# Logistic regression

binary classification, $\underline{X}$ — predictors ($\in \mathbb{R}^n$), $Y \in \{0, 1\}$ — class

model for $p(X) = Pr(Y=1 | X)$

If $X$ is a discrete random variable, then it's simple: suppose $X \in \{a_1, ..., a_m\}$

$$p(a_m) = Pr(Y=1 | X=a_m) = \frac{Pr(Y=1 \wedge X=a_m)}{Pr(X=a_m)}$$



otherwise,

§ heuristically: $p(\ast) \, Pr(X=x) = Pr(Y=1 \wedge X=x)$

$$\oint_A p(x)\,dx \cdot P(X \in A) \overset{?}{=} \int_A p(x)\, P(X=x)\, dx \overset{?}{=} \int_A Pr(Y=1 \wedge X=x)\, dx \overset{?}{=} \underline{Pr(Y=1, X \in A)}$$

$$\int_A p(x)\,dx = \frac{Pr(Y=1 | X \in A)}{Pr(X \in A)}$$

$$p(X) = \varphi(\beta_0 + \beta_1 X) = \frac{1}{1 + e^{-\beta_0 - \beta_1 X}}$$
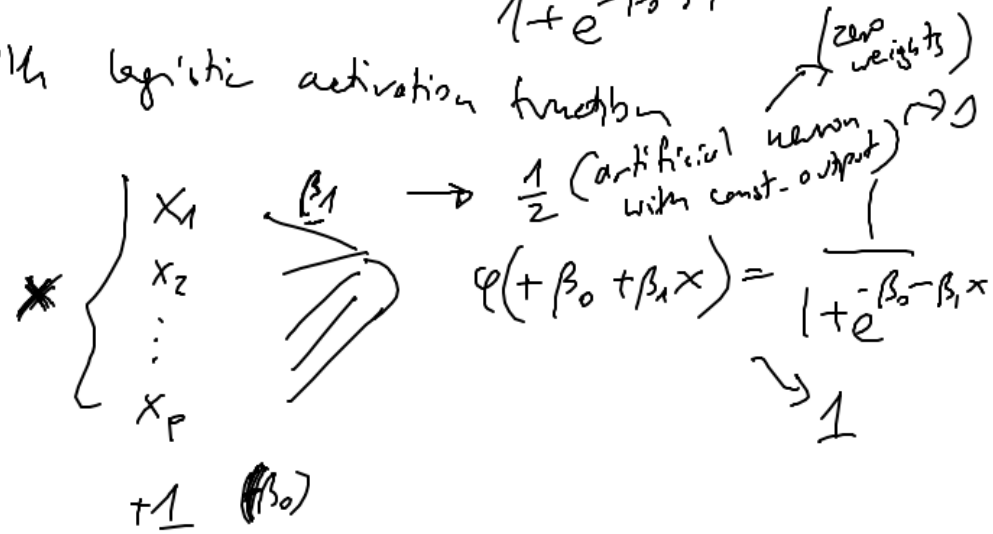
$$\varphi(x) = \frac{1}{1 + e^{-x}} \in (0, 1)$$

$$\implies \beta_0 + \beta_1 X = \log \frac{p(x)}{1 - p(x)}$$

$\beta_0, \beta_1$ are chosen to maximise the likelyhood function:

$$\mathcal{L}(\beta_0, \beta_1) = \prod_{i: y_i=1} p(x_i) \cdot \prod_{i: y_i=0} \left(1 - p(x_i)\right)$$

$$p(x) = \Pr(Y = 1 \mid X = x)$$
$$\parallel$$
$$\frac{1}{1 + e^{-\beta_0 - \beta_1 x}}$$

This corresponds to a single neuron with logistic activation function (zero weights)

(and $-\ell$ as the loss function)

$$\ast \left\{ \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_p \end{array} \right. \quad \xrightarrow{\beta_1} \quad \frac{1}{2} \text{ (artificial neuron with const. output)}$$

$$\varphi(+\beta_0 + \beta_1 x) = \frac{1}{1 + e^{-\beta_0 - \beta_1 x}}$$

$$+1 \quad (\beta_0)$$

$$\searrow 1$$

# Multinomial Logistic regression

$$Y \in \{0, 1, \ldots, k-1\} - \text{classes}$$

model:
$$\log \frac{Pr(Y=j \mid X)}{Pr(Y=k-1 \mid X)} = \beta_0^{(j)} + \beta_1^{(j)} X \qquad , j = 0, 1, \ldots, k-2$$

