

Lecture 17:

Step 1: Devise a parametric model:

$$\Pr\{0|x\} = \frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}} \quad \leftarrow$$

$$\Pr\{1|x\} = 1 - \Pr\{0|x\}$$

Step 2:

Define an objective for the "goodness of fit" → minimize the objective → find \hat{B}_0, \hat{B}_1

$$\Pr\{y_i | x_i; B_0, B_1\} \quad \begin{cases} y_i = 0 \rightarrow \frac{e^{B_0 + B_1 x_i}}{1 + e^{B_0 + B_1 x_i}} \\ y_i = 1 \rightarrow \frac{1}{1 + e^{B_0 + B_1 x_i}} \end{cases}$$

$$\Pr\{x_i, y_i | B_0, B_1\} = \Pr\{x_i\} \cdot \underbrace{\Pr\{y_i | x_i, B_0, B_1\}}$$

$$\Pr\{A, B\} = \Pr\{A\} \cdot \Pr\{B | A\}$$

$$\Pr\{x_i, y_i | B_0, B_1\} = \Pr\{x_i\} \cdot \Pr\{y_i | x_i, B_0, B_1\}$$

as we learned in module 1, we can be used to estimate parameters of the distribution that generates the data

mle:

$$\begin{aligned} \text{lik}(\beta_0, \beta_1) &= \prod_{i=1}^n \Pr\{x_i, y_i \mid \beta_0, \beta_1\} \\ &= \prod_{i=1}^n \Pr\{x_i\} \Pr\{y_i \mid x_i; \beta_0, \beta_1\} \\ &= \prod_{i=1}^n \Pr\{x_i\} \cdot \underbrace{\prod_{i=1}^n \Pr\{y_i \mid x_i; \beta_0, \beta_1\}}_{\text{objective to be optimized}} \end{aligned}$$

$$\begin{aligned} \arg \max_{\beta_0, \beta_1} & \text{lik}(\beta_0, \beta_1) \\ & \text{the same value for any choice of } \beta_0, \beta_1 \\ = \arg \max_{\beta_0, \beta_1} & \underbrace{\prod_{i=1}^n \Pr\{x_i\}}_C \cdot \underbrace{\prod_{i=1}^n \Pr\{y_i \mid x_i; \beta_0, \beta_1\}}_{} \\ = \cancel{\times} \arg \max_{\beta_0, \beta_1} & \underbrace{\prod_{i=1}^n \Pr\{y_i \mid x_i; \beta_0, \beta_1\}}_{} \\ = \arg \max_{\beta_0, \beta_1} & \sum_{i=1}^n \log \Pr\{y_i \mid x_i; \beta_0, \beta_1\} \\ = \arg \max_{\beta_0, \beta_1} & \sum_{i: y_i=0} \log \left(\frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) + \sum_{i: y_i=1} \log \left(\frac{1}{1 + e^{\beta_0 + \beta_1 x_i}} \right) \end{aligned}$$

$$= \underset{B_0, B_1}{\operatorname{argmax}} \left\{ \sum_{i:y_i=0} (B_0 + B_1 x_i) - \sum_{i=1}^n \log (1 + e^{B_0 + B_1 x_i}) \right\}$$

$$\underline{g(B_0, B_1)}$$

↳ convex in terms of
 B_0, B_1

$$\left\{ \begin{array}{l} \frac{\partial g}{\partial B_0} = 0 \\ \frac{\partial g}{\partial B_1} = 0 \end{array} \right.$$

- two equations in terms
of the two unknowns B_0, B_1 .

$$\Rightarrow$$

- no closed-form solution.
⇒ we need to solve it
numerically / iteratively

- the solution is unique.

