

Lecture 10 : Tikhinov Regularization / Intro to SVD

Recall : $X \in \mathbb{R}^{n \times p}$ (n training samples, p features)
 $y \in \mathbb{R}^n$ (n labels)

Model : $y \approx Xw$, $y_i \approx x_i^T w$ for some $w \in \mathbb{R}^p$

Least Squares : $\hat{w}_{LS} = \underset{w}{\operatorname{argmin}} \|y - Xw\|_2^2 = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n (y_i - x_i^T w)^2$

- if X is full rank (cols are linearly indep),
 then \hat{w}_{LS} is unique and $\hat{w}_{LS} = (X^T X)^{-1} X^T y$

- if X not full rank, $X^T X$ is not invertible
 \hat{w}_{LS} is unique.

2

Alternative: Tikhonov Regularization / Ridge Regressions

$$\hat{w} = \arg \min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

measures fit to data

"regularizer" measures "energy" in w

$\lambda > 0$

"regularization parameter" or
"tuning parameter"

- \hat{w} unique even when no least squares soln exists.
 - even when X is full ~~full~~ rank, it can be "badly behaved", and regularization adjusts for this.

(3)

e.g. $X^T X = \begin{bmatrix} 1 & 0 \\ 0 & 10^{-12} \end{bmatrix}$ $(X^T X)^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 10^{12} \end{bmatrix}$

$$\begin{aligned} \text{let } f(\bar{w}) &= \|y - Xw\|_2^2 + \lambda \|w\|_2^2 \\ &= y^T y - 2w^T X^T y + \underline{w^T X^T X w + \lambda w^T w} \\ &= y^T y - 2w^T X^T y + w^T (X^T X + \lambda I) w \end{aligned}$$

$$\nabla_w f = 0 - 2X^T y + 2(X^T X + \lambda I)w = 0$$

if $(X^T X + \lambda I)$ is invertible, then $\hat{w} = (X^T X + \underline{\lambda I})^{-1} X^T y$
always true

note: $X^T X + \lambda I$ is always invertible

recall Q is invertible if it is positive definite

i.e. $a^T Q a > 0$ for all $\boxed{a \neq 0}$

$$\begin{aligned} a^T (X^T X + \lambda I) a &= a^T X^T X a + \lambda a^T a \\ &= \underbrace{\|Xa\|_2^2}_{\geq 0} + \lambda \underbrace{\|a\|_2^2}_{> 0} > 0 \end{aligned}$$

alternative derivation of \hat{w} :

$$\begin{aligned} \text{1st note } \|a\|_2^2 + \|b\|_2^2 &= \left\| \begin{bmatrix} a \\ b \end{bmatrix} \right\|_2^2 \\ &= \sum a_i^2 + \sum b_i^2 \end{aligned}$$

$$\begin{aligned}
 \|y - Xw\|_2^2 + \lambda \|w\|_1^2 &= \|y - Xw\|_2^2 + \|\sqrt{\lambda} w\|_2^2 \\
 &= \left\| \begin{bmatrix} y - Xw \\ \sqrt{\lambda} w \end{bmatrix} \right\|_2^2 = \left\| \begin{bmatrix} y \\ 0 \end{bmatrix} - \begin{bmatrix} Xw \\ \sqrt{\lambda} w \end{bmatrix} \right\|_2^2 \\
 &= \left\| \begin{bmatrix} y \\ 0 \end{bmatrix} - \underbrace{\begin{bmatrix} X \\ \sqrt{\lambda} I \end{bmatrix}}_{\tilde{X}} w \right\|_2^2 = \|\tilde{y} - \tilde{X}w\|_2^2
 \end{aligned}$$

$$\hat{w} = \underset{w}{\operatorname{argmin}} \| \tilde{y} - \tilde{X} w \|_2^2 \Rightarrow \hat{w} = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{y}$$

$$\tilde{X}^T \tilde{X} = [X^T \quad \sqrt{\lambda} I] \begin{bmatrix} X \\ \sqrt{\lambda} I \end{bmatrix} = X^T X + \lambda I \quad ; \quad \tilde{X}^T \tilde{y} = [X^T \quad \sqrt{\lambda} I] \begin{bmatrix} y \\ 0 \end{bmatrix} = X^T y$$

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

example: $y = Xw + \varepsilon$ $\iff y_i = x_i^T w + \varepsilon_i$
 \dagger noise or errors in y or our model
~~egen~~

$$X = \begin{bmatrix} 1 & 1 & .01 \\ 1 & -1 & .01 \\ 1 & 1 & -.01 \\ 1 & -1 & -.01 \end{bmatrix}$$

cols of X are lin. indep. (orthogonal)
 $\Rightarrow X$ is full rank, $X^T X$ is invertible.

