# Lecture 9:

# SVD in Machine Learning

# Ridge Regression

$(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}, \quad i = 1, \ldots, n$

until now, we considered $n > p$, $X = \begin{bmatrix} -x_1^T- \\ -x_2^T- \\ \vdots \\ -x_n^T- \end{bmatrix}$ w/ $p$ LI columns. Then $\hat{w}_{LS} = (X^TX)^{-1}X^Ty$

These assumptions were necessary to ensure there was a unique set of weights minimizing squared error — and that $X^TX$ had an inverse.

Let $X = U\Sigma V^T \implies X^TX = (U\Sigma V^T)^T(U\Sigma V^T) = V\Sigma^T U^T U \Sigma V^T = V\Sigma^T\Sigma V^T$

Then $(X^TX)^{-1} = (V\Sigma^T\Sigma V^T)^{-1} = (V^T)^{-1}(\Sigma^T\Sigma)^{-1}V^{-1} = V(\Sigma^T\Sigma)^{-1}V^T$

Now $\Sigma^T\Sigma = \begin{bmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \text{\Large O} \\ & & \ddots & \\ \text{\Large O} & & & \sigma_r^2 \end{bmatrix}$, so we can only invert $\Sigma^T\Sigma$ if all the $\sigma_i$'s $> 0$

That is,

> $X$ has $p$ LI columns $\iff$ $X^TX$ invertible $\iff$ $X$ is positive definite $\iff$ all singular values $> 0$

If $X$ has $r < p$ LI columns, then $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > 0 = \sigma_{r+1} = \sigma_{r+2} = \cdots = \sigma_p$

> Can we still learn a predictor when $X$ has $r < p$ LI columns?
> e.g., if $n < p$ (more features than training samples), then there are at most $n$ LI cols

Recall ridge regression:

$$\hat{\underline{w}}_\lambda = \underset{\underline{w}}{\text{argmin}} \; \|y - X\underline{w}\|_2^2 + \lambda \|\underline{w}\|_2^2 = (X^TX + \lambda I)^{-1} X^T y$$

Let $X = U\Sigma V^T$, as before.

Then $X^TX + \lambda I = V\Sigma^T\Sigma V^T + \lambda I = V\Sigma^T\Sigma V^T + \lambda \underbrace{VV^T}_{I} = V\Sigma^T\Sigma V^T + V(\lambda I)V^T$

$$= V(\Sigma^T\Sigma + \lambda I)V^T$$

$$\Sigma^T\Sigma = \begin{bmatrix} \sigma_1^2 & & & \bigcirc \\ & \sigma_2^2 & & \\ & & \ddots & \\ \bigcirc & & & \sigma_p^2 \end{bmatrix} \Rightarrow \Sigma^T\Sigma + \lambda I = \begin{bmatrix} \sigma_1^2 + \lambda & & & \bigcirc \\ & \sigma_2^2 + \lambda & & \\ & & \ddots & \\ \bigcirc & & & \sigma_p^2 + \lambda \end{bmatrix}$$

Even if $\sigma_p = 0$,
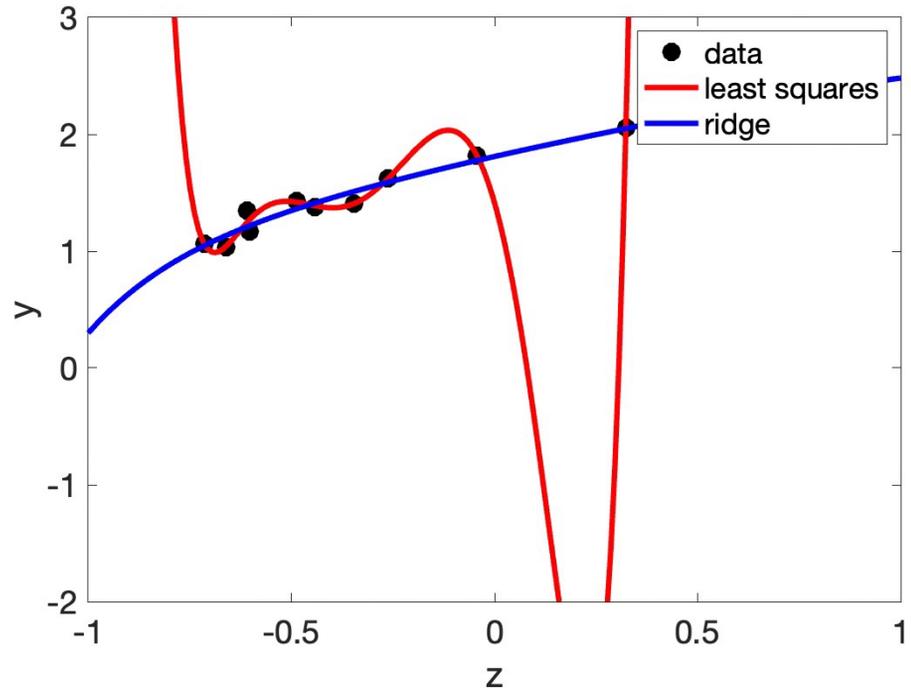$\sigma_p^2 + \lambda > 0$
$\Rightarrow \Sigma^T\Sigma$ may not be invertible

