# Lecture 6

- Overfitting & Regularization & model complexity



training data

overfitting

testing data



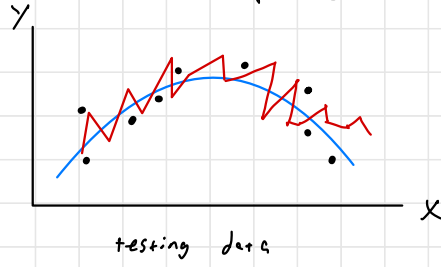$$x_k = \max(0, x - a_k) \quad \text{design}$$

knots

ReLU

$a_k$

$$S = f(x) = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k + \beta_p x_p$$



slope

$\beta_0$

$\beta_1$

$\beta_1 + \beta_2$

$\beta_2$

$\beta_1 + \beta_2 + \beta_3$

$\beta_3$

$a_1 \quad a_2 \quad a_2 \quad a_3$

$\beta_k \equiv$ change of slope
"curvature"

piecewise linear / spline

|   | $p$ | |
|---|-----|---|
| $i$ | $x_i^T$ | $y_i$ |
| $n$ | $X$ | $Y$ |

$$p \gg n$$

$l_2$-regularization

Ridge Regression : $|Y - X\beta|^2 + \lambda |\beta|^2$

we want the curvature to be small.

$$\text{Loss}(\beta) = (Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta$$

$$= Y^T Y - \beta^T X^T Y - Y^T X\beta + \beta^T X^T X\beta + \lambda \beta^T \beta$$

$$= \text{const} - 2\langle X^T Y, \beta \rangle + \beta^T (X^T X + \lambda I_p)\beta$$

$$\text{Loss}'(\beta) = -2X^T Y + 2(X^T X + \lambda I_p)\beta = 0$$

$$\hat{\beta} = (X^T X + \lambda I)^{-1} X^T Y$$

$\lambda > 0$

ridge

Shrinkage estimator

- Note:
$$|\beta|_{\ell_2}^2 = \sum_{k=1}^{p} \beta_k^2 \qquad (\text{We didn't penalize } \beta_0)$$

$$= \beta^T D \beta$$

$$\beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix} \qquad D = \begin{pmatrix} 0 & & & \\ & \tau & & \\ & & \ddots & \\ & & & \tau \end{pmatrix}$$

So the more general version of ridge regression is

$$Loss'(\beta): -2x^T y + 2(x^T x + \underline{D})\beta = 0$$

- Suppose now: $Loss(\beta) = \frac{1}{2} |y - X\beta|^2$

$$= \frac{1}{2} \sum_{i=1}^{n} (y_i - x_i^T \beta_i)^2$$

- Gradient Descent: $\beta_0 = 0 \leftarrow$ <span style="color:red">starting point of gradient descent</span>

$$Loss'(\beta) = - \sum_{i=1}^{n} e_i x_i$$

$$\beta_{t+1} = \beta_t - \eta_t L'(\beta_t)$$

$$\longrightarrow \hat{\beta} = \sum_{i=1}^{n} c_i x_i$$
$$\underset{p \times 1}{} \qquad \underset{p \times 1}{}$$

Representer

$L(\hat{\beta}) = 0$, but there can be other solutions: $L(\tilde{\beta}) = 0$

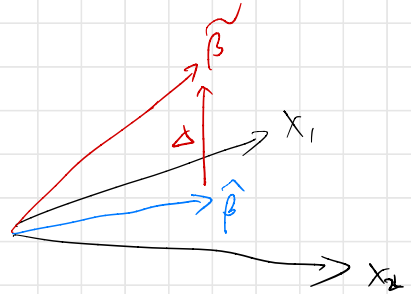- This implies $x_i^T \hat{\beta} = Y_i \quad \forall i$

