

# Parameter estimation for linear regression

$$y_i = \vec{w}^T \vec{x}_i + \epsilon \quad \text{where } \epsilon \sim \mathcal{N}(0, \sigma^2)$$

Here,  $y_i$  is a real scalar and  $\vec{w}$  and  $\vec{x}_i$  are vectors of dimensions  $d \times 1$

Given  $n$  pairs of  $(x_i, y_i)$ , we can estimate  $\vec{w}$  that best fits the model.

These notes cover the following methods of parameter estimation:

- (1) least squares estimation
  - (2) M C L E
  - (3) Regularized least squares estimation
  - (4) MAP estimate
- } Similar or equivalent
- } Similar or equivalent

# (1) Least Squares Estimation

We consider the sum of squared errors as our loss function:

$$f(\vec{w}) = (\vec{y} - X\vec{w})^T (\vec{y} - X\vec{w})$$

$\downarrow$        $\downarrow$        $\downarrow$   
 $n \times 1$     $n \times d$     $d \times 1$

$$= (\vec{y}^T - \vec{w}^T X^T) (\vec{y} - X\vec{w})$$

$$= (\vec{y}^T \vec{y} - \vec{w}^T X^T \vec{y} - \vec{y}^T X \vec{w} + \vec{w}^T X^T X \vec{w})$$

Same scalar value

$$= (\vec{y}^T \vec{y} - 2\vec{y}^T X \vec{w} + \vec{w}^T X^T X \vec{w})$$

$$\nabla_{\vec{w}} f = 0 - 2X^T \vec{y} + 2X^T X \vec{w}$$

Put to zero:

$$2X^T X \vec{w} = 2X^T \vec{y}$$

$$\boxed{\therefore \vec{w} = (X^T X)^{-1} X^T \vec{y}}$$

## (2) Maximum Conditional Likelihood Estimation

$$y_i = \vec{w}^T x_i + \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\therefore y \sim \mathcal{N}(\vec{w}^T x, \sigma^2)$$

$$\text{MLE} \equiv \arg \max_{\vec{w}} \prod_{i=1}^n P(y_i | x_i, \vec{w})$$

$$\log l(\vec{w}) = \log \left[ \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma^2} \times e^{-\frac{1}{2} \left( \frac{y_i - \vec{w}^T x_i}{\sigma} \right)^2} \right]$$

$$= \sum_{i=1}^n \left[ -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2} \left( \frac{y_i - \vec{w}^T x_i}{\sigma} \right)^2 \right]$$

Ignoring the 1<sup>st</sup> term because it does not involve  $\vec{w}$

$$\log l(\vec{w}) = \sum_{i=1}^n \left[ -\frac{1}{2\sigma^2} \times (y_i - \vec{w}^T x_i)^2 \right]$$

$$= -\frac{1}{2\sigma^2} \times (\vec{y} - X\vec{w})^T (\vec{y} - X\vec{w})$$

$\therefore$  We are maximizing <sup>the negative of</sup> a function proportional to sum of squared errors

$\therefore$  MLE view is equivalent to sum of squared errors estimation.

### (3) Regularized least squares estimation

We now consider our loss function to be sum of squared errors PLUS square of  $l_2$  norm of the parameter vector  $\vec{w}$ :

$$f(\vec{w}) = (\vec{y} - X\vec{w})^T (\vec{y} - X\vec{w}) + \frac{\lambda}{2} \|\vec{w}\|_2^2$$

Here,  $\lambda$  is a parameter that helps decide the importance attached to minimized the squares of the weight vector  $\vec{w}$

$$\nabla_{\vec{w}} f = 0 - X^T \vec{y} - X^T \vec{y} + 2X^T X \vec{w} + \lambda \vec{w}$$

Put to zero:

$$2X^T X \vec{w} + \lambda \vec{w} = 2X^T \vec{y}$$

$$\boxed{\therefore \vec{w} = (X^T X + \frac{\lambda}{2} I)^{-1} X^T \vec{y}}$$

#### (4) MAP estimate for linear regression

$$P(\vec{w}|D) = \frac{P(D|\vec{w}) \times P(\vec{w})}{P(D)} \rightarrow \begin{array}{l} \text{likelihood} \\ \text{Here, } D \text{ is data} \end{array}$$

$$\begin{aligned} \vec{w}_{\text{MAP}} &= \arg \max_{\vec{w}} \left[ \ln[P(D|\vec{w})] + \ln[P(\vec{w})] - \ln[P(D)] \right] \\ &= \arg \max_{\vec{w}} \left[ \ln[P(\vec{w})] + \ln[P(D|\vec{w})] \right] \end{aligned}$$

Let us assume that  $\vec{w}$  is drawn from a

Gaussian  $\delta$

$$\vec{w} \sim \mathcal{N}(\vec{0}, \frac{1}{\lambda} \times I) \quad \text{where } I \text{ is identity matrix}$$

$$\therefore \vec{w}_{\text{MAP}} = \arg \max_{\vec{w}} \left[ \ln \left[ \frac{1}{(2\pi)^{d/2} |I/\lambda|^{1/2}} \times e^{-\frac{1}{2} (\vec{w}^T (\lambda I) \vec{w})} \right] + \ln[P(D|\vec{w})] \right]$$

$$= \arg \max_{\vec{w}} \left[ \ln \left[ \frac{1}{(2\pi)^{d/2} |I/\lambda|^{1/2}} \right] - \frac{\lambda}{2} \vec{w}^T \vec{w} + \ln[P(D|\vec{w})] \right]$$

$$\begin{array}{c} \downarrow \\ \text{Ignore} \end{array} = \arg \max_{\vec{w}} \left[ -\frac{\lambda}{2} \|\vec{w}\|^2 + \text{log-likelihood} \right]$$

$\equiv$   $l_2$  regularization for linear regression

For  $\sigma =$  some constant, this is equivalent to  $l_2$  regularization