

<> Code

Issues

Pull requests

Actions

Projects

Wiki

Security

 master ▾

...

[ece523\\_doc21](#) / [ece523](#) / [hw](#) / [hw2](#) / hw2ml.pdf



SafwanElmadani update

 History

 1 contributor

4.78 MB

...

# ① Linear classifier with a Margin :

a data set has 2 data points

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

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

Setup the minimization problem w/

constraints on  $w^T x_1 + b$

$$w^T x_2 + b$$

To find the hyperplane, we need to solve

$$\arg \min_{w \in \mathcal{R}^p} \|w\|_2^2 = \arg \min_{w \in \mathcal{R}^p} w^T w$$

subject to :

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

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

Using Lagrange multiplier  $\lambda_1$  and  $\lambda_2$ , we can write the following:

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

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

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

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

$$\lambda_1 = -\lambda_2$$

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

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

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

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

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

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

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

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

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

$$2 = w^T (x_1 - x_2) \quad \text{We now substitute for } w.$$

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

$$2 = \frac{\lambda_2}{2} (x_1^T x_1 + x_1^T x_2 - x_2^T x_1 + x_2^T x_2) \quad x_{n \times d}$$

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

$$\lambda_1 = -\lambda_2$$

## ② Linear Regression with Regularization:

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

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

sum of squared  
errors from lecture notes

Adding the penalty:

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

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

$$0 = -2 X^T y + 2 X^T X w + 2 \lambda w \Rightarrow 0 = -X^T y + X^T X w + \lambda w$$

$$\frac{\partial L}{\partial W} = 2X^T(W - y) + 2\lambda W = 0$$

$$X^T y = X^T X W + \lambda W$$

$$X^T y = (X^T X + \lambda I) W$$

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

parameter of linear regression with penalty.

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.



