

# Solutions to Selected Problems in *Machine Learning: An Algorithmic Perspective*

Alex Kerr  
email: ajkerr0@gmail.com

## Chapter 2

### Problem 2.1

Let's say  $S$  is the event that someone at the party went to the same school,  $R$  is the event that someone at the party is vaguely recognizable. By Bayes' rule:

$$P(S|R) = \frac{P(R|S)P(S)}{P(R)} \quad (1)$$

We are given:

$$\begin{aligned} P(R|S) &= 1/2 = 0.5 \\ P(S) &= 1/10 = 0.1 \\ P(R) &= 1/5 = 0.2 \end{aligned}$$

The result is

$$P(S|R) = \frac{0.5 \times 0.1}{0.2} = 0.25 \quad (2)$$

## Chapter 3

### Problem 3.1

$$\mathbf{x} \cdot \mathbf{w}^T - \text{bias} = 0 \quad (3)$$

$$x_1 + x_2 = -1.5 \quad (4)$$

This function would not correspond to a logic gate, but a constant voltage source.

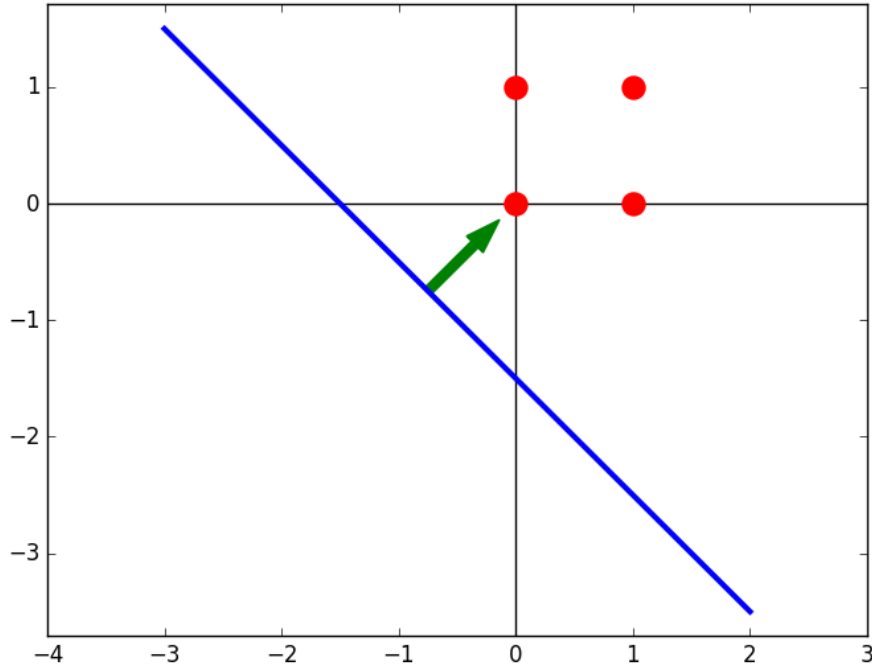
### Problem 3.7

#### Lagrange Multipliers

We are trying to minimize  $f(\mathbf{x}) = \|\mathbf{x} - \mathbf{x}'\|^2$  while constraining  $\mathbf{x}$  to the hyperplane  $y = \sum_i w_i x_i + b$ . First look at partial derivatives:

$$\nabla f = 2 \sum_i (x_i - x'_i) \hat{\mathbf{x}}_i \quad (5)$$

$$\nabla g = \sum_i w_i \hat{\mathbf{x}}_i \quad (6)$$



**Figure 3.1:** The *blue* line is the decision boundary ( $x_2 = -x_1 - 1.5$ ), the *green* arrow is the weight vector  $\mathbf{w} = (1, 1)$ , and the *red* points are the standard logic gate inputs.