$\Rightarrow \hat{B}_0, \hat{B}_1 \leftarrow$ the best fit
to the data.

given a new
data point x

$$\hat{P}_{\text{o}}\{0|x\} = \frac{e^{\hat{B}_0 + \hat{B}_1 x}}{1 + e^{\hat{B}_0 + \hat{B}_1 x}}$$

$$\hat{P}_{\text{i}}\{1|x\} = \frac{1}{1 + e^{\hat{B}_0 + \hat{B}_1 x}}$$

$$\rightarrow \hat{y}_{\text{logistic-reg}}(x) = \underset{y \in \{0,1\}}{\operatorname{argmax}} \hat{P}_{\text{y|x}}$$

- logistic-regression for $\underline{x} \in \mathbb{R}^P \rightarrow \underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_P \end{pmatrix} \in \mathbb{R}^P$

- binary labels: $y \in \{0, 1\}$

$$\Pr\{y=0|x\} = \frac{e^{B_0 + B_1 x_1 + \dots + B_p x_p}}{1 + e^{B_0 + B_1 x_1 + \dots + B_p x_p}}$$

$$\Pr\{y=1|x\} = 1 - \Pr\{y=0|x\}$$

- general label set:

$y \in \{0, 1, 2\}$

$$\Pr\{y=i|x\} = \frac{e^{\hat{B}_0^i + \hat{B}_1^i x_1 + \dots + \hat{B}_p^i x_p}}{\sum_{j=0,1,2} e^{\hat{B}_0^j + \hat{B}_1^j x_1 + \dots + \hat{B}_p^j x_p}}$$

logistic regression is typically used for settings with binary labels ($y \in \{0, 1\}$), and it does not perform well when we have multiple classes.

- Linear / Quadratic Discriminant Analysis:

- Linear Discriminant Analysis (LDA):

$$\Pr\{y \mid x\}$$

logistic regression directly models this conditional probability.

LDA/QDA try to estimate the conditional probabilities through a different route \rightarrow Bayes rule.

$$\Pr\{y=k \mid x=x\} = \frac{\text{Bayes rule}}{\Pr\{x=x \mid y=k\} \cdot \Pr\{y=k\}}$$

LDA/QDA aim at estimating these probabilities.

$$\Pr\{x=x\}$$

$$\Pr\{x=x\} = \sum_{k=1}^K \underbrace{\Pr\{y=k\}}_{\Pr\{y=k\mid x=x\}} \underbrace{\Pr\{x=x \mid y=k\}}_{\Pr\{x=x\mid y=k\}}$$

Assume K classes:

$$\Pr\{y=k \mid x=x\} =$$

$$\frac{f_k(x) = \text{gaussian}(\mu_k) e^{-\frac{(x-\mu_k)^2}{2\sigma^2}}}{\Pr\{x=x\}}$$