$$x_i^T \tilde{\beta} = Y_i \quad \forall i$$

Let $\tilde{\beta} = \hat{\beta} + \Delta$

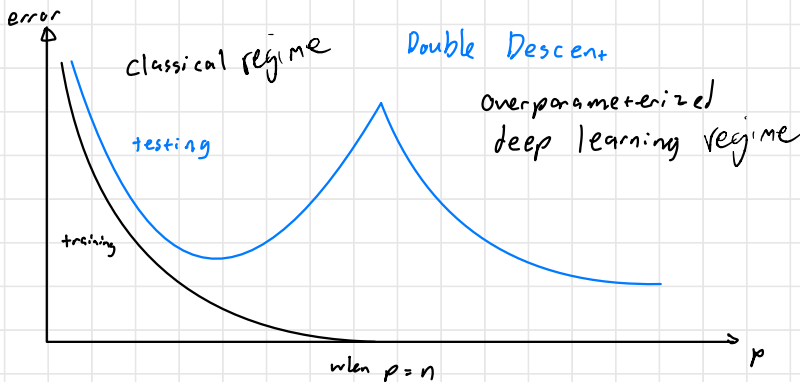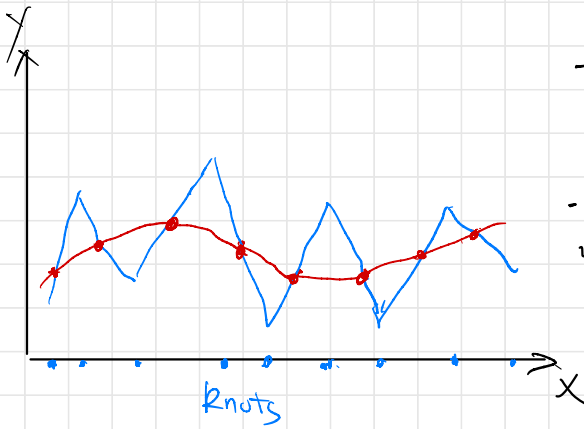Then $x_i^T \Delta = 0 \quad \forall i$

So $\Delta \perp \hat{\beta}$.

- Now $|\tilde{\beta}|^2 = |\hat{\beta}|^2 + |\Delta|^2$

  therefore $\hat{\beta}$ min $|\beta|^2$ among all $\beta$ s.t. $L(\beta) = 0$

- In other words gradient descent self-regularizes.

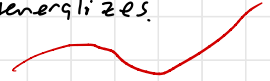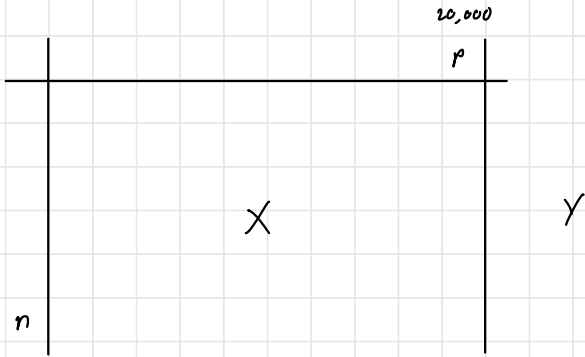  - implicit regularization
  - benign overfitting



classical regime   Double Descent

testing   Overparameterized deep learning regime

training

when $p = n$

- small # of nodes

- with many nodes + GD we get a smooth overfit, so hopefully this generalizes.

knots

• Gene :

$$20,000$$
$$p$$

$X$  $Y$

$n$

"$|\beta|_{\ell_0}$"

Not convex or smooth
$$\text{Loss}(\beta) = \frac{1}{2}|Y - X\beta|^2 + \lambda \#(\beta_j \neq 0)$$

remove irrelivant variables

relax ↓

relax (approximation) ⇓

Convex
$$\text{loss}(\beta) = \frac{1}{2}|Y - X\beta|^2 + \lambda|\beta|_{\ell_1}$$

$$|\beta|_{\ell_1} = \sum_{k=1}^{p}|\beta_k|$$

- Lasso: Least absolute shrinkage selection operator.

$$\text{Loss}(\beta) \;=\; \tfrac{1}{2}\,\Big|\,Y - \sum_{j=1}^{p} x_j \beta_j \,\Big|^2 \;+\; \lambda \sum_{j=1}^{p} |\beta_j|$$

$$Y - \sum_{j \neq k} x_j \beta_j \;-\; X_k \beta_k$$

- Coordinate Descent : $\hat{y}$

  Each iteration :
  for $k$ in $1:p$
  $$\hat{\beta}_k \;=\; \min_{\beta_k} \; |\,\hat{Y} - x_k \beta_k\,|^2 \;+\; \lambda\,|\beta_k|$$

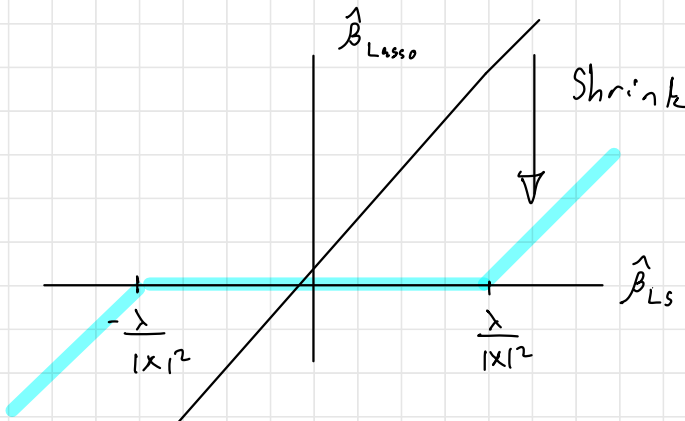  $$\hat{Y} \;=\; Y - \sum_{j \neq k} x_j \beta_j$$