In accordance to  $\nabla f - \lambda \nabla g = 0$  we have

$$x_i = x'_i + \frac{\lambda}{2} w_i \quad (7)$$

Sub this back into the constraint  $g$ :

$$\sum_i w_i x'_i + \frac{\lambda}{2} \sum_i w_i^2 + b = 0 \quad (8)$$

$$\lambda = -2 \frac{\sum_i w_i x'_i + b}{\sum_i w_i^2} \quad (9)$$

Sub Equation 7 back into  $f$  and use our value for  $\lambda$ :

$$f_{\min} = \frac{\lambda^2}{4} \sum_i w_i^2 = \frac{(\sum_i w_i x'_i + b)^2}{\sum_i w_i^2} = \frac{y(\mathbf{x}')^2}{\|\mathbf{w}\|^2} \quad (10)$$

Finally the distance:

$$\|\mathbf{x} - \mathbf{x}'\|_{\min} = \sqrt{f_{\min}} = \frac{|y(\mathbf{x}')|}{\|\mathbf{w}\|} \quad (11)$$

### Projection

The shortest vector between  $\mathbf{x}'$  and the hyperplane is parallel to the vector orthogonal to the hyperplane. This is the vector that points ‘most directly’ to  $\mathbf{x}$ . In the case of the Perceptron boundary, this vector is  $\mathbf{w}^T$ . The boundary is ‘flat’ so if we ensure we project  $\mathbf{x}' - \mathbf{x}$  to  $\mathbf{w}^T$  we are good to go. For the projection, I like to use the Gram-Schmidt process as a reference.

$$\text{proj}_{\mathbf{w}}(\mathbf{x}' - \mathbf{x}) = \frac{(\mathbf{x}' - \mathbf{x}) \cdot \mathbf{w}}{\mathbf{w} \cdot \mathbf{w}} \mathbf{w} = \frac{\mathbf{x}' \cdot \mathbf{w} - \mathbf{x} \cdot \mathbf{w}}{w^2} \mathbf{w} \quad (12)$$

Substitute in the discriminant function:

$$\mathbf{w} \cdot \mathbf{x} = -b \quad (13)$$

$$\frac{\mathbf{x}' \cdot \mathbf{w} + b}{w^2} \mathbf{w} = \frac{y(\mathbf{x}')}{w^2} \mathbf{w} = \frac{y(\mathbf{x}')}{w} \hat{\mathbf{w}} \quad (14)$$

The magnitude is

$$\frac{|y(\mathbf{x}')|}{\|\mathbf{w}\|} \quad (15)$$

## Chapter 4

### Problem 4.5

Refer to the `train_seq` method of the multi-level Perceptron class MLP on my [github repository](#).

### Problem 4.12

$$\begin{aligned} \tanh(h) &= \frac{e^h - e^{-h}}{e^h + e^{-h}} = \frac{1 - e^{-2h}}{1 + e^{-2h}} \\ &= \frac{2 - 1 - e^{-2h}}{1 + e^{-2h}} \\ &= \frac{2}{1 + e^{-2h}} - \frac{1 + e^{-2h}}{1 + e^{-2h}} \\ &= 2g(2h) - 1 \end{aligned} \quad (16)$$

## Chapter 6

### Problem 6.1

Refer to my `lda` function on [github](#).

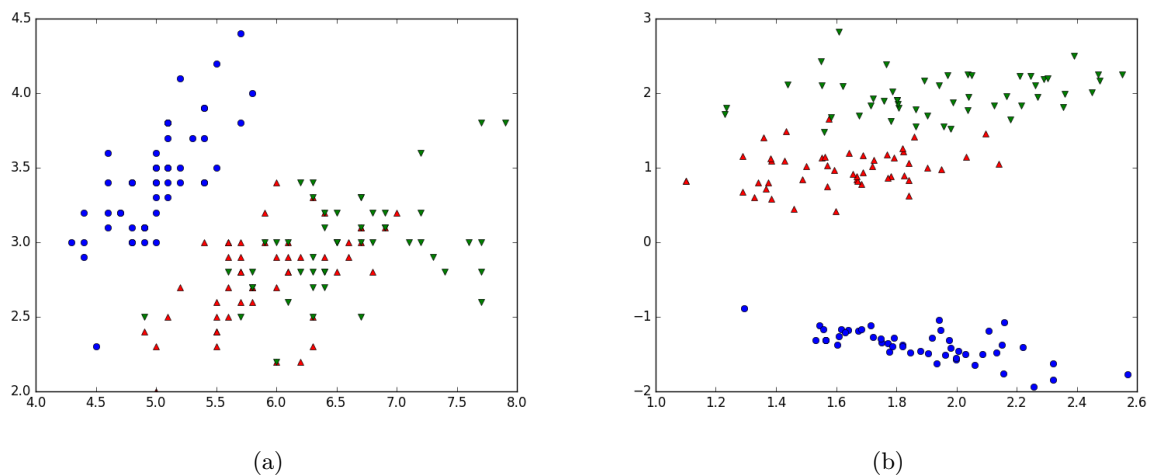
The author may have made mistakes in the text relating the scatter matrices. At the very least he wasn't clear. The covariance matrix is estimated by the total scatter matrix in the `numpy` routine:

