

Bayesian Regression

Mon Apr 12

Linear Discriminant Analysis

Regression

Discriminative model

✓
 $p(y|x)$

vs.

Generative model

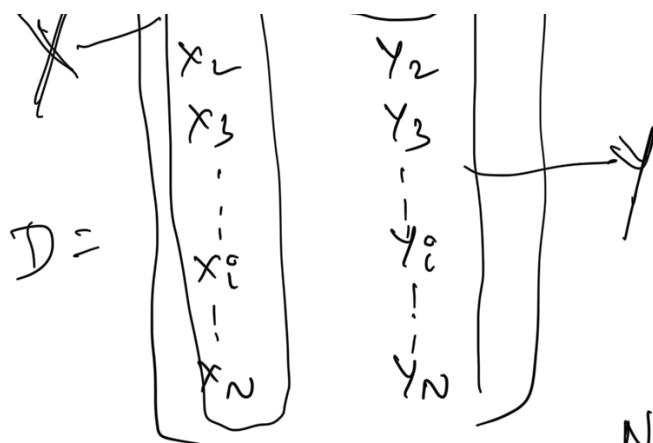
$p(y)$ $p(x|y)$

$w - ?$

$$y|x, w = \mathcal{N}(w^T x, \sigma^2)$$

↘ scalar

$\boxed{x_i}$ $\boxed{y_i}$



$$L(D|w) = \prod_{i=1}^N p(y_i | x_i, w)$$

$$LL(D|w) = \sum_{i=1}^N \log p(y_i | x_i, w)$$

$$= \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \right] \right]$$

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

$$= \underbrace{-\frac{N}{2} \log 2\pi - N \log \sigma}_{\text{constant}} - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - w^T x_i)^2$$

Find (w, σ^2) that maximizes $LL(D|w)$

$$\frac{\partial LL(D|w)}{\partial w} = 0 \quad \bigg| \quad \frac{\partial LL(D|w)}{\partial \sigma^2} = 0$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0$$

$$\hat{\mathbf{w}}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{N} (\mathbf{y} - \mathbf{X} \hat{\mathbf{w}}_{MLE})^T (\mathbf{y} - \mathbf{X} \hat{\mathbf{w}}_{MLE})$$

Putting a prior on \mathbf{w}

\mathbf{w} is a $(D+1)$ length vector

$$p(\mathbf{w}) \sim \mathcal{N}(\mathbf{w} | \mu_0, \Sigma_0)$$

$$p(\mathbf{w} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathbf{w}) p(\mathbf{w})}{\int_{\mathbf{w}'} p(\mathcal{D} | \mathbf{w}') p(\mathbf{w}') d\mathbf{w}'}$$

Posterior ~~is~~ of \mathbf{w} will also be a Gaussian.

$$p(\mathbf{w}) \sim \mathcal{N}(\mathbf{w} | \mathbf{0}, \tau^2 \mathbf{I})$$

↘ scalar

Special but often-used prior on \mathbf{w}

Posterior:

$$\begin{aligned}\bar{\mathbf{w}} &= \left(\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y} \\ \bar{\Sigma} &= \sigma^2 \left(\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I} \right)^{-1}\end{aligned}$$

Think ridge regression