- One-dimensional Problem :

  $$L(\beta) = |\,Y - X\beta\,|^2 + \lambda\,|\beta|$$

  So $\quad \hat{\beta}_{LS} \;=\; \dfrac{\langle x, Y \rangle}{|x|^2} \qquad (\lambda = 0)$

  $$\hat{\beta}_{Lasso} \;=\; \text{sign}(\hat{\beta}_{LS})\,\max\!\Big(0,\;\Big|\hat{\beta}_{LS}\Big| - \frac{\lambda}{|x|^2}\Big)$$

  soft-thresholding , selection



$\hat{\beta}_{Lasso}$

Shrink

$-\dfrac{\lambda}{|x|^2}$ $\qquad \dfrac{\lambda}{|x|^2}$ $\qquad \hat{\beta}_{LS}$
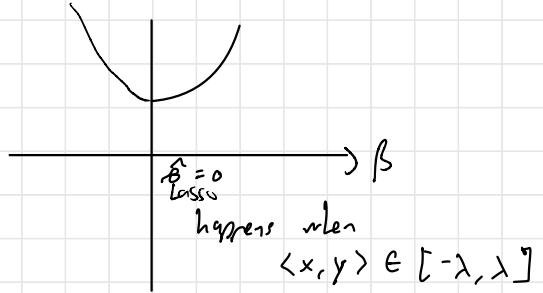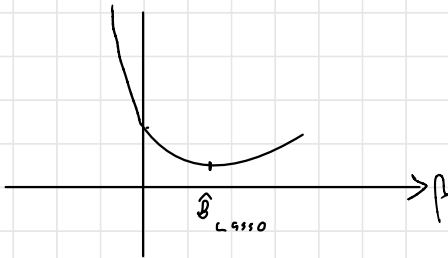
$H_0 : \beta = 0$

$H_1 : \beta \neq 0$

- Consider :

$$L'(\beta) = \begin{cases} -\langle x,y \rangle + |x|^2 \beta + \lambda & \text{if } \beta \geq 0 \\ -\langle x,y \rangle + |x|^2 \beta - \lambda & \text{if } \beta < 0 \end{cases}$$

$L'(0) = -\langle x, y \rangle - \lambda$
left

$L'(0) = -\langle x, y \rangle + \lambda$
right



$\hat{\beta} = 0$
Lasso
happens when
$\langle x,y \rangle \in [-\lambda, \lambda]$

If $L'_{left} < 0$ & $L'_{right} < 0$



$\hat{\beta}_{Lasso}$

$-\langle x,y \rangle + |x|^2 \beta + \lambda = 0$

$$\hat{\beta}_{Lasso} = \hat{\beta}_{LS} - \frac{\lambda}{|x|^2}$$

If $L'_{left} > 0$ & $L'_{right} > 0$



$\hat{\beta}_{Lasso}$

$-\langle x,y \rangle + |x|^2 \beta - \lambda = 0$

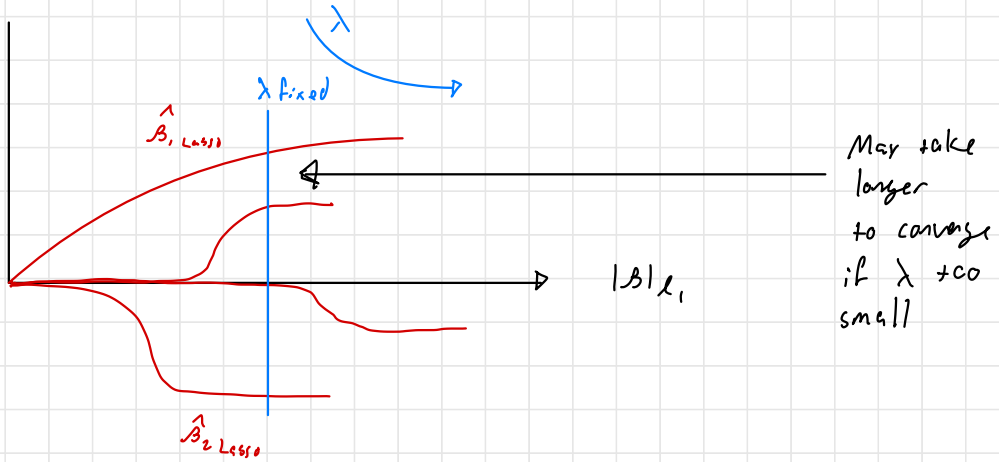$$\hat{\beta}_{Lasso} = \hat{\beta}_{LS} + \frac{\lambda}{|x|^2}$$

○ Solution Path :

　　- Start with $\lambda = \max_j |\langle x_j, Y \rangle|$

　　- then gradually reduce $\lambda$.

　　- $\hat{\beta}_{j \, Lasso} = \max\left(0, \; \hat{\beta}_{j}^{LS} - \dfrac{\lambda}{|x_j|^2}\right)$



May take
larger
to converge
if $\lambda$ too
small

Forward Selection, related to boosting