

Review

if columns of X are Linearly independent,

then $\underset{w}{\operatorname{argmin}} \|Xw - y\|_2^2 = \underbrace{(X^T X)^{-1} X^T y}$

pseudoinverse of X
 X^+

$$XX^+ = \text{projection matrix onto } \text{span}(\text{cols}(X))$$

if $X = U\Sigma V^T$ (SVD of X), then

$$X^+ = \underbrace{(V\Sigma U^T U\Sigma V^T)^{-1} V\Sigma U^T}_{= V^T \Sigma^{-2} U^T U\Sigma V^T} =$$

$$X^+ = (V\Sigma U^T U\Sigma V^T)^{-1} V\Sigma U^T$$

$$= V\Sigma^{-2} V^T V\Sigma U^T$$

$$= V\Sigma^{-1} U^T$$

Now if cols of X might not be linearly independent

$$X = \begin{array}{|c|} \hline \text{ } \\ \hline \end{array} = \begin{array}{|c|c|} \hline U_1 & \begin{array}{|c|} \hline \text{ } \\ \hline \end{array} \\ \hline \end{array} + \begin{array}{|c|c|} \hline \Sigma_1 & \begin{array}{|c|} \hline \text{ } \\ \hline \end{array} \\ \hline \end{array} + \begin{array}{|c|c|} \hline V_1^T & \begin{array}{|c|} \hline \text{ } \\ \hline \end{array} \\ \hline \end{array}$$

$n \times p \qquad n \times n \qquad n \times p \qquad p \times p$

if X has r L.I. cols, then Σ has r nonzero singular values.

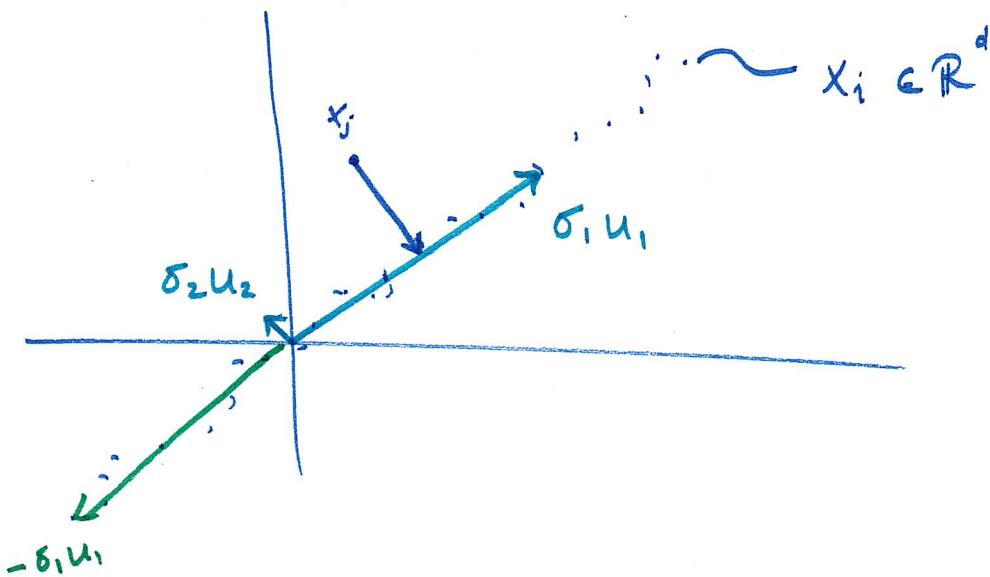
let Σ^+ : take reciprocal of non-zero elements of Σ , then transpose result.

e.g. $\Sigma = \begin{bmatrix} 4 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \rightarrow \Sigma^+ = \begin{bmatrix} \frac{1}{4} & 0 & 0 & 0 \\ 0 & \frac{1}{3} & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$

then

$$\boxed{X^+ = V \Sigma^+ U^T} = \text{pseudoinverse of } X$$

$$X^+ = V_i \Sigma_i^{-1} U_i^T$$



if $X = \begin{bmatrix} 1 & 0 \\ 0 & 3 \\ 0 & 0 \end{bmatrix}$,

$$U = I, V = I, \Sigma = \begin{bmatrix} 1 & & \\ & 2 & \\ & & 3 \\ & & & 4 \end{bmatrix}$$

