

## Regularized or Constrained Linear Regression

- goal: restrict the possible coefficients  $\beta$  to avoid overfitting
- overfitting occurs when features are (almost) redundant, i.e. (almost) linearly dependent: 
$$X_j \approx \sum_{j' \neq j} w_{j'} X_{j'} \quad (*)$$

happens in two cases: (1) if system of equations is under determined,  $D > N$  (more features than instances!)

Theorem of linear algebra: (\*) is always fulfilled  

$$[\text{rank}(X) \leq \min(N, D) = N < D]$$

(2) if condition  $\kappa(X)$  is bad (large), (\*) is approximately true  
 why does this cause overfitting: for simplicity, suppose  $X_j = X_{j_1}$  are exactly equal,  $\beta_j X_j + \beta_{j_1} X_{j_1} = 0$  if  $\beta_j = -\beta_{j_1}$   
 now, if  $X$  are slightly noisy,  $\beta_j X_j \approx \beta_{j_1} X_{j_1}$  with ~~very~~ coefficients of very big magnitude  $|\beta_j|, |\beta_{j_1}| \gg 0$   
 if ~~that~~ the test data have different noise, the balance between  $\beta_j X_j$  and  $\beta_{j_1} X_{j_1}$  no longer works  $\Rightarrow$  if  $|\beta_j|, |\beta_{j_1}| \gg 0$ , we get huge error

- conclusion: we want to prevent  $|\beta_j| \gg 0$  "constraint" or "regularization"

$\Rightarrow$  Ridge Regression

$$\hat{\beta} = \arg \min_{\beta} (Y - X\beta)^2$$

$$\text{s.t. } \underbrace{\beta^T \beta = \|\beta\|_2^2}_{\text{constraint}} \leq \epsilon$$

hyperparameter  $\epsilon$



if  $\hat{\beta}_{OLS}$  of OLS solution already has small norm  $\|\hat{\beta}_{OLS}\|_2^2 \leq \epsilon$ ,  
 then the constraint is "inactive" and nothing changes.

Otherwise, the optimal solutions will differ. Rewrite objective with  
 Lagrange multiplier  $\tau$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} (Y - X\beta)^2 + \tau \beta^T \beta \quad \tau \geq 0$$

Lagrange multipl.

Theorem: the two variants of the objective gives the same solution  $\hat{\beta}$

if  $\tau$  is chosen appropriately for each given  $\epsilon$ .

$$\frac{d \text{Loss}}{d\beta} = 2 [X^T(Y - X\beta) + \tau\beta] \stackrel{!}{=} 0$$

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

regularized pseudo-inverse

• for  $\tau = 0$ , this reduces to ordinary pseudo-inverse  $(X^T X)^{-1} X^T$

• for  $\tau > 0$ , the diagonal of  $X^T X + \tau I$  is bigger than the diagonal

of  $\underbrace{X^T X}_{\text{scatter matrix } S}$   $\Rightarrow \kappa(X^T X + \tau I) < \kappa(X^T X)$ , i.e. better condition  
 regularized scatter matrix  $S_\tau$

[reminder:  $X$  and  $Y$  must be centered,  $\bar{X} = 0$ ,  $\bar{Y} = 0$ ]

The effect of adding  $\tau$  to diagonal element  $(X^T X)_{jj}$  depends on the  
 magnitude of  $(X^T X)_{jj}$   $\Rightarrow$  scale the features beforehand to have

unit variance  $\Rightarrow$  standardization  $\tilde{X}_j = \frac{X_j - \bar{X}_j}{\text{std}(X_j)} \Rightarrow (X^T X)_{jj} = 1$  for all  $j$



solutions of ridge regression:

- ① compute and apply regularized pseudo-inverse: easy, but numerically less accurate
- ② reduce it to ordinary least squares: augmented data matrix

$$X' = \begin{bmatrix} X & & \\ \sqrt{\tau} & & 0 \\ & \sqrt{\tau} & \dots \\ 0 & & \sqrt{\tau} \end{bmatrix} \begin{matrix} N \text{ rows} \\ \\ 0 \text{ rows} \end{matrix}$$

$$Y' = \begin{bmatrix} Y \\ 0 \end{bmatrix} \begin{matrix} N \text{ elements} \\ 0 \text{ zeros} \end{matrix}$$

$$\underbrace{(Y' - X'\beta)^T (Y' - X'\beta)}_{\text{OLS objective for } X', Y'} = \underbrace{(Y - X\beta)^T (Y - X\beta) + \tau \beta^T \beta}_{\text{ridge regression objective}}$$

- ③ use singular value decomposition (SVD):

$$X = U \Lambda V^T$$

$$X^T = V \Lambda U^T$$

$$X^T X = V \Lambda^2 V^T$$

$$X^T X + \tau I = V \Lambda^2 V^T + \tau I = V \Lambda^2 V^T + \tau V I V^T$$

$$V V^T = I \approx V I V^T = I$$

$$= V (\Lambda^2 + \tau I) V^T$$

$$(X^T X + \tau I)^{-1} X^T = (V (\Lambda^2 + \tau I) V^T)^{-1} V \Lambda U^T = \underbrace{V (\Lambda^2 + \tau I)^{-1} \Lambda}_{\Lambda'} U^T$$

