

Probabilistic Interpretation of linear models.

$$\boxed{y = w^T x}$$

↘ weight vector

Probabilistic linear regression

y is a random variable

x — is not a random variable,

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

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

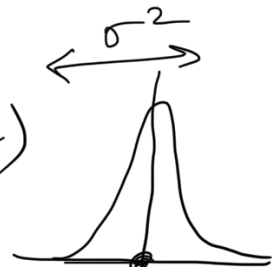
$$y = \underbrace{w^T x + \epsilon}_{\text{random variable}}$$

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

If we have w and σ^2

Given a new x^* ,

$$y^* = \mathcal{N}(\underbrace{w^T x^*}_{\text{mean}}, \sigma^2)$$



$$\begin{array}{c}
 \text{Training data} \\
 D = \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \\ \vdots & \vdots \\ x_N & y_N \end{bmatrix} \quad \begin{array}{l} y_1 \sim \mathcal{N}(w^T x_1, \sigma^2) \\ y_2 \sim \mathcal{N}(w^T x_2, \sigma^2) \\ \vdots \\ y_N \sim \mathcal{N}(w^T x_N, \sigma^2) \end{array}
 \end{array}$$

Likelihood of the dataset:

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

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

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

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

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

$$\arg \max_{w, \sigma^2} \ell L(D)$$

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

$\underbrace{w}_{\text{const}}$
 Equivalent to maximizing: $-\frac{1}{2} \sum_{i=1}^n (y_i - w^T x_i)^2$
 which is equiv. to minimizing $\frac{1}{2} \sum_{i=1}^n (y_i - w^T x_i)^2$

Squared loss for
geometric linear
regression.

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

where $X \rightarrow$ data matrix $N \times d$

$y \rightarrow$ vector of target values $N \times 1$

$$\begin{aligned}
 \hat{\sigma}_{MLE}^2 &= \frac{1}{N} \sum (y_i - w^T x_i)^2 \\
 &= \frac{1}{N} (y - Xw)^T (y - Xw)
 \end{aligned}$$

Set a prior distribution for w $p(w)$

Use the likelihood $p(D|w) \equiv p(x, y|w)$

to calculate the posterior distribution for

$$w: \underline{p(w|D)} \equiv p(w|x, y)$$

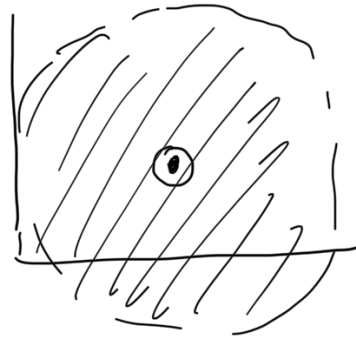
w is a $d \times 1$ vector.

$$W \sim \mathcal{N}(\underbrace{\mu_0}_{d \times 1}, \underbrace{\Sigma_0}_{d \times d})$$

Encodes prior information.

Example let $d=2$

$$\underline{p(w)} \sim \mathcal{N}(\mu_0 = (3, 2), \Sigma_0 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix})$$



$$p(w|D) \equiv p(w|x, y)$$

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

Simple prior

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

↑ variance

$$\begin{pmatrix} \tau^2 & 0 \\ 0 & \tau^2 \end{pmatrix}$$

$$p(w|x, y) \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu = \left(X^T X + \frac{\sigma^2}{\gamma^2} I_D \right)^{-1} X^T y \rightarrow \hat{w}_{MAP}$$

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

looks a little bit like ridge-regression

$$J(w) = \frac{1}{2} (y - Xw)^T (y - Xw) + \frac{\lambda}{2} w^T w$$

$$w = (X^T X + \lambda I_D)^{-1} X^T y$$

For a general setting, when

$$\begin{cases} \mu_0 \neq 0 \\ \Sigma_0 \neq \gamma^2 I_D \end{cases}$$

How to predict y^* for a given x^*

$$① y^* \sim \mathcal{N}(\hat{w}_{MLE}^T x^*, \sigma^2)$$

$$② y^* \sim \mathcal{N}(\hat{w}_{MAP}^T x^*, \sigma^2) \quad \underline{p(w|x, y)}$$

For a Gaussian r.v.

Mode \equiv Mean

③ Bayesian

Things to take away:

- ① Prob. linear regression $y^* \sim \mathcal{N}(w^T x, \sigma^2)$
 - ② $\hat{w}_{MLE} \equiv w_{\text{least squares}}$
 - ③ \hat{w}_{MAP} using a $\mathcal{N}(0, \tau^2 I)$ prior on w ,
is same as ridge regression estimate.
-

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

$$\left[\exp \frac{-1}{2\sigma^2} (y_i - w^T x_i)^2 \right]$$

\mathcal{N} can be replaced by other distributions

Generalized Linear Models (GLM)

$$y|x \sim \text{Laplace}(\underline{w^T x}, b)$$

↳ Robust regression.

$$(y_i - w^T x_i)$$
