

Fisher's linear discriminant

$$L(w) = \frac{(\text{separation of projected means})^2}{\text{sum of within-class variances}} \quad (\text{want to maximize})$$

Let $m_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n$ (class means in the original space),

$m_k = w^T m_k$ (means in projection space),

and $y_n = w^T \mathbf{x}_n$ (projections / labels of each data point).

Then the separation of projected means² is:

$$(m_2 - m_1)^2 = w^T \underbrace{(m_2 - m_1)(m_2 - m_1)^T w}_{S_B}$$

and the sum of within-class variances is:

$$\sum_k \sum_{n \in C_k} (y_n - m_k)^2 = w^T \underbrace{\sum_k \sum_{n \in C_k} (\mathbf{x}_n - m_k)(\mathbf{x}_n - m_k)^T w}_{S_w}$$

So the objective is:

$$L(w) = \frac{w^T S_B w}{w^T S_w w}$$

We set the derivative to 0 to find the solution w :

$$\begin{aligned} \frac{\partial}{\partial w} L(w) &= w^T S_w w \frac{d}{dw} [w^T S_B w] - w^T S_B w \frac{d}{dw} [w^T S_w w] \\ &\stackrel{\text{quotient rule}}{=} w^T S_w w \cdot 2 S_B w - w^T S_B w \cdot 2 S_w w = 0 \\ &\Rightarrow \underbrace{(w^T S_w w)}_{\text{scalar}} S_B w = \underbrace{(w^T S_B w)}_{\text{scalar}} S_w w \end{aligned}$$

①

$$\Rightarrow S_w w \propto S_B w$$

$$\Rightarrow S_w^{-1} S_w w \propto S_w^{-1} S_B w$$

$$\Rightarrow w \propto S_w^{-1} (m_2 - m_1) \underbrace{(m_2 - m_1)^T w}_{\text{scalar}}$$

Perceptron

- * If $t_n = +1$, want $w^T \phi_n > 0 \Rightarrow (t_n) w^T \phi_n > (t_n) 0 = 0$.
If $t_n = -1$, want $w^T \phi_n < 0 \Rightarrow (t_n) w^T \phi_n < (t_n) 0 = -1$
i.e. we want $t_n w^T \phi_n > 0$ for all samples.
- * If this condition is violated for some sample, we want it to be violated by as little as possible, i.e. for misclassified samples want $t_n w^T \phi_n \leq 0$ as close to 0 as possible, i.e. want to maximize $t_n w^T \phi_n$ / minimize $-t_n w^T \phi_n$.
- * $\frac{\partial}{\partial w} [-t_n w^T \phi_n] = -t_n \phi_n$

Bayes theorem : $P(A|B) = \frac{P(B|A) P(A)}{P(B)} = \frac{P(B|A) P(A)}{\sum_i P(B|A_i) P(A_i)}$

Maximum likelihood estimation : $w^* = \operatorname{argmax}_w P(\text{Data} | w)$

* Example: We want to model coin tosses. The underlying model $w = P(H|w)$.

Then $w^* = \operatorname{argmax}_w L(w)$ where $L(w) = P(\{H, T, T, H, H\} | w) = P(H|w)^{N_H} P(T|w)^{N_T}$.

$$\frac{\partial}{\partial w} [\log L(w)] = 0 = \frac{\partial}{\partial w} [N_H \log P(H|w) + N_T \log P(T|w)] = \frac{\partial}{\partial w} [N_H \log w + N_T \log(1-w)]$$

$$= \frac{N_H}{w} - \frac{N_T}{1-w} \Rightarrow N_H w - N_T w = 0 \Rightarrow w = \frac{N_H}{N_H + N_T} \quad (\text{as expected!})$$

(2)

Logistic regression

* $P(y_i=1|x_i) = \frac{1}{1+e^{-w^T x_i}} = \boxed{\mathcal{G}(w^T x_i)}$ where \mathcal{G} is the logistic sigmoid.

* Decision boundary: $P(y=1|x) \stackrel{?}{>} P(y=0|x) \Rightarrow \frac{P(y=1|x)}{P(y=0|x)} \stackrel{?}{>} 1$

$$\Rightarrow \log \left[\frac{P(y=1|x)}{P(y=0|x)} \right] \stackrel{?}{>} 0 \Rightarrow \log \left[\frac{P(y=1|x)}{1-P(y=1|x)} \right] \stackrel{?}{>} 0$$

$$\Rightarrow \log \left[\frac{\frac{1}{1+e^{-w^T x}}}{1 - \frac{1}{1+e^{-w^T x}}} \right] \stackrel{?}{>} 0 \Rightarrow \log \left[\frac{\frac{1}{1+e^{-w^T x}}}{\frac{(1+e^{-w^T x})-1}{1+e^{-w^T x}}} \right] \stackrel{?}{>} 0 \Rightarrow$$

$$\Rightarrow \log \frac{1}{e^{-w^T x}} \stackrel{?}{>} 0 \Rightarrow -\log e^{-w^T x} \stackrel{?}{>} 0 \Rightarrow w^T x \stackrel{?}{>} 0$$

(linear classifier)

* Solution for w : $w^* = \underset{w}{\operatorname{argmax}} L(w)$ where

$$L(w) = P(\text{Data} | w) = \prod_{i=1}^N P(y_i | w, x_i) = \prod_{i=1}^N \mathcal{G}(w^T x_i)^{y_i} (1-\mathcal{G}(w^T x_i))^{(1-y_i)}$$

$$\log L(w) = \sum_{i=1}^N y_i \log \mathcal{G}(w^T x_i) + (1-y_i) \log (1-\mathcal{G}(w^T x_i)) \quad (y_i=1 \text{ if pos, } y_i=0 \text{ if neg})$$

$$\begin{aligned} \frac{d L(w)}{d w} &= 0 = \sum_{i=1}^N y_i \frac{1}{\mathcal{G}(w^T x_i)} \frac{d \mathcal{G}(w^T x_i)}{d w} + (1-y_i) \frac{1}{1-\mathcal{G}(w^T x_i)} \frac{-d \mathcal{G}(w^T x_i)}{d w} \\ &= \nabla_w L(w) \\ &= \sum_{i=1}^N \frac{y_i - y_i \mathcal{G}(w^T x_i) - \mathcal{G}(w^T x_i) + y_i \mathcal{G}(w^T x_i)}{\mathcal{G}(w^T x_i)(1-\mathcal{G}(w^T x_i))} \frac{d \mathcal{G}(w^T x_i)}{d w} \\ &= \sum_{i=1}^N \frac{(y_i - \mathcal{G}(w^T x_i))}{\mathcal{G}(w^T x_i)(1-\mathcal{G}(w^T x_i))} \underbrace{\mathcal{G}(w^T x_i)(1-\mathcal{G}(w^T x_i))}_{\text{derivative of sigmoid}} \underbrace{x_i}_{\text{chain rule}} \\ &= \sum_{i=1}^N (y_i - \mathcal{G}(w^T x_i)) x_i \end{aligned}$$

prediction error = 0 if: $y_i=1, P(y_i=1|x_i)=1$ or
 $y_i=0, P(y_i=1|x_i)=0$

* Gradient ascent update: $w^{(t+1)} = w^{(t)} + \gamma \frac{d L(w)}{d w}$

$$w^{(t+1)} = w^{(t)} + \gamma \sum_{i=1}^N (y_i - \mathcal{G}(w^{(t)T} x_i)) x_i \quad (3)$$