if $x_i \in$ subspace w/ basis U ,
then $x_i = U a_i$ for some a_i

σ_i 's are increasing

If wanted σ_i 's decreasing,

$$\Sigma = \begin{bmatrix} 4 & & & \\ & 3 & & \\ & & 2 & \\ & & & 1 \end{bmatrix}$$

$$X = U V^T$$

$$\begin{bmatrix} 1 & & & \\ 2 & 0 & & \\ 0 & 3 & 0 & \\ 0 & 0 & 4 & \end{bmatrix} = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 & & & \\ & 3 & & \\ & & 2 & \\ & & & 1 \end{bmatrix} \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$


=

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

Lecture 20: Stochastic Gradient Descent

minimize $\underbrace{l(w)}_w + \lambda \underbrace{r(w)}_{\text{regularization function}}$

$f_i(w) = (y_i - x_i^T w)^2$

loss, measures fit to data

e.g. $l(w) = \|y - Xw\|_2^2 = \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{n}$

$\nabla l(w) = \sum_{i=1}^n -2x^T(y - Xw) = \sum_{i=1}^n -2(y_i - x_i^T w)x_i$

e.g. $l(w) = \sum_{i=1}^n (1 - y_i x_i^T w)_+ \quad (\text{hinge loss})$

$f_i(w) = (1 - y_i x_i^T w)_+$

$\nabla l(w) = \sum_{i=1}^n I_{\{1 > y_i x_i^T w\}} (-y_i x_i)$

$r(w) = \|w\|_2^2$
(ridge)

$r(w) = \|w\|_1$
(lasso)

in either case

$$\text{let } f(w) = l(w) + \lambda r(w)$$

goal: minimize $f(w)$

$$\text{write } f(w) = \sum_{i=1}^n f_i(w)$$

before: Gradient Descent

$$\hat{w}^{(k+1)} = \hat{w}^{(k)} - \frac{\tau}{2} \nabla f(\hat{w}^{(k)})$$

(e.g., w/ squared error loss:

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

today: Stochastic Gradient Descent

| @ iteration k , choose (SGD)
| $i_k \in \{1, \dots, n\}$

$$|\hat{w}^{(k+1)} = \hat{w}^{(k)} - \frac{\tau}{2} \nabla f_{i_k}(\hat{w}^{(k)})$$

SGD

- each iteration is easier/faster to compute
- need more iterations.

How to choose i_k ?

A. Cyclical (incremental gradient descent)

$$i_k = k \bmod n$$

$$\text{if } n = 3, \quad i_k = 1, 2, 3, \underbrace{1, 2, 3}, \underbrace{1, 2, 3}$$

B. random permutations

every n rounds, reshuffle training data

$$i_k's = \underbrace{1, 3, 2}, \underbrace{3, 1, 2}, \underbrace{2, 1, 3} \dots$$

C. Stochastic gradient descent

$$i_k \sim \underline{\text{unif}}(1, \dots, n)$$

$$i_k's = 1, 3, 3, 2, 3, 1, 2, 2, 2, 1, 3$$

$$\mathbb{E}[\nabla f_i(w)] = \nabla f(w)/n$$

$$\text{Ex 1 : } f(w) = \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

$$= \sum_{i=1}^n (y_i - x_i^T w)^2 + \underline{\lambda \|w\|_2^2}$$

$$\Rightarrow f_i(w) = (y_i - x_i^T w)^2 + \frac{\lambda}{n} \|w\|_2^2$$

$$\Rightarrow \sum_{i=1}^n f_i(w) = f(w)$$

$$\nabla f_i(w) = -2(y_i - x_i^T w)x_i + 2\frac{\lambda}{n}w$$

SGD :

$$\hat{w}^{(k+1)} = \hat{w}^{(k)} - \frac{\tau}{2} \left(-2(y_i - x_i^T \hat{w}^{(k)})x_i + 2\frac{\lambda}{n}\hat{w}^{(k)} \right)$$

$$= \hat{w}^{(k)} + \tau(y_i - x_i^T \hat{w}^{(k)})x_i - \frac{\tau\lambda}{n}\hat{w}^{(k)}$$