

① Linear classifier with a Margin :

a data set has 2 data points

$$x_1 \in C_1 (y_1 = +1)$$

$$x_2 \in C_2 (y_2 = -1)$$

Set up the minimization problem w/

constraints on $W^T x_1 + b$

$$W^T x_2 + b$$

To find the hyperplane, we need to solve

$$\arg \min_{W \in \mathbb{R}^P} \|W\|_2^2 = \arg \min_{W \in \mathbb{R}^P} W^T W$$

subject to :

$$W^T x_1 + b = 1$$

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

Using Lagrange multiplier λ_1 and λ_2 , we can write the following:

$$L = \arg \min_{W \in \mathbb{R}^P} \left\{ \|W\|_2^2 + \lambda_1 (W^T x_1 + b - 1) + \lambda_2 (W^T x_2 + b + 1) \right\}$$

Taking the derivative w.r.t. W and b and make them equal to 0.

$$\frac{\partial L}{\partial W} = 0 \Rightarrow 2W + \lambda_1 x_1 + \lambda_2 x_2 = 0$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \lambda_1 + \lambda_2 = 0$$

$$\lambda_1 = -\lambda_2$$

$$W = -\frac{1}{2} (\lambda_1 x_1 + \lambda_2 x_2) \Rightarrow \text{substituting } \lambda_1 = -\lambda_2$$

$$W = -\frac{1}{2} (-\lambda_2 x_1 + \lambda_2 x_2)$$

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

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

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

$$2b = -w^T(x_1 + x_2)$$

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

$$2 = 1+1 \Rightarrow 2 = (w^T x_1 + b) - (w^T x_2 + b)$$

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

$2 = w^T(x_1 - x_2)$ We now substitute
for w .

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

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

$$\lambda_2 = 4 (x_1^T x_1 + x_1^T x_2 - x_2^T x_1 + x_2^T x_2)^{-1}$$

$$\lambda_1 = -\lambda_2$$

x
 $n \times d$

② Linear Regression with Regularization:

The loss function $L(w) = \sum_{i=1}^n (y_i - w^T x_i)^2$

sum of squared errors from lecture notes

$$= (y - Xw)^T (y - Xw)$$

Adding the penalty:

$$\begin{aligned} L(w) &= (y - Xw)^T (y - Xw) + \lambda w^T w \\ &= (y^T - w^T x^T)(y - Xw) + \lambda w^T w \\ &= y^T y - y^T Xw - w^T x^T y + w^T x^T Xw + \lambda w^T w \\ &= y^T y - 2y^T Xw + w^T X^T Xw + \lambda w^T w \end{aligned}$$

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

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

$$x^T y = x^T x w + \lambda w$$

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

$$w = (x^T x + \lambda I)^{-1} x^T y \quad \begin{matrix} \text{parameter of linear regression} \\ \text{with penalty.} \end{matrix}$$

It penalizes w for taking large values. It makes w small to prevent the coefficients from overfitting.

(4) Conceptual.

$$P(w|x)$$

$$\frac{P(x|w) P(w)}{P(x)}$$

- Easier to model because it doesn't require modeling the joint distribution $P(w, x)$.
- Estimate the posterior directly.
- Can't detect outlier in the data.
- $P(x|w)$ become hard to model if the dimension x is large.
- uses the available data to estimate the prior $P(w)$, likelihood $P(x|w)$, and evidence $P(x)$.
- Known the evidence term $P(x)$ is useful because it normalizes the term and changes the posterior into probability $[0, 1]$.
The likelihood term can be bigger than 1.