$\Lambda'$  is a diagonal matrix with elements  $\lambda_j' = \frac{\lambda_j}{\lambda_j^2 + \tau}$

for  $\tau = 0$ ,  $\lambda_j' = \frac{1}{\lambda_j^2} \Rightarrow \Lambda' = \Lambda^{-1}$ , reduces to ordinary pseudo-inverse



• if  $\tau > 0$ , rewrite

$$\lambda_j' = \frac{\lambda_j}{\lambda_j^2 + \tau} = \frac{1}{\lambda_j + \frac{\tau}{\lambda_j}}$$

two cases: (1)  $\lambda_j \gg 0$  (no redundancy)  $\Rightarrow \frac{\tau}{\lambda_j} \approx 0 \Rightarrow \lambda_j' \approx \frac{1}{\lambda_j}$

$\Rightarrow$  if  $\lambda_j \gg \tau$ , regularization has no effect

(2)  $\lambda_j \approx 0$  (redundant features)  $\Rightarrow \lambda_j' = \frac{0}{0^2 + \tau} = 0$

$\Rightarrow$  in ordinary least squares  $\lambda_j' = \frac{1}{\lambda_j} \gg 0$  pseudo-inverse "explodes"

ridge regression

$\lambda_j' = \frac{\lambda_j}{\tau} \approx 0 \Rightarrow$  no "explodes"

(4) solve via dual optimization problem  $\Rightarrow$  later

• What's the draw back of regularization? Does improved conditioning come at a price?

### Bias - Variance - Trade-off

Suppose, we have  $M$  training sets of size  $N$ . What happens when  $M \rightarrow \infty$ ?

Repeat experiment  $M$  times and check how much  $\beta$  varies: if it varies a lot  $\Rightarrow$  overfitting, otherwise no little/no overfitting.

Mathematically:  $\hat{\beta}_m$  = OLS-solution for TS in  $(m=1 \dots M)$

expectation of all TS:  $E_m[\hat{\beta}_m]$ ,  $\beta^*$ : unknown true solution

bias (systematic error) =  $\beta^* - E_m[\hat{\beta}_m]$ , variance =  $E_m[(\hat{\beta}_m - E_m[\hat{\beta}_m])^2]$



- calculate the expected mean-square error of  $\hat{\beta}_m$

$$\begin{aligned} \text{MSE} &= E_m[(\hat{\beta}_m - \beta^*)^2] = E_m[(\hat{\beta}_m - E_m[\hat{\beta}_m] + (E_m[\hat{\beta}_m] - \beta^*))^2] \\ &= \underbrace{E_m[(\hat{\beta}_m - E_m[\hat{\beta}_m])^2]}_{\text{expected variance}} + \underbrace{E_m[(E_m[\hat{\beta}_m] - \beta^*)^2]}_{\text{expected squared bias}} \end{aligned}$$

[cross terms vanish:

$$\begin{aligned} &-2 E_m[(\hat{\beta}_m - E_m[\hat{\beta}_m])(E_m[\hat{\beta}_m] - \beta^*)] \\ &= -2 \left( E_m[\hat{\beta}_m](E_m[\hat{\beta}_m] - \beta^*) + \underbrace{E_m[E_m[\hat{\beta}_m]](E_m[\hat{\beta}_m] - \beta^*)}_{E_m[\hat{\beta}_m]} \right) = 0 \end{aligned}$$

- apply bias-variance decomposition to OLS:  $\hat{\beta}_m = (X_m^T X_m)^{-1} X_m^T \underbrace{y_m}_{= X_m \beta^* + \epsilon_m}$

$$\begin{aligned} E_m[\hat{\beta}_m] &= E_m[(X_m^T X_m)^{-1} X_m^T (X_m \beta^* + \epsilon_m)] \\ &= E_m[\underbrace{(X_m^T X_m)^{-1} (X_m^T X_m)}_{= I} \beta^* + \underbrace{(X_m^T X_m)^{-1} X_m^T \epsilon_m}_{= E_m[(X_m^T X_m)^{-1} X_m^T] E_m[\epsilon_m]}] \\ &= E_m[\beta^*] = \beta^* \end{aligned}$$

