)^{-1} X^T g$$

$$\|A\|_2 = \sqrt{\sum_{i=1}^d a_i^2}$$

Lasso

$$\hat{A}_s = \arg \min_A \left(\|XA - g\|_2 + s \|A\|_1 \right)$$

$$\|A\|_1 = \sum_{i=1}^d |a_i|$$

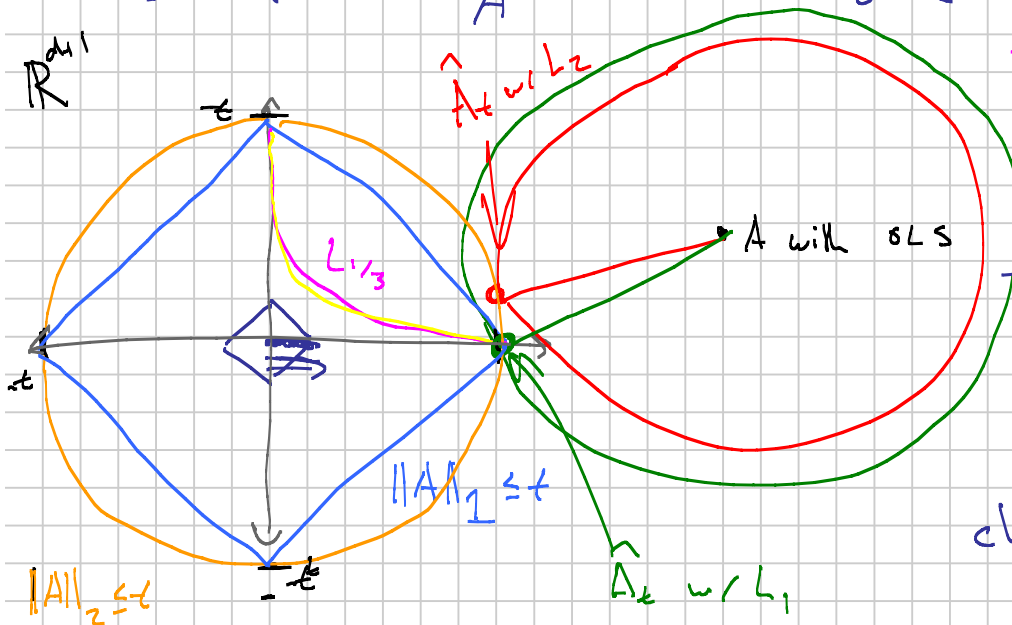
- in high dim \Rightarrow bias towards A sparse
- force to consider value s .

$$\forall s \Rightarrow \exists t \text{ s.t.}$$

s increase
 t decrease

$$t_s = \|\hat{A}_s\|_1$$

$$\hat{A}_s = \hat{A}_t = \arg \min_A \|XA - g\|_2 \text{ s.t. } \|A\|_1 \leq t$$



t large
 \hookrightarrow OLS A dense

t very small
 \hookrightarrow very sparse
 \hookrightarrow 1 non-zero

choose coord j_1
 $j_1 = \arg \max_j |\langle X_j, y \rangle|$
 $a_{j_1}(t) = t$

$$r(A) = y - \sum_{j=0}^d X_j a_j(t)$$

$$a_{j_1}(t) = a_{j_1}(t_2) + (t - t_2) b_1$$

$$a_{j_2}(t) = (t - t_2) b_2$$

until

$$|\langle X_{j_1}, r(t) \rangle| = |\langle X_{j_2}, r(t) \rangle|$$

- Values of a_j can be negative
 - as t increases a value of a_j can oscillate between neg and pos.
 - If a_j hits 0. Snap.
 ↳ only change a_j if $|<, >| = K, >|$
-

Is this better than OLS?

generalization error

