

# 1. Linear Classification with a Margin

- Show that a dataset with two point is sufficient to determine the maximum-margin

Hint:  $x_1 \in C_1 [y_1 = +1]$   $x_2 \in C_2 [y_2 = -1]$

- Setup minimization problem with constraints on  $w^T x_i + b \dots$  Subject to:

$$w^T x_1 + b = +1 \quad \text{for optimization problem:}$$

$$w^T x_2 + b = -1$$

$$\underset{w \in \mathbb{R}^p}{\operatorname{argmin}} \|w\|_2^2$$

- What is  $w$ ?  $b$ ?

Lagrange function and lagrange multipliers:  $L(\lambda, w) = f(w) - \lambda g(x)$

$$L(\lambda, w) = w^T x_1 + b - 1 \quad \lambda \geq 0$$

$$g(x_2) = w^T x_2 + b + 1$$

$$L(\lambda_1, \lambda_2, w, b) = \underset{w}{\operatorname{argmin}} \left[ \|w\|_2^2 + \lambda_1 (w^T x_1 + b - 1) + \lambda_2 (w^T x_2 + b + 1) \right]$$

① Take derivatives:

$$\frac{\partial L}{\partial w} = 2w + \lambda_1 x_1 + \lambda_2 x_2 = 0; \quad \frac{\partial L}{\partial b} = \lambda_1 + \lambda_2 = 0$$

② Find solution for terms:

$$\lambda_1 = -\lambda_2$$

$$0 = 2w + \lambda_1 x_1 + \lambda_2 x_2$$

$$= 2w - \lambda_2 x_1 + \lambda_2 x_2$$

$$= 2w + \lambda_2 (x_2 - x_1)$$

$$-2w = -\lambda_2 (x_1 - x_2)$$

$$w = \frac{\lambda_2 (x_1 - x_2)}{2}$$

$$w^T x_1 + b = 1 \quad w^T x_2 + b = -1$$

$$w^T x_1 + b = -(w^T x_2 + b)$$

$$w^T x_1 + b = -w^T x_2 - b$$

$$2b = -w^T x_1 - w^T x_2$$

$$b = -\frac{w^T}{2} (x_1 + x_2)$$

$$b = -\frac{\lambda_2}{2} (x_1 - x_2)^T (x_1 + x_2)$$

Constraints:

$$b = -\frac{\lambda_2}{4} (x_1 - x_2)^T (x_1 + x_2)$$

## Problem 2

- Find the result of minimizing the loss of sum of squared errors.
- Add a penalty for an L2 penalty

$$\arg \min_w \left\{ \sum (w^T x_i - y_i)^2 + \lambda \|w\|_2^2 \right\}$$

- How does this change the solution to the original regression solution?
- What is the impact of adding the penalty?

$$\text{Logistic Function: } \frac{1}{1 + \exp(-w^T x_i)}$$

$$\text{Cross Entropy: } L(w) = - \sum_{i=1}^n \left\{ y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right\}$$

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

$$\arg \min_w \left\{ \sum (w^T x_i - y_i)^2 + \lambda \|w\|_2^2 \right\}$$

$$(w x - y)(w x - y)^T \quad \lambda w^T w$$

$$L(w) = (w x - y)^T (w x - y) + \lambda w^T w$$

$$= \bar{w}^T X^T \bar{w} X - \bar{w}^T X^T \bar{y} - \bar{y}^T \bar{w} X + \bar{y}^T \bar{y} + \lambda \bar{w}^T \bar{w}$$

$$L(\bar{w}) = \bar{w}^T X^T \bar{w} X - 2 \bar{y}^T X \bar{w} + \bar{y}^T \bar{y} + \lambda \bar{w}^T \bar{w}$$

$$\frac{\partial L}{\partial w} = 2 X^T X \bar{w} - 2 \bar{y}^T X + 0 + 2 \lambda \bar{w} = 0$$

$$2 (X^T X \bar{w} - \bar{y}^T X + \lambda \bar{w}) = 0$$

$$\bar{w} (X^T X + \lambda I) - \bar{y}^T X = 0$$

① Rework equation

$$A^T \bar{x}^T \bar{b} + A \bar{x} \bar{b}^T = 2 \bar{b}^T A \bar{x}$$

② Take derivative w.r.t w

$$\frac{d}{dw} (2 \bar{y}^T X \bar{w}) = 2 \bar{y}^T X$$

$$\frac{d}{dw} (X^T \bar{w}^T X \bar{w}) = 2 X^T X \bar{w}$$

$$w(X^T X + \lambda I) = \bar{y}^T X$$

$$\bar{w} = \frac{\bar{y}^T X}{(X^T X + \lambda I)}$$

③ New solution with penalty

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

$$\bar{w}^T X^T \bar{w} X + 2 \bar{y}^T X \bar{w} + y^T y$$

④ Original solution without

penalty

$$\frac{\partial L}{\partial w} = \cancel{2} X^T X \bar{w} - \cancel{2} y^T X = 0 \rightarrow \bar{w} = \frac{\bar{y}^T X}{X^T X}$$

The add penalty prevents  $\bar{w}$  values from getting too high.

4.

- The Bayes decision rule describes an approach we take for choosing a  $w$  for  $x$ .
- Can be achieved by modeling  $P(w|x)$ , or  $\frac{P(x|w)P(w)}{P(x)}$
- What are the Pros/Cons for each approach?
- Why is know  $P(x)$  useful?

$P(w|x)$  :- Directly estimate the probability data belongs to a class, however difficult to determine outliers.

$\frac{P(x|w)P(w)}{P(x)}$  :- Obtain useful information of the data by finding the prior, likelihood, evidence, probabilities

$P(x)$  :- Since the likelihood  $P(x|w)$  can be  $> 1$ ,  $P(x)$  acts as a normalize factor to give an appropriate probability for the posterior