

Feb 18, 2020.

Today : Connections.

- Optimization  $\leftrightarrow$  Probability.
- Principal Components Regression.
- TLS..

Ridge regression: minimize:  $\|X\vec{w} - \vec{y}\|_2^2 + \lambda \|\vec{w}\|_2^2$

Hyperparameter.

How can we use probabilistic information about our data?

How does this connect to optimization models?

$(\vec{x}_i, y_i)$  are my data points.

$y_i = g(\vec{x}_i) + z_i$  iid Independent  
Identically  
distributed.

$z_i \sim N(0, \sigma_i^2)$

$f(z_i) = \frac{e^{-z_i^2/2\sigma_i^2}}{\sqrt{2\pi} \sigma_i}$

Consider  $\otimes$  linear model:

$y_i = \vec{x}_i^T \vec{w} + z_i$

$\vec{w}$  is "our model"  
 What we want to learn. (unknown)

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_n^T \end{bmatrix} \vec{w} + \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{bmatrix}$$

$$\vec{y} = X \vec{w} + \vec{z}$$

Probabilistic solution: Maximum likelihood estimator.

Find that  $\vec{w}$  that makes the observed data most likely.

argmax  $\vec{w}_0$   $f_{y_1, y_2, \dots, y_n} (y_1 = y_1, y_2 = y_2, \dots, y_n = y_n \mid \vec{w} = \vec{w}_0)$  (Maximum Likelihood)

= argmax  $\vec{w}_0$   $\prod_{i=1}^n f(y_i = y_i \mid \vec{w} = \vec{w}_0)$  (Because all of my ~~z\_i~~  $z_i$ 's are independent.)

Consider:  $f(y_i = y_i \mid \vec{w} = \vec{w}_0) = f(\vec{x}_i^T \vec{w}_0 + z_i = y_i \mid \vec{w} = \vec{w}_0)$   
 $= f(z_i = y_i - \vec{x}_i^T \vec{w}_0 \mid \vec{w} = \vec{w}_0) = \frac{e^{-(y_i - \vec{x}_i^T \vec{w}_0)^2 / 2\sigma_i^2}}{\sqrt{2\pi} \sigma_i}$

$$\underset{\vec{w}_0}{\text{argmax}} \frac{\prod_{i=1}^n e^{- (y_i - \vec{x}_i^T \vec{w}_0)^2 / 2\sigma_i^2}}{\sqrt{2\pi} \sigma_i}$$

$$= \underset{\vec{w}_0}{\text{argmax}} \frac{1}{(\sqrt{2\pi})^n} \frac{1}{\prod_{i=1}^n \sigma_i} \exp \left\{ \sum_{i=1}^n - (y_i - \vec{x}_i^T \vec{w}_0)^2 / 2\sigma_i^2 \right\}$$

$$= \underset{\vec{w}_0}{\text{argmin}} \sum_{i=1}^n (y_i - \vec{x}_i^T \vec{w}_0)^2 / 2\sigma_i^2$$

$$S^2 = \begin{bmatrix} \frac{1}{2\sigma_1^2} & & & 0 \\ & \frac{1}{2\sigma_2^2} & & \\ & & \dots & \\ 0 & & & \frac{1}{2\sigma_n^2} \end{bmatrix}$$

$$= \underset{\vec{w}_0}{\text{argmin}} \underbrace{\| S (\vec{y} - X \vec{w}_0) \|^2}_{\text{Weighted Least Squares.}}$$

$$S = \begin{bmatrix} \frac{1}{\sqrt{2}\sigma_1} & & & 0 \\ & \dots & & \\ & & \dots & \\ 0 & & & \frac{1}{\sqrt{2}\sigma_n} \end{bmatrix}$$

what if we had a prior on  $\vec{w}$ ?

"side-information"

MAP: Maximum a. posteriori

$$y_i = \vec{x}_i^T \vec{w} + z_i$$

$$z_i \sim N(0, \sigma_i^2)$$

$$w_i \sim N(\mu_i, \beta_i^2) \quad \text{"Prior" on } \vec{w}$$

$$\star \vec{w} \sim N(\vec{\mu}, \Sigma_w)$$

$$\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

$$\Sigma_w = \begin{bmatrix} \beta_1^2 & & & 0 \\ & \beta_2^2 & & \\ & & \ddots & \\ 0 & & & \beta_n^2 \end{bmatrix}$$

$$\operatorname{argmax}_{\vec{w}} f(\vec{w} | Y_1=y_1, Y_2=y_2, \dots, Y_n=y_n) \quad (*)$$

What is the most likely  $\vec{w}$ , given the data?

$$f(\vec{w} | Y_1=y_1, \dots, Y_n=y_n) = \frac{f(Y_1=y_1, Y_2=y_2, \dots, Y_n=y_n | \vec{w}) \cdot f(\vec{w})}{f(Y_1=y_1, Y_2=y_2, \dots, Y_n=y_n)} \quad \text{(Bayes Rule)}$$

