# Bayesian Regression

Linear Discriminant Analysis

Regression

---

Discriminative Model $\boxed{p(y|x)}$ ✓

vs.

Generative Model

$$p(y) \quad p(x|y)$$
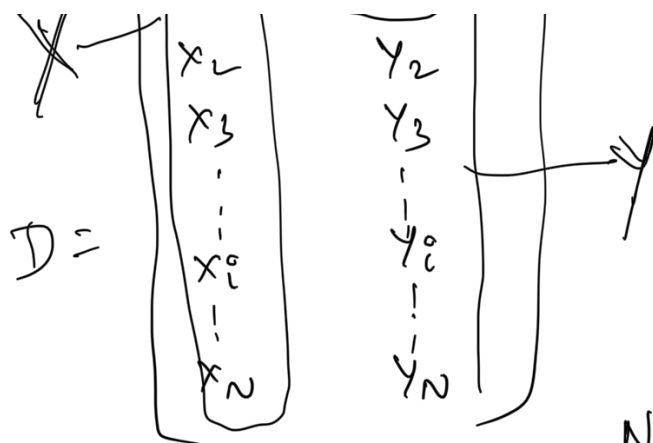
---

$$w - ?$$

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

$\rightarrow$ scalar

---

$\begin{bmatrix} x_1 \end{bmatrix}$ $(y_1)$ $\begin{bmatrix} \\ \end{bmatrix}$

$$D = \begin{pmatrix} X \end{pmatrix} \begin{bmatrix} x_2 \\ x_3 \\ \vdots \\ x_i^\sigma \\ \vdots \\ x_N \end{bmatrix} \begin{bmatrix} y_2 \\ y_3 \\ \vdots \\ y_i^\sigma \\ \vdots \\ y_N \end{bmatrix} \quad \rightarrow y$$

$$L(D|w) = \prod_{i=1}^{N} p(y_i | \cancel{x}_i; w)$$

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

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

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

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

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

$$\frac{\partial\, LL(D|w)}{} = 0 \qquad \Bigg| \qquad \frac{\partial LL(D|w)}{} = 0$$

$$\partial w \qquad\qquad | \quad \partial \sigma$$

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

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

Putting a prior on **w**

**w** is a $(D+1)$ length vector

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

$$p(w|D) = \frac{p(D|w)\, p(w)}{\int_{w'} p(D|w')\, p(w')\, dw}$$

Posterior ~~will~~ of $w$ will also be a Gaussian.

$$p(w) \sim \mathcal{N}(w | 0, \nu^2 I)$$
$\longrightarrow$ scalar

Special but often-used
prior on $w$

## Posterior:

$$\overline{w} = \left( X^T X + \frac{\sigma^2}{\gamma^2} I \right)^{-1} X^T y$$

$$\overline{\overline{\Sigma}} = \sigma^2 \left( X^T X + \frac{\sigma^2}{\gamma^2} I \right)^{-1}$$

### Think ridge regression

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

( Generalized linear Model )

$$pdf(y|x) = \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left[ -\frac{1}{2} \frac{(y_i - w^T x_i)^2}{\sigma^2} \right]$$

$$Laplace(y|x) = \frac{1}{\exp} \| y - w^T x_i \|$$

Estimating $w$ using Laplace
distribution is more
...

$$y|x \sim \text{Bernoulli} ( \theta )$$

$$0 \le \theta \le 1$$

Binary classification.

$$\theta = \sigma(w^T x)$$

$$= \frac{1}{1 + \exp(-w^T x)}$$

---

If $y_i = 1$

$$p(y_i) = \theta_i$$

$$= \frac{1}{1 + \exp(-w^T x_i)}$$

If $y_i = 0$

$$p(y_i) = 1 - \theta_i$$

$$= \frac{1}{1 + \exp(w^T x_i)}$$

In general

$$p(y_i) = \frac{1}{1 + \exp(-y_i w^T x_i)}$$

$$l(w) = \prod^N \left[ \frac{1}{\phantom{1 + \exp}} \right]$$

$$\sum_{i=1}^{'} \left[ 1 + \exp\left(-y_i w' x_i\right) \right]$$

$$\boxed{LL(w) = -\sum_{i=1}^{N} \log\left(1 + \exp\left(-y_i w^T x_i\right)\right)}$$

No closed-form solution

Gradient Descent

Regularizing LR

$$w \sim N(0, \tau^2 I)$$

$$p(w \mid X, y) \propto p(w)\, p(X, y \mid w)$$

Multinoulli

$$X \in \{1, 2, ---, C\}$$

$$\Theta = [0.1, 0.2, 0.05, \cdots\cdots 0.2]$$

$$\Theta_j = \frac{\exp(w_j^T x)}{}$$

$$\sum_{k=1}^{c} \exp\left(w_k^{i} x\right)$$