but $\Sigma^T\Sigma + \lambda I$ is _always_

invertible

$\Rightarrow$ We can _always_ compute the ridge regression estimate, even when a unique least squares estimate does not exist

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

$$= V(\Sigma^T\Sigma + \lambda I) V^T V \Sigma^T U^T y$$

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

$$\begin{bmatrix} \dfrac{\sigma_1}{\sigma_1^2 + \lambda} & & \bigcirc \\ & \dfrac{\sigma_2}{\sigma_2^2 + \lambda} & \\ & & \ddots & \\ & & & \dfrac{\sigma_p}{\sigma_p^2 + \lambda} \\ \bigcirc & & \end{bmatrix}$$

Ridge regression can help prevent overfitting



WHY ?

Imagine $X$ has some very small singular values, e.g. $10^{-6}$.

Also, imagine $y = X\underline{w}^* + \underline{\varepsilon}$ where $\underline{\varepsilon}$ is random error or noise added to each sample

Then
$$\hat{\underline{w}}_{LS} = (X^TX)^{-1}X^Ty = (X^TX)^{-1}X^T(X\underline{w}^* + \underline{\varepsilon})$$

$$= (X^TX)^{-1}X^TX\underline{w}^* + (X^TX)^{-1}X^T\underline{\varepsilon}$$

$$= \underline{w}^* + V^T(\Sigma^T\Sigma)^{-1}\Sigma^TU^T\underline{\varepsilon}$$

$\underbrace{\qquad}_{\substack{\text{good!}\\\text{what we}\\\text{want}}}$

$$\begin{bmatrix} 1/\sigma_1 & & & \\ & 1/\sigma^2 & & \\ & & \ddots & \\ & & & 1/\sigma_p \end{bmatrix}$$

$\leftarrow \sim 10^6 \Rightarrow$ a little noise in $\underline{\varepsilon}$ gets multiplied by a huge number $\Rightarrow$ least squares might fit observations very closely but give strange predictions on test data

In contrast, $\hat{\underline{w}}_R = (X^TX + \lambda I)^{-1}X^T(X\underline{w}^* + \underline{\varepsilon})$

$$= (X^TX + \lambda I)^{-1}X^TX\underline{w}^* \qquad + (X^TX + \lambda I)^{-1}X^T\underline{\varepsilon}$$

$$= V\underbrace{(\Sigma^T\Sigma + \lambda I)^{-1}\Sigma^T\Sigma}_{}V^T\underline{w}^* \qquad + V(\Sigma^T\Sigma + \lambda I)^{-1}\Sigma^TU^T\underline{\varepsilon}$$

$$\begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & \\ & \ddots & \\ & & \frac{\sigma_p^2}{\sigma_p^2 + \lambda} \end{bmatrix}$$

$\approx \underline{w}^*$

for small $\lambda$

$\dfrac{\sigma_i}{\sigma_i^2 + \lambda} \approx \begin{cases} 1/\sigma_i & \sigma_i \gg \lambda \\ 1/\lambda & \sigma_i \approx 0 \end{cases}$

none of these values "blow up" by being close to dividing by 0, so we do not magnify noise $\Rightarrow$ better predictions on test data

The Elephant In the Room — how do we choose $\lambda$?

1. Split data into 2 sets $(\underline{x}_i, y_i), i=1,...,m$ = training set

$(\underline{x}_i, y_i), i=m+1,...,n$ = validation set.

2. For each $\lambda \in \{\lambda_1, \lambda_2, ..., \lambda_g\}$, find $\underline{\hat{w}}_\lambda$ using training data.