→ does not depend on  $\vec{w}$ .

$$(*) \text{ MAP} = \operatorname{argmax}_{\vec{w}} f(Y_1=y_1, \dots, Y_n=y_n | \vec{w}) \cdot f(\vec{w})$$

$$\vec{Y} = \vec{y}$$

$$= \operatorname{argmax}_{\vec{w}} f(\vec{Y}=\vec{y} | \vec{w}) f(\vec{w})$$

$$= \operatorname{argmax}_{\vec{w}} \left( \prod_{i=1}^n f(Y_i=y_i | \vec{w}) \right) f(\vec{w}).$$

$$= \operatorname{argmax}_{\vec{\omega}} \left( \prod_{i=1}^n \frac{\exp\left(-\frac{(\vec{x}_i^T \vec{\omega} - y_i)^2}{2\sigma_i^2}\right)}{\sqrt{2\pi} \cdot \sigma_i} \right) \cdot \frac{e^{-\frac{(\vec{\omega} - \vec{\mu})^T \Sigma_{\omega}^{-1} (\vec{\omega} - \vec{\mu})}{2\sigma_i}}}{(\sqrt{2\pi})^n (\prod \sigma_i)}$$

$$= \operatorname{argmax}_{\vec{\omega}} \cdot \exp \left\{ \sum_{i=1}^n -\frac{(\vec{x}_i^T \vec{\omega} - y_i)^2}{2\sigma_i^2} + -(\vec{\omega} - \vec{\mu})^T \Sigma_{\omega}^{-1} (\vec{\omega} - \vec{\mu}) \right\}$$

$$= \operatorname{argmin} \quad \left\| S(X\vec{\omega} - \vec{y}) \right\|_2^2 + \left\| \sqrt{\Sigma_{\omega}^{-1}} (\vec{\omega} - \vec{\mu}) \right\|_2^2$$

like the  $\lambda$  terms.

What happens if  $\sigma_i$  is large?

choose less penalty for deviation from the mean.

# Principal Components Regression.

$$X \in \mathbb{R}^{m \times n}$$

$X$  is full column rank.

$$\min \|X\vec{w} - \vec{y}\|_2^2$$

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

$$\text{LS: } \vec{\hat{w}} = (X^T X)^{-1} X^T \vec{y}$$

$$= ((U \Sigma V^T)^T (U \Sigma V^T))^{-1} \cdot (U \Sigma V^T)^{-1} \cdot \vec{y}$$

⋮

Do usual math.

$$= V \begin{array}{|ccc|c} \frac{1}{\sigma_1} & & 0 & 0 \\ & \dots & & \\ & & \frac{1}{\sigma_n} & \\ 0 & & & \end{array} U^T \vec{y}.$$

For PCR: Only consider top  $k$  principal components instead of all of  $X$ .

# "Ridge regression as soft PCA"

$\vec{w}$  in the  $V$  basis.

$$\vec{w} = V \cdot \vec{z}$$

$$\underset{\vec{w}}{\text{argmin}} \quad \|X\vec{w} - \vec{y}\|_2^2 + \lambda \|\vec{w}\|_2^2$$

$$= \underset{\vec{w}=V\vec{z}}{\text{argmin}} \quad \|X \cdot V \cdot \vec{z} - \vec{y}\|_2^2 + \lambda \|V \cdot \vec{z}\|_2^2$$

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

$$= \underset{\vec{z}}{\text{argmin}} \quad \underbrace{\|XV\vec{z} - \vec{y}\|_2^2}_A + \lambda \|\vec{z}\|_2^2 \quad (\text{Ridge}).$$

$$\begin{aligned} \vec{z}_{\text{ridge}} &= ((XV)^T(XV) + \lambda I)^{-1} (XV)^T \vec{y} \\ &= (V^T X^T X V + \lambda I)^{-1} (XV)^T \vec{y} \end{aligned}$$

$$(A^T A + \lambda I) A^T \vec{b}$$

~~$$(V^T V \Sigma^T \Sigma U^T U)$$~~

$$= (V^T (U \Sigma V^T)^T (U \Sigma V^T) V + \lambda I)^{-1} (XV)^T \vec{y} \quad \text{HW: do this cancellation.}$$

$$= \left( \begin{array}{c} \Sigma^T \Sigma + \lambda I \\ n \times m \quad m \times n \end{array} \right)^{-1} \Sigma^T U^T \vec{y} = \left( \begin{array}{c|c} \sigma_1^2 + \lambda & \\ \sigma_2^2 + \lambda & \\ \vdots & \\ \sigma_n^2 + \lambda & \end{array} \right)^{-1} \left( \begin{array}{c|c} \sigma_1 & \\ \sigma_2 & \\ \vdots & \\ \sigma_n & \end{array} \right) U^T \vec{y}$$