

Quiz 2 : next Wednesday (Oct 25th)
Hw 3 will be out this week

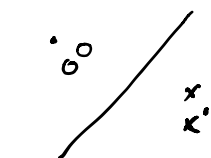
Outline ML estimation / MAP estimation / Laplaceian Smooth
Logistic Regression

Recap: LDA / QDA / Naïve Bayes

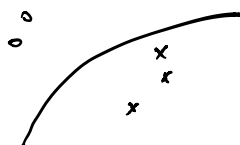
$$\hat{f}(x) = \arg \max_{k=1, \dots, K} P(Y=k | X=x)$$

$$= \arg \max_{k=1, \dots, K} P(Y=k) \cdot P(X=x | Y=k)$$

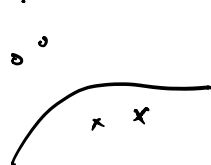
↑
class conditional dist



LDA $N(\mu_k, \Sigma)$



QDA $N(\mu_k, \Sigma_k)$



NB $\prod_{j=1}^d g^{(j)}(x)$

(ex. Gaussian NB $g^{(j)} \sim N(\mu_j, \sigma_j^2)$)

NB discrete model (Laplaceian Smoothing)

$$P(\text{"Michigan"} | Y = \text{"sports"}) \approx \frac{\sum_{i \in \text{sports}} (c(d_i, w_1) + 1)}{\sum_j \sum_{i \in \text{sports}} (c(d_i, w_j) + 1)}$$

Why add 1?

Example coin tossing Estimate $P(\text{heads}) = \theta$ from n independent coin tosses

$$P(x_i | \theta) = \begin{cases} \theta & x_i = 1 \\ 1 - \theta & x_i = 0 \end{cases} = \theta^{x_i} \cdot (1 - \theta)^{1 - x_i}$$

Likelihood function: $p(\underline{x} | \theta) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1 - x_i}$

$$= \theta^{\sum_{i=1}^n x_i} \cdot (1 - \theta)^{\sum_{i=1}^n 1 - x_i}$$

N_1 : # heads
 N_0 : # tails

$$= \theta^{N_1} \cdot (1 - \theta)^{N_0}$$

Maximum Likelihood (ML) estimate of θ

$$\arg \max_{\theta} p(\underline{x} | \theta)$$

$$= \arg \max_{\theta} \log p(\underline{x} | \theta) = \arg \max_{\theta} N_1 \log \theta + N_0 \log(1 - \theta)$$

Concave $\Rightarrow \nabla \log p(\underline{x} | \theta) = 0$

$$\frac{N_1}{\theta} + \frac{N_0}{1 - \theta} = 0 \Rightarrow \theta = \frac{N_1}{N_0 + N_1} = \frac{\# \text{heads}}{n}$$

\rightarrow Now suppose θ is random

$\theta \sim$ Beta distribution (α_1, α_0)

$$p(\theta) = \theta^{\alpha_1 - 1} \cdot (1 - \theta)^{\alpha_0 - 1}$$

$$p(\theta | x_1, \dots, x_n) \propto p(x_1, \dots, x_n | \theta) \cdot p(\theta)$$

$$= \theta^{N_1} \cdot (1 - \theta)^{N_0} \cdot \theta^{\alpha_1 - 1} \cdot (1 - \theta)^{\alpha_0 - 1}$$

$$= \theta^{N_1 + \alpha_1 - 1} \cdot (1 - \theta)^{N_0 + \alpha_0 - 1}$$

$$= \theta \cdot (1-\theta)$$

Maximum a posteriori estimation (MAP)

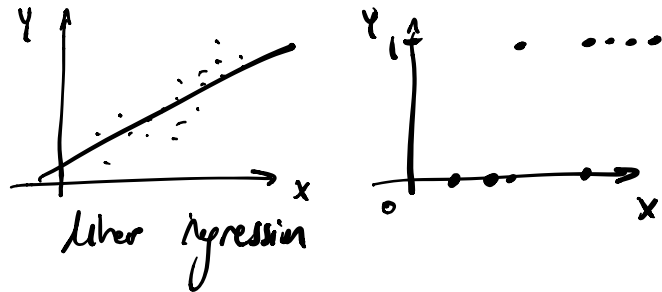
$$\arg \max_{\theta} p(\theta | x_1, \dots, x_n) = \frac{N_1 + \alpha_1 - 1}{N_0 + N_1 + \alpha_0 + \alpha_1 - 2}$$

add $\alpha_1 - 1$ to heads
 $\alpha_0 - 1$ to tails

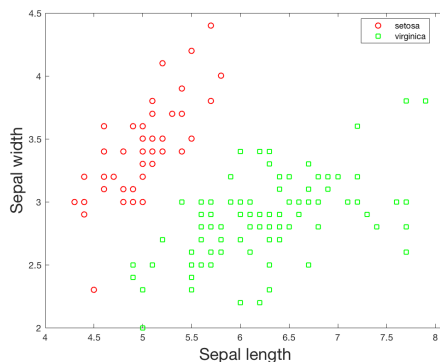
$\alpha_1 = \alpha_0 = 2 \Rightarrow$ add 1
Laplace Smoothing

Multinomial case, Dirichlet prior implies adding 1 to counts.