$$\text{cov}(\mathbf{x}_a, \mathbf{x}_b) = \text{E} [(\mathbf{X} - \text{E}[\mathbf{X}])(\mathbf{X} - \text{E}[\mathbf{X}])^T] \simeq \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T \quad (17)$$

where the mean data point is understood to be the sample mean:

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad (18)$$

We can expand the total scatter matrix as:



**Figure 6.1** *Left*: The original iris data. *Right*: The reduced iris data after LDA.

$$\begin{aligned}
\frac{1}{n} \sum_{i=0}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T &= \frac{1}{n} \sum_c \sum_{j \in c} (\mathbf{x}_j - \boldsymbol{\mu})(\mathbf{x}_j - \boldsymbol{\mu})^T \\
&= \frac{1}{n} \sum_c \sum_{j=1}^{m_c} (\mathbf{x}_j - \boldsymbol{\mu}_c)(\mathbf{x}_j - \boldsymbol{\mu}_c)^T + (\boldsymbol{\mu}_c - \boldsymbol{\mu})(\boldsymbol{\mu}_c - \boldsymbol{\mu})^T \\
&\quad + (\mathbf{x}_j - \boldsymbol{\mu}_c)(\boldsymbol{\mu}_c - \boldsymbol{\mu})^T + (\boldsymbol{\mu}_c - \boldsymbol{\mu})(\mathbf{x}_j - \boldsymbol{\mu}_c)^T
\end{aligned} \tag{19}$$

Let's look at Equation 19 piece by piece. The third and fourth terms goes to zero:

$$\sum_{j=1}^m (\mathbf{x}_j - \boldsymbol{\mu}_c) = \sum_{j=1}^m \mathbf{x}_j - m\boldsymbol{\mu}_c = \sum_{j=1}^m \mathbf{x}_j - \frac{m}{m} \sum_{j=1}^m \mathbf{x}_j = 0 \tag{20}$$

The within-scatter matrix as written in the text is invalid. It comes from the our first term in Equation 19:

$$\begin{aligned}
\frac{1}{n} \sum_c \sum_{j=1}^m (\mathbf{x}_j - \boldsymbol{\mu}_c)(\mathbf{x}_j - \boldsymbol{\mu}_c)^T &= \sum_c \frac{m_c}{n} \frac{1}{m_c} \sum_{j=0}^m (\mathbf{x}_j - \boldsymbol{\mu}_c)(\mathbf{x}_j - \boldsymbol{\mu}_c)^T \\
&= \sum_c p_c \text{cov}_c(\mathbf{x}_a, \mathbf{x}_b)
\end{aligned} \tag{21}$$

where the covariance is calculated as the total scatter matrix of class  $c$  ala Equation 17. This is where the  $p_c = m_c/n$  term comes from in the code, but the expression is not exactly right in the text. I think there is a similar error in book equation 6.2. It comes from the second term of Equation 19 and it should read:

$$S_B = \frac{1}{n} \sum_c \sum_{j=1}^m (\boldsymbol{\mu}_c - \boldsymbol{\mu})(\boldsymbol{\mu}_c - \boldsymbol{\mu})^T = \frac{1}{n} \sum_c m_c (\boldsymbol{\mu}_c - \boldsymbol{\mu})(\boldsymbol{\mu}_c - \boldsymbol{\mu})^T = \sum_c p_c (\boldsymbol{\mu}_c - \boldsymbol{\mu})(\boldsymbol{\mu}_c - \boldsymbol{\mu})^T \tag{22}$$