$$X^T X = 4 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 10^{-4} \end{bmatrix} \Rightarrow (X^T X)^{-1} = \frac{1}{4} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 10^4 \end{bmatrix}$$

If least squares:

$$\begin{aligned}
 \hat{w} &= (X^T X)^{-1} X^T y \\
 &= (X^T X)^{-1} X^T (Xw + \varepsilon) \\
 &= \cancel{(X^T X)^{-1}} \cancel{X^T X} w + (X^T X)^{-1} X^T \varepsilon \\
 &= \underbrace{w}_{\text{good}} + \underbrace{(X^T X)^{-1} X^T \varepsilon}_{\begin{bmatrix} 1 \\ 5 \cdot 10^4 \end{bmatrix}}
 \end{aligned}$$

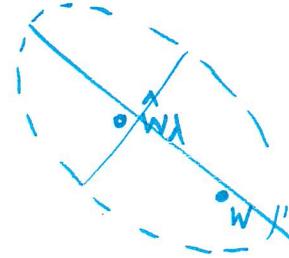
LS soln exists
(unique)
but small noise in ε
gets amplified by
large values in $(X^T X)^{-1}$

$$(X^T X + \lambda I)^{-1} = \begin{bmatrix} 4+\lambda & & \\ & 4+\lambda & \\ & & \underline{4 \cdot 10^{-4} + \lambda} \end{bmatrix}$$

$$(X^T X + \underline{\lambda I})^{-1} = \begin{bmatrix} \cancel{4+\lambda} & & \\ & \cancel{4+\lambda} & \\ & & \underline{\approx \frac{1}{\lambda}} \end{bmatrix}$$

↳ avoids amplifying noise

Singular Value Decomposition



Finding λ

- split data into training and validation set. $(X^{(t)}, y^{(t)}, X^{(v)}, y^{(v)})$

consider a set of possible λ 's : $\Lambda = \{\lambda_1, \lambda_2, \dots, \lambda_M\}$

- for each $\lambda \in \Lambda$, compute $\hat{w}_\lambda = (X^{(t)T} X^{(t)} + \lambda I)^{-1} X^{(t)T} y$

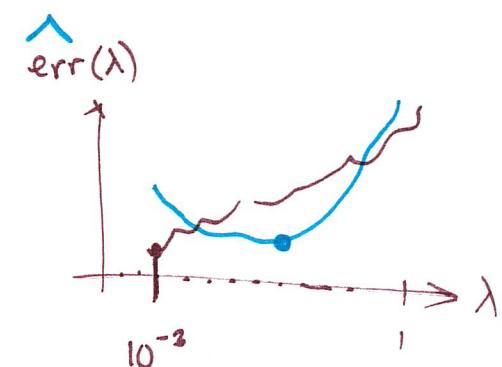
measure error of \hat{w}_λ on validation set

$$\text{err}(\lambda) = \|y^{(v)} - X^{(v)} \hat{w}_\lambda\|_2^2$$

choose $\lambda_{\text{opt}} = \underset{\lambda}{\operatorname{argmin}} \text{err}(\lambda)$

$$\hat{w} = (X^T X + \lambda_{\text{opt}} I)^{-1} X^T y$$

$$\hat{w} = \underset{w}{\operatorname{argmin}} \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$



$$\text{err}(\lambda) = \mathbb{E} (y - X^T \hat{w}_\lambda)^2$$

Sing. Val. Decomp. (SVD)

$$\underbrace{\begin{bmatrix} 1 & 1 \\ 1 & 1.01 \end{bmatrix}}_X \underbrace{\begin{bmatrix} w_1 \\ w_2 \end{bmatrix}}_y = \underbrace{\begin{bmatrix} 2 \\ 2 \end{bmatrix}}_z \Rightarrow X^{-1}y = \begin{bmatrix} z \\ 0 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$\underbrace{\begin{bmatrix} 1 & 1 \\ 1 & 1.01 \end{bmatrix}}_X \underbrace{\begin{bmatrix} w_1 \\ w_2 \end{bmatrix}}_y = \underbrace{\begin{bmatrix} 2 \\ 2.01 \end{bmatrix}}_z \Rightarrow X^T y = \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

System is very sensitive to small changes in y .

14

SVD will help us quantify how well-behaved X is
in terms of its inverse or rank

another ex