Logistic Regression



In LDA / ODA / NB we made assumptions on $(x_i^{(0)}, \dots, x_i^{(L)})$



prob. model of
 $P(Y=1 | X=x)$

Consider a binary classification problem w. labels $y \in \{0, 1\}$

The Bayes Classifier

$$f^*(x) = \begin{cases} 1 & \text{if } \eta(x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

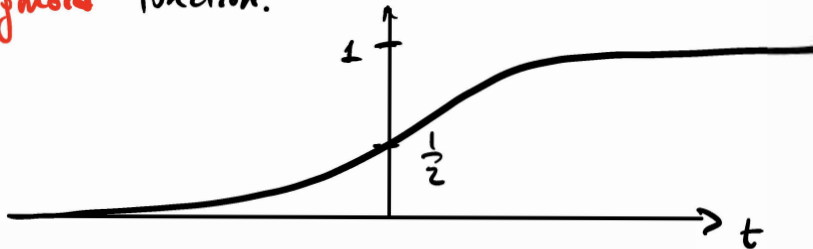
$$\eta(x) = P(Y=1 | X=x) \quad \text{posterior dist.}$$

In LR we assume $\eta(x) = \frac{1}{1 + e^{-(w^T x + b)}}$

$w \in \mathbb{R}^d, b \in \mathbb{R}$ $\theta = \begin{bmatrix} w \\ b \end{bmatrix}$. Then we find ML estimates of w and b : \hat{w}, \hat{b} , plug in to the Bayes Classifier rule

$$\hat{\eta}(x) = \frac{1}{1 + e^{-(\hat{w}^T x + \hat{b})}}$$

The function $\frac{1}{1 + e^{-t}}$ is called a **logistic** or **Sigmoid** function.



$$P(Y=1 | x) = \frac{1}{1 + e^{-(w^T x + b)}}$$

$$\text{odds} = \frac{P(Y=1 | x)}{1 - P(Y=1 | x)} = e^{+(w^T x + b)}$$

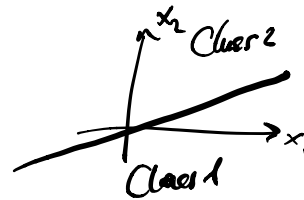
$$\Rightarrow -\log \frac{P(Y=1 | x)}{1 - P(Y=1 | x)} = w^T x + b = w_1 x_1 + \dots + w_d x_d + b$$

$$\text{LR classifier } \hat{f}(x) = 1_{\{\hat{\eta}(x) \geq \frac{1}{2}\}}$$

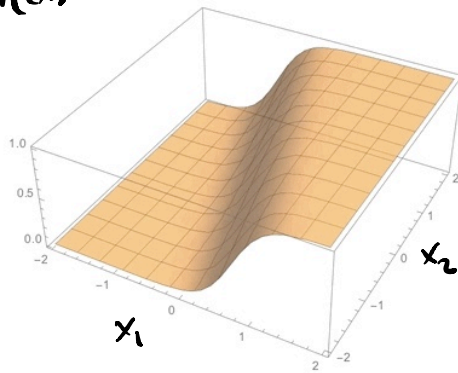
$$\hat{f}(x) = 1 \Leftrightarrow \frac{1}{1 + e^{-\hat{w}^T x - \hat{b}}} \geq \frac{1}{2}$$

$$\Leftrightarrow \hat{w}^T x + \hat{b} \geq 0$$

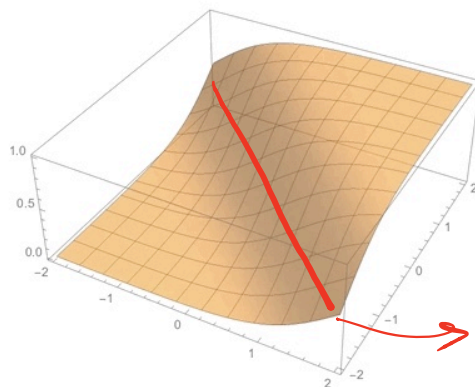
$\hat{f}(x)$ is a linear classifier



$\eta(x)$



$$\begin{pmatrix} w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 4.7 \\ 1.70587 \end{pmatrix}$$



$$\begin{pmatrix} w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 2. \\ 2. \end{pmatrix}$$

$$\rightarrow \eta(x) = \frac{1}{2}$$

ML estimation

Let $(x_1, y_1), \dots, (x_n, y_n)$. LR does not model the distribution of x . We will treat x as fixed and maximize Conditional likelihood

$p(y|x;\theta)$ denotes conditional pmf of y given x .

Conditional likelihood of θ is

$$L(\theta) = \prod_{i=1}^n p(y_i | x_i; \theta)$$

(assumed labels are cond. ind. given x_i)

$$L(\theta) = \prod_{i=1}^n p(x_i)^{y_i} \cdot (1-p(x_i))^{1-y_i}$$

$$\text{where } p(x_i) = P(Y=1 | X=x_i)$$