$$\Pr\{x=x \mid y=k\}$$

$$\Pr\{y=k\}$$

$$\Pr\{x=x\}$$

$$\rightarrow \Pr\{y=k\}$$

$$\text{Data} = \left\{ (x_i, y_i) \right\}_{i=1}^n$$

$$\text{example } (x_i \in \mathbb{R}, y_i \in \{-, +\})$$

$$\Pr\{x=x \mid y=+\}$$



$$\Pr\{y=+\} \stackrel{\text{estimate}}{=} \frac{\sum_{i: y_i=+} 1}{n}$$

$$\Pr \{ y = k \} \approx \frac{\# \{ i : y_i = k \}}{n}$$

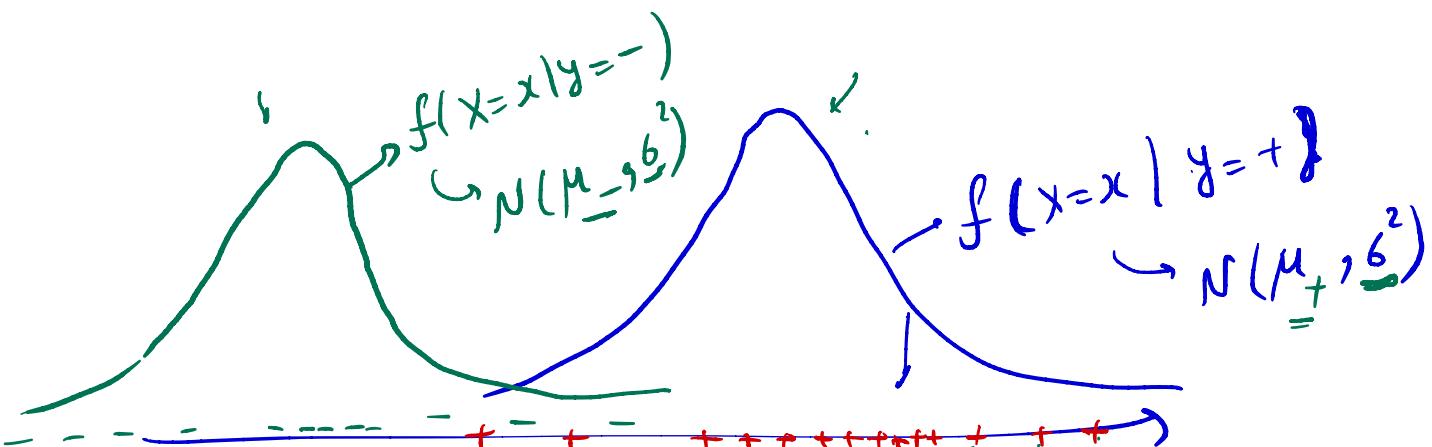
↓

Unbiased estimate of
 $\Pr \{ y = k \}$

$$f \{ x = x \mid y = k \} \stackrel{\text{LDA}}{=} f_k(x)$$

$$= \frac{1}{\sqrt{2\pi}^6} e^{-\frac{(x - \mu_k)^2}{2\sigma^2}}$$

LDA "assumes" / "models" that data, given a class k , is generated according to the gaussian distribution.



In LDA, we assume that the variance of the data x is the same over all the classes :

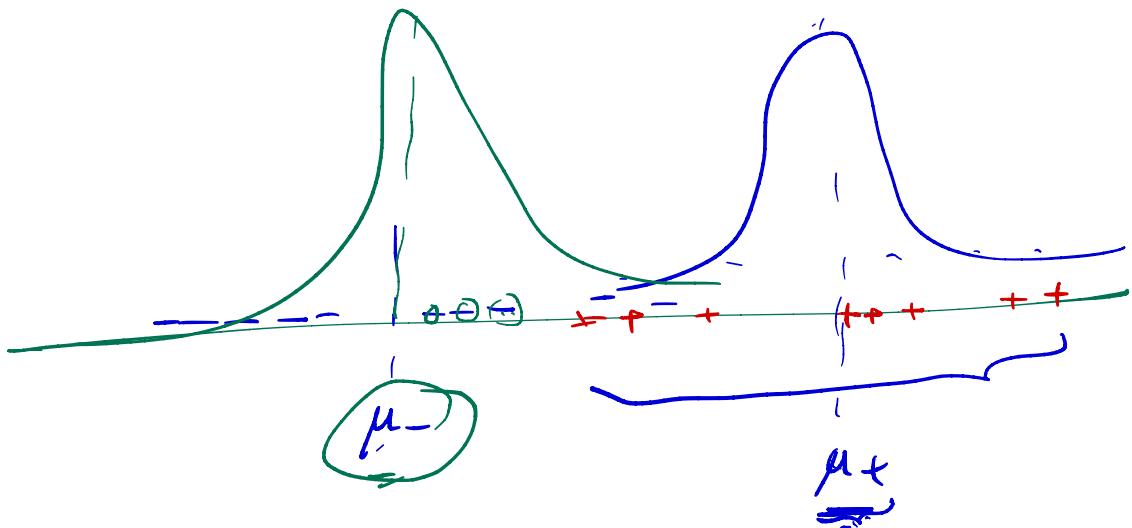
For a gaussian, we need to estimate two parameters : μ_k, σ .

$$\hat{\mu}_k \xrightarrow{\text{mean of the density}} \approx \frac{\sum_{i: y_i=k} x_i}{n_k} \rightarrow \text{Sample mean}$$

$$n_k = \#\{ i : y_i=k \}$$

Let's now estimate σ :

$$\hat{\sigma}^2 = \frac{1}{n-k} \sum_{k=1}^K \underbrace{\sum_{i: y_i=k} (x_i - \hat{\mu}_k)^2}_{\text{variance across class } k}$$



* LDA uses the same σ^2 for each class, \Rightarrow LDA requires to learn less parameters than the case where we assign different σ^2 's to the classes.

\Rightarrow QDA \rightarrow uses different σ^2 's for the classes.

LDA :