$$\Rightarrow \quad Pr(Y=j \mid X) = e^{\beta_0^{(j)} + \beta_1^{(j)} \cdot X} \cdot Pr(Y=k-1 \mid X) \qquad \Big| \sum_{j=0}^{k-2}$$

$$1 - Pr(Y=k-1 \mid X) = \left( \sum_{j=0}^{k-2} e^{\beta_0^{(j)} + \beta_1^{(j)} X} \right) Pr(Y=k-1 \mid X)$$

$$\begin{cases} Pr(Y=j \mid X) = \dfrac{e^{\beta_0^{(j)} + \beta_1^{(j)} \cdot X}}{1 + \sum_{j=0}^{k-1} e^{\beta_0^{(j)} + \beta_1^{(j)} x}} \qquad , j = 0, \ldots, k-2 \\[5mm] Pr(Y=k-1 \mid X) = \dfrac{1}{1 + \sum_{j=0}^{k-1} e^{\beta_0^{(j)} + \beta_1^{(j)} x}} \end{cases}$$

$\rightarrow$ neural net with $k-1$ neurons
with logistic function and 1 additional
neuron with 0 weights and also logistic act-f.

$$x : \nearrow \quad \frac{1}{1 + e^{-\beta_0^{(j)} - \beta_1^{(j)} x}} \quad , j = 0, \ldots, k-2$$

$$1 : \geq \frac{1}{2} \text{ for } k-1$$

$$x \left\{ \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right.$$

$$1 \cdot$$

$$\left[ \frac{1}{1 + e^{-\beta_0^{(j)} - \beta_1^{(j)} x}} \right] \quad \text{for} \quad j = 0, \ldots, k-2 \quad \text{class}$$

$$\left[ \varphi(0) = \frac{1}{2} \right] \quad \text{for} \quad k-1 \quad \text{class}$$

$$\overset{"}{=} \frac{1}{1 + 1}$$

$\rightarrow$ classification: the
the class with smallest
value wins

$\rightarrow e^{\beta_0^{(j)} + \beta_1^{(j)} x}$, $j = 0, \ldots, k-2$

or $1$ for $k-1$

has to be the largest

# Bayes classifier

assign each observation to the most likely class, given its predictor values

$$\underset{k}{argmax} \; Pr\left(Y=k \mid X=x_0\right)$$

Thm. The average test error rate, $E\left(\frac{1}{N} \mathbb{1}\left(y_0 \neq \hat{y}_0\right)\right)$, is the smallest for the Bayes classifier.
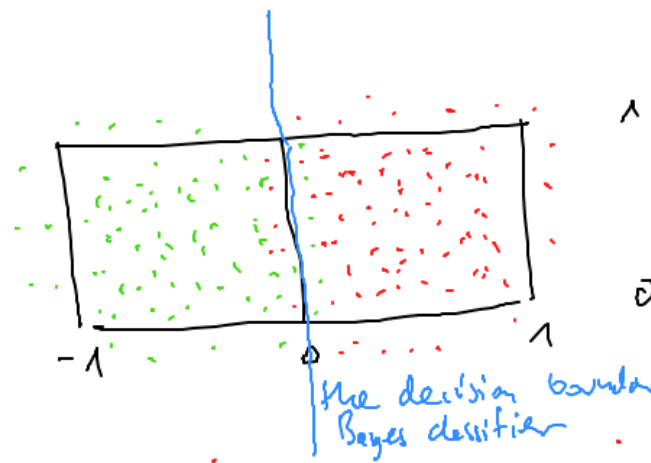
$\uparrow$ actual value (class) in the obs. $(x_0, y_0)$

$\uparrow$ predicted class

The problem is, in practice we don't know $Pr\left(Y=k \mid X=x_0\right)$

Example: Take $(\underline{x}, y)$ in the following way: $y \in \{0,1\}$

If $y=0$, then $\underline{x}=u_0+n$, with $u_0 \sim Uniform\left([-1,0] \times [0,1]\right)$, $n \sim N\left(0, \sigma^2 \underline{Id}\right)$
$u_0 \perp n$ $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

If $y=1$, then $\underline{x}=u_1+n$, with $u_1 \sim Uniform\left([0,1] \times [-1,1]\right)$, $u_1 \perp n$

$\sigma^2$ - small

the decision boundary for the
Bayes classifier

Here, one could compute the densities of the random variables $u_0 + u$ on $u_1 + u$

Bayes Theorem:

$$P(A|B) = \frac{P(B|A)\, P(A)}{P(B)}$$

# Linear Discriminant Analysis (LDA)

Model the distribution of $X$ separately in each of the response classes (i.e., given $Y$)

and then use Bayes' theorem:

$$Pr(Y = k \mid X = x) = \frac{\pi_k \cdot f_k(x)}{\sum_{\ell=1}^{K} \pi_\ell f_\ell(x)} \quad , \quad \text{where}$$

$$\pi_k = Pr(Y = k)$$

$$f_k(x) = \underbrace{Pr(X = x \mid Y = k)}_{\text{density}}$$

i.e.