Measure the loss of each $\lambda$ using validation set: $L_\lambda = \sum\limits_{i=m+1}^{n} (y_i - \underline{x}_i^T \underline{\hat{w}}_\lambda)^2$

Choose $\lambda$ with smallest $L_\lambda$.

---

Alternative to ridge regression (when some singular values $= 0$ so many $\underline{w}$ fit $\underline{y}$ perfectly)

with least squares, we needed to compute $(\Sigma^T \Sigma)^{-1} \Sigma^T = \begin{bmatrix} 1/\sigma_1 & & & \\ & 1/\sigma & & \\ & & \ddots & \\ & & & 1/\sigma_p \end{bmatrix}$ ?

which was problematic for any $\sigma_j = 1$

Alternative: define pseudoinverse of $\Sigma$ as

$(\Sigma^+)_{ii} = \begin{cases} 1/\sigma_i & \text{if } \sigma_i > 0 \\ 0 & \text{otherwise} \end{cases}$

that is, $\Sigma^+$ corresponds to taking the transpose of $\Sigma$ and only inverting the nonzero diagonal entries

Special case: $p > n$, X has $n$ linearly independent rows.

$\Rightarrow \Sigma \sim n \times p$

$$X_{n \times p} = U_{n \times n} \; \Sigma_{n \times p} \; V^T_{p \times p}$$

$$\Sigma = \begin{bmatrix} \sigma_1 & & & & \\ & \sigma_2 & & & \mathcal{O} \\ & & \ddots & & \\ & & & \sigma_n & \end{bmatrix}$$

$\underbrace{\qquad}_{\substack{n \text{ nonzero} \\ \text{singular vals.}}}$ $\underbrace{\qquad}_{\substack{p-n \text{ columns} \\ \text{of zeros.}}}$

$$\Rightarrow \Sigma^+ = \begin{bmatrix} 1/\sigma_1 & & & \\ & 1/\sigma_2 & & \\ & & \ddots & \\ & & & 1/\sigma_n \\ \hline & & \mathcal{O} & \end{bmatrix} \quad \left. \vphantom{\begin{matrix} a \\ b \end{matrix}} \right\} \substack{p-n \text{ rows of} \\ \text{zeros}}$$

$= \Sigma^T (\Sigma \Sigma^T)^{-1}$

pseudoinverse solution $\underline{\hat{w}} = V \Sigma^+ U^T \underline{y}$

$\qquad = V \Sigma^T (\Sigma \Sigma^T)^{-1} U^T \underline{y}$

$\qquad = X^T (X X^T)^{-1} \underline{y}$

**Claim:** when $X$ has $n$ LI Rows, the choice $\hat{w} = V \Sigma^{+} U^{T} y = X^{T}(XX^{T})^{-1} y$

has the smallest norm $\|w\|_{2}^{2}$ of any $w$ satisfying $y = Xw$.

**Proof:** for any $w$,

$$\|w\|_{2}^{2} = \|w - \hat{w} + \hat{w}\|_{2}^{2} = \|w - \hat{w}\|_{2}^{2} + 2(w-\hat{w})^{T}\hat{w} + \|\hat{w}\|_{2}^{2}$$

Assume $Xw = y = X\hat{w} \Rightarrow \underline{X(w-\hat{w}) = 0}$

Then $(w-\hat{w})^{T}\hat{w} = (w-\hat{w})^{T} X^{T}(XX^{T})^{-1} y$

$$= \left(\underline{X(w-\hat{w})}\right)^{T} (XX^{T})^{-1} y$$

$$= \underline{0}$$

$$\Rightarrow \|w\|_{2}^{2} = \|w-\hat{w}\|_{2}^{2} + \|\hat{w}\|_{2}^{2} \geq \|\hat{w}\|_{2}^{2} \Longrightarrow \hat{w} \text{ has smallest norm!}$$

**Why does minimum norm make sense?**

e.g. $X = \begin{bmatrix} 1 & 0 & 0.1 \\ 0 & 1 & 0 \end{bmatrix}$, $y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$.

possible $\underline{w}$'s include $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 10 \end{bmatrix}, \begin{bmatrix} -10 \\ 0 \\ 100 \end{bmatrix}, \begin{bmatrix} -1000 \\ 0 \\ 10000 \end{bmatrix}$

$\uparrow$ minimum norm!