# Probabilistic Interpretation of linear models.

$$y = w^T x$$

$\rightarrow$ weight vector

---

## Probabilistic linear regression

$y$ is a random variable

$\cancel{x}$ — is not a random variable,

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

$$\boxed{w, \sigma^2}$$

---

$$\underline{y = w^T x + \epsilon}$$
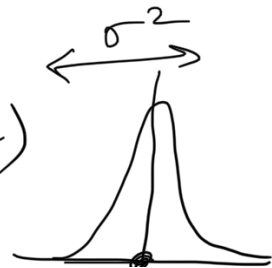
$\uparrow$ random variable

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

---

If we have $w$ and $\sigma^2$

Given a new $x^*$,

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

### Training data

$$D = \begin{bmatrix} x_1 & & Y_1 \\ x_2 & & Y_2 \\ \vdots & & \vdots \\ x_N & & Y_N \end{bmatrix}$$

$$Y_1 \sim \mathcal{N}(w^T x_1, \sigma^2)$$
$$Y_2 \sim \mathcal{N}(w^T x_2, \sigma^2)$$
$$\vdots$$
$$Y_N \sim \mathcal{N}(w^T x_N, \sigma^2)$$

2

Likelihood of the dataset:

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

$$\ell\ell(D) = \sum_{i=1}^{N} \log p(y_i)$$

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

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

MLE estimate $\equiv$ the values of $w$ & $\sigma^2$ at which $\ell\ell(D)$ is max.

$$\underset{w,\sigma^2}{\arg\max}\ \ell\ell(D)$$

$$\ell\ell(D) = const - N\log\sigma - \left(\frac{1}{2\sigma^2}\right)\sum_{i=1}^{N}(y_i - w^T x_i)^2$$

$\underline{\underline{W}}$

$\overbrace{\qquad}$ $\overbrace{\qquad}$ $\overbrace{(\quad)}^{\text{Cost}/} \Big|_{i=1}$

Equivalent to maximizing: $-\frac{1}{2} \sum (y_i - w^T x_i)^2$

which is equiv. to minimizing $\frac{1}{2} \sum (y_i - w^T x_i)^2$

squared loss for geometric linear regression.

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

where $X \to$ data matrix $\underline{N \times d}$

$y \to$ vector of target values $\underline{N \times 1}$

$$\hat{\sigma}^2_{MLE} = \frac{1}{N} \sum (y_j - w^T x_j)^2$$

$$= \frac{1}{N} (y - Xw)^T (y - Xw)$$

## Imposing a prior on $w$

$w \to$ a d-dimensional vector

MVN or a multivariate Gaussian

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

Prior $\qquad\qquad\qquad d \times 1 \qquad d \times d$

Simple case: $\mu_0 = \underset{d \times 1}{0}$

$$\Sigma_0 = \gamma^2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= \gamma^2 I \quad (d \times d)$$

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

$$= \frac{p(X, y \mid w) \, p(w)}{\int p(X, y \mid w') \, p(w') \, d(w')}$$

$X, y \to$ Training

$\bar{w}$

$$p(X, y \mid w) = \prod_{i=1}^{N} N(y_i \mid w^T x_i, \sigma^2)$$

Posterior is also Gaussian $p(w \mid X, y)$

$\sim N(w \mid \bar{w}, \bar{\Sigma})$

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

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

$$y \sim N(y \mid w^T x, \sigma^2)$$
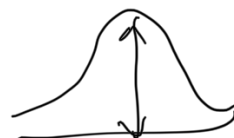
$\sigma^2$ - fixed and known

# Relationship with ridge regression

$$W_{Ridge} = (X^T X + \lambda I)^{-1} X^T y$$

regularization parameter

$$\bar{W} = (X^T X + \frac{\sigma^2}{\gamma^2} I)^{-1} X^T y$$

What is the MAP estimate for $\omega$

$$\boxed{\hat{W}_{MAP} = \bar{W}}$$

Inference or prediction task:

$$x^* \longrightarrow y^*$$

$$y^* = \hat{W}_{MLE}^T x^*$$

$$y^* = \hat{W}_{MAP}^T x^*$$

Full Bayesian Treatment.