$$Pr(X \in A \mid Y = k) = \int_A f_k(x) \, dx$$

$\pi_k$ are simple to estimate: just take $\dfrac{\#\{i \in \{1,\dots,N\} : \cancel{y} Y_i = k\}}{N}$

For $f_k$: assume that $X \mid Y = k$ ~~comes from some~~ are drawn from some

distribution with some parameters and then estimate the parameters

In LDA we assume that the distribution is Gaussian with the covariance matrix being the same for all classes

# LDA for $p=1$ predictors: $X \in \mathbb{R}$

- assume that $f_k(x) = \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_k)^2\right)$  for some $\mu_k, \sigma \in \mathbb{R}$, $\sigma > 0$

(i.e., we assume that all the variance $\sigma^2$ does not depend on $k$)

Then

$$\Pr(Y=k \mid X=x) = \frac{\pi_k \cdot \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{1}{2\sigma^2}(x-\mu_k)^2\right)}{\sum_{\ell=1}^{K} \pi_\ell \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_\ell)^2\right)} \qquad (*)$$

and we classify the sample to be of class $k$ for which $(*)$ is the largest.
Since the denominator is the same for all $k$, so we classify to be of class $k$
for which the foll. expression is the largest.

$$\underbrace{\log \pi_k + \log \frac{1}{\sqrt{2\pi}\sigma}}_{\text{does not depend on } k} - \frac{1}{2\sigma^2}\underbrace{(x-\mu_k)^2}_{(x^2 - 2x\mu_k + \mu_k^2)} \qquad \text{or} \qquad \underbrace{\log \pi_k + \frac{\mu_k}{\sigma^2}x - \frac{\mu_k^2}{2\sigma^2}}$$

$\uparrow$ does not depend on $k$

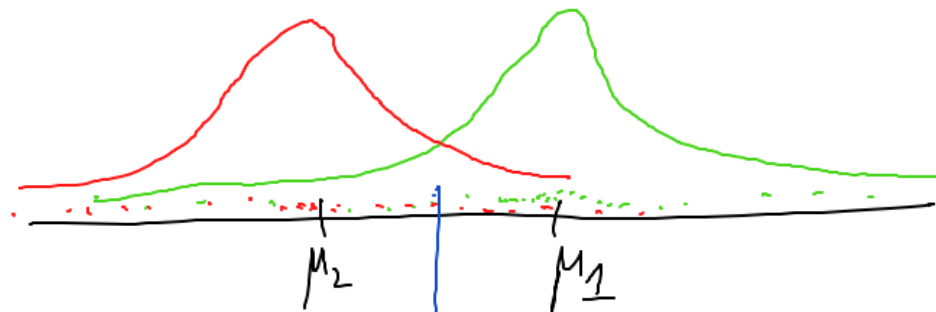Ex. Suppose that $K=2$ , $\pi_1 = \pi_2 = \frac{1}{2}$

$$\log \pi_1 + \frac{\mu_1}{\sigma^2}x - \frac{\mu_1^2}{2\sigma^2} \overset{?}{<} \log \pi_2 + \frac{\mu_2}{\sigma^2}x - \frac{\mu_2^2}{2\sigma^2}$$

$$\frac{\mu_1 - \mu_2}{\sigma^2}x < \frac{\mu_1^2 - \mu_2^2}{2\sigma^2} \qquad \left( \because (\mu_1 - \mu_2) > 0 \right.$$

Suppose $\mu_1 > \mu_2$ :

$$\frac{x}{\sigma^2} < \frac{\mu_1 + \mu_2}{2\sigma^2}$$

if $\quad x < \frac{\mu_1 + \mu_2}{2} \quad$ we predict to be in class 2



$\mu_2$     $\mu_1$

class 1 (green)
class 2 (red)

predict class 2

here predict class 1

Estimation of $\hat{\mu}_k, \hat{\sigma}^2, \hat{\pi}_k$ :

$$\hat{\pi}_k = \frac{n_k}{n} \qquad \text{, } k=1,\ldots,K$$

the number of observation in class $k$

the number of all observations

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:\, y_i = k} x_i$$

$$\hat{\sigma}^2 = \frac{1}{n-K} \sum_{k=1}^{K} \sum_{i:\, y_i = k} \left( x_i - \hat{\mu}_k \right)^2$$