$= 0 \quad \epsilon_m \sim \mathcal{N}(0, \sigma^2)$

$$\Rightarrow \cancel{E_m[(\hat{\beta}_m - \beta^*)^2]} = E_m[(E_m[\hat{\beta}_m] - \beta^*)^2] = E_m[(\beta^* - \beta^*)^2] = 0$$

squared bias of OLS is zero: "OLS is unbiased"



expected co-variance matrix of  $\hat{\beta}_m$  across infinitely many TS

$$E_m[(\hat{\beta}_m - E[\hat{\beta}_m])^2] = E_m[(\hat{\beta}_m - \beta^*)^2]$$

$$(\hat{\beta}_m - \beta^*)^2 = (\hat{\beta}_m - \beta^*)(\hat{\beta}_m - \beta^*)^T$$

$$\begin{aligned}\hat{\beta} &= (X_m^T X_m)^{-1} X_m^T Y_m \\ &= \beta^* + (X_m^T X_m)^{-1} X_m^T \varepsilon_m\end{aligned}$$

$$\hat{\beta}_m - \beta^* = (X_m^T X_m)^{-1} X_m^T \varepsilon_m$$

$$(\hat{\beta}_m - \beta^*)(\hat{\beta}_m - \beta^*)^T = (X_m^T X_m)^{-1} X_m^T \varepsilon_m \varepsilon_m^T X_m (X_m^T X_m)^{-1}$$

$$E_m[(\hat{\beta}_m - \beta^*)(\hat{\beta}_m - \beta^*)^T] = E_m[(X_m^T X_m)^{-1} X_m^T] \underbrace{E_m[\varepsilon_m \varepsilon_m^T]}_{= \sigma^2 I} E_m[X_m (X_m^T X_m)^{-1}]$$

noise of all instances is independent  $\varepsilon_m \sim \mathcal{N}(0, \sigma^2)$

$$= \sigma^2 E_m[\underbrace{(X_m^T X_m)^{-1} X_m^T X_m (X_m^T X_m)^{-1}}_{= I}]$$

$$E_m[(\hat{\beta}_m - \beta^*)^2] = E_m[(X_m^T X_m)^{-1}] \sigma^2, \text{ for single TS: } (X^T X)^{-1} \sigma^2$$

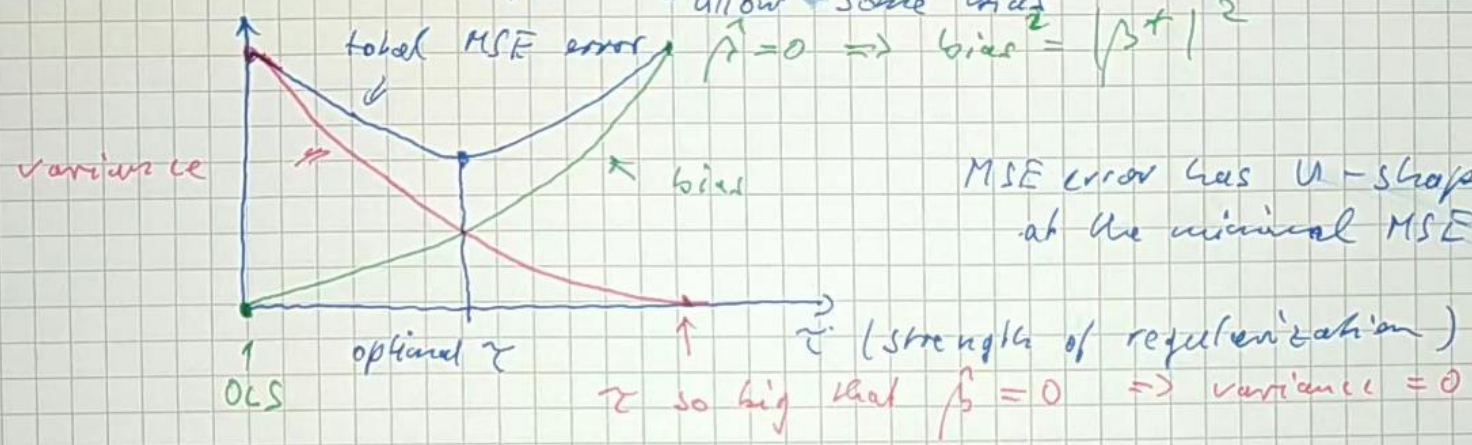
variance of solution  $\hat{\beta}$  is proportional to noise variance  $\sigma^2$  and  
to inverse scatter matrix  $(X_m^T X_m)^{-1}$

$\Rightarrow$  if features are redundant,  $X_m^T X_m$  has bad condition,  $(X_m^T X_m)^{-1}$  explodes

$\Rightarrow$  high variance of  $\hat{\beta}_m$  over TS  $\Rightarrow$  we cannot profit from unbiased res



⇒ bias-variance trade-off: sacrifice some bias, if it reduces the variance a lot  
"allow some bias"



bias and variance of ridge regression:

$$\begin{aligned} E_m [\hat{\beta}_n(\tau)] &= (X^T X + \tau I)^{-1} X^T X \cdot \beta^* \quad \text{for } \tau=0: E_m[\hat{\beta}_n(0)] = \beta^* \\ &= V \operatorname{diag} \left\{ \frac{\lambda_j^2}{\lambda_j^2 + \tau} \right\} V^T \beta^* \end{aligned}$$

variance ~~the~~  $\operatorname{Cov}[\hat{\beta}_\tau] = (X^T X + \tau I)^{-1} X^T X (X^T X + \tau I)^{-1} \sigma^2$

$$= V \operatorname{diag} \left( \frac{\lambda_j^2}{(\lambda_j^2 + \tau)^2} \right) V^T \sigma^2$$