CSE610 - fall 2020

Non-parametric Bayesian Methods

$$y \sim N(y \mid w^T x, \sigma^2)$$

$$y \sim F(y \mid \theta = f(w^T x))$$

Generalized linear Models. (GLM)

Replace $\quad N(\ ) \to$ Laplace $(\ )$

Robust Regression

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

Impacted by outliers

$$p(y) = \text{const} \exp\left[-\text{Cst}|y_i - w^T x|\right]$$

$$\frac{1}{2}\sum_{i=1}^{N}(y_i - w^T x_i)^2 \qquad \text{least squares}$$

least absolute
MAD. Deviation

$$\sum_{i=1}^{N} \cdots w^T x \mid$$

$$\frac{1}{2} \sum_{i=1}^{} |y_i - w \, x_i)$$

logistic Regression    is a member
of GLM family.

## Probabilistic Logistic Regression

① **Generalized Linear Models (GLM)**

② **Bayesian Logistic Regression**

③ **Handling multiple classes**

---

### GLM

$$w^T x$$

$$\boxed{y \sim F(f(w^T x), --)}$$

**Classification**  $y \in \{0, 1\}$

Bernoulli distribution

$$y \sim Ber(y \mid \theta)$$

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

$$0 \leq \theta \leq 1$$

$$x^* \longrightarrow w^T x^*$$

$$\text{sigmoid}(w^T x^*)$$

If Sigmoid $(w^Tx^*) \geq 0.5$ $y^* = 1$

else $y^* = 0$

Given training data, $D$ $\langle x_i, y_i \rangle_{i=1}^{N}$

learn $w$

MLE: $x_i, \boxed{y_i} \rightarrow$ random variable

$$L(D) = \prod_{i=1}^{N} P(Y = y_i)$$

$P(Y = y_i)$?

If $y_i = 1$, then $P(Y = y_i) = \theta_i = \dfrac{1}{1 + \exp(-w^Tx_i}$

If $y_i = 0$, then $P(Y = y_i) = 1 - \theta_i$

$$= 1 - \left( \dfrac{1}{1 + \exp(-w^Tx_i)} \right)$$

$$= \dfrac{1}{1 + \exp(w^Tx_i)}$$

In general
$$P(Y = y_i) = \theta_i^{y_i} (1 - \theta_i)^{1 - y_i}$$

$$L(D) = \prod_{i=1}^{N} \theta_i^{y_i} (1-\theta_i)^{1-y_i}$$

$$LL(D) = \sum_{i=1}^{N} y_i \log \theta_i + (1-y_i) \log(1-\theta_i)$$

where $\quad \theta_i = \dfrac{1}{1 + \exp(-w^T x_i)}$

We can maximize this LL(D) w.r.t w

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

No closed-form expression.

Have to use a gradient based method.

→ Gradient Descent

→ Newton's method

Regularization and prior are some what equivalent.

$p(w)$

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

$\longrightarrow LL(v)$

$\quad\quad\longrightarrow p(w|D)$

---

I could get the MAP estimate for $w$ by adding a L2-penalty to the $LL(w)$

But it is not easy to get $p(w|D)$

posterior.

---

Generalize to multi-class classification.

$$y \in \{1, 2, 3, \cdots, 10\}$$

Bernoulli X

Multinoulli $\rightarrow$ 

$$\boxed{P(y = k) \quad 1 \leq j \leq C}$$

$$\theta_j = P(y = j) = \boxed{\frac{\exp(w_j^T x)}{\sum_{k=1}^{C} \exp(w_\cdot^T x)}}$$

$C$ weight vector $\begin{matrix} w_1 \\ w_2 \\ | \\ h \end{matrix}$

$wc$

Training $\langle x_i, y_i \rangle_{i=1}^{N}$

$$\prod_{i=1}^{N} P(y_i)$$

$w \sim N\langle\ \rangle_{\alpha}$   $\dfrac{p(w)\,p(D|w)}{\int_{w'} p(w')\,p(D|w')\,dw'}$

$Ber(\ )$

$p(w|D)$ is not easy