# 5 Classification

Classification refers to the case when $Y$ is <u>categorical</u> or <u>qualitative</u>. (response $Y$)

**Ex**
- Will a person default on their mortgage?
- "    "    "    have heart disease as they age?
- Is an email spam?
- Will a certain ad make it more likely for someone to buy a product?

**Note** As before, still have predictors, and the response can be either binary (0/1) or multinomial.

**Methods**
- Logistic regression
- Discriminant analysis
- k-nearest neighbors
- support vector machines

Remark If $Y \in \{0, 1\}$, don't want to do usual regression

$$\underline{Y} = x\underline{\beta} + \varepsilon$$

[why?] $\quad x\hat{\underline{\beta}}$ never 0 or 1

Alternative Model $P(Y=0)$ or $P(Y=1)$

The classifier / classification rule is a rule that assigns some probabilities or 1 and others to 0.

Ex: If $P(Y=1) \geq \frac{1}{2} \Rightarrow \hat{Y} = 1$
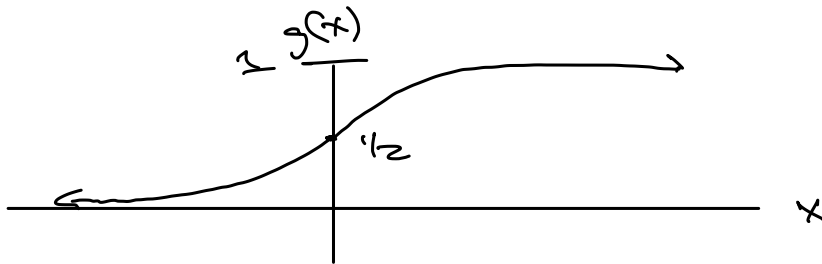
# 5.1 Logistic Regression

**Setup** $Y \in \{0, 1\}$ is binary response & $X$ is a single predictor.

**Goal** Model $P(Y=1 \mid x) := p(x)$

$p(x) = \beta_0 + \beta_1 x$ ? Bad idea b/c $\beta_0 + \beta_1 x \in \mathbb{R}$

**DEF** The <u>logistic function</u> is

$$g(x) = \frac{e^x}{1 + e^x} \in (0, 1)$$

logistic regression uses

$$p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} \qquad \left[ = P(y = 1 \mid x) \right] = \text{function of } x$$

[ Aside: other functions could be used, e.g. Probit regression
$$p(x) = \Phi(\beta_0 + \beta_1 x) \text{ where } \Phi \text{ is the cdf of } N(0,1) ]$$

[Note] The logit transform is the inverse of the logistic fun.

$$f(p) = \log\left(\frac{p}{1-p}\right)$$

So for logistic regression $\Rightarrow$

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

$\Rightarrow$ we assume the <u>log odds ratio</u> is linear in $x$.

$$\frac{p}{1-p} = \text{odds ratio}$$

$p = \frac{1}{2} \qquad \Rightarrow \text{odds} = 1 = \text{even}$

$p = 0.2 \qquad \Rightarrow \text{odds} = \frac{1}{4}$

Assumption is log odds ratio increase by $\beta_1$ for a unit increase in $x$.

[Note] The sign of $\beta_1$ tells us the effect of $x$ on $p(y=1)$

$\beta_1 > 0 \Rightarrow$ prob grows with $x$

$\beta_1 < 0 \Rightarrow$ " shrinks " "

$\boxed{\text{Estimation}}$   In setup   $y_i \in \{0,1\}$   $\Rightarrow y \sim \text{Bernoulli}(p(x))$

Suppose we have i.i.d.es data  $(x_1, y_1), \dots, (x_n, y_n)$, then

the likelihood fcn for  $\not{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

$$ f(\not{Y}) = \prod_{i=1}^{n} p(x_i)^{y_i} (1-p(x_i))^{1-y_i} \qquad \left[ = f(\beta_0, \beta_1) \right] $$

Estimate $\beta_0 + \beta_1$ by maximizing $f(\not{Y})$ (AKA maximum likelihood estimation). Closed form MLEs are not available, so maximize $f$ numerically.

$\boxed{\text{Sanity check}}$     No features: $p(x) = p = \dfrac{e^{\beta_0}}{1 + e^{\beta_0}}$

Then

$$f(x) = \prod_{i=1}^{n} p^{y_i} (1-p)^{1-y_i} = p^{n\bar{y}} (1-p)^{n-n\bar{y}}$$

take deriv wrt $p$

$$\frac{d}{dp} \log f = \frac{n\bar{y}}{p} - \frac{n-n\bar{y}}{1-p} \stackrel{set}{=} 0$$

$$\Rightarrow \hat{p} = \bar{y} = \text{proportion of 1s in data}$$

$$\hat{\beta}_0 = \log\left(\frac{\bar{y}}{1-\bar{y}}\right) = \text{log odds ratio of empirical}$$
$$P(\text{success}) \text{ to } P(\text{failure})$$

Multiple logistic regression follows:

$$P(Y=1 \mid x_1, \dots, x_p) = p(x_1, \dots, x_p) = p(\underline{x})$$

$$= \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}} = \frac{e^{\beta_0 + \beta^T \underline{x}}}{1 + e^{\beta_0 + \beta^T \underline{x}}}$$

where
$$\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$
no intercept

$$\implies \log\left(\frac{p(\underline{x})}{1 - p(\underline{x})}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

Note Can use z-statistics $\dfrac{\hat{\beta}_j}{\widehat{SE}(\hat{\beta}_j)}$ to do hypothesis testing.

Can construct CIs.

$\boxed{\text{prediction}}$ For new set of features $\underline{x} = (f_1, \ldots, x_g)^T$

our predictor for $p(\underline{x})$ is

$$\hat{p}(\underline{x}) = \hat{P}(Y=1 \mid \underline{x}) = \frac{e^{\hat{\beta}_0 + \hat{\beta}^T \underline{x}}}{1 + e^{\hat{\beta}_0 + \hat{\beta}^T \underline{x}}} \quad \in (0,1)$$

Need a way to convert $\hat{p}$ to $\hat{Y} \in \{0, 1\}$, the usual classification rule is

$$\hat{Y} = \begin{cases} 1 & \hat{p}(\underline{x}) > \frac{1}{2} \\ \\ 0 & \hat{p}(\underline{x}) < \frac{1}{2} \end{cases}$$

with randomisation at $\hat{p} = \frac{1}{2}$.

$$\hat{y} = 1 \text{ if}$$

$$\frac{1}{2} < \hat{p} \quad \Longleftrightarrow \quad \frac{1}{2} < \frac{e^{\hat{\beta}_0 + \hat{\beta}^T x}}{1 + e^{\hat{\beta}_0 + \hat{\beta}^T x}}$$

$$\Longleftrightarrow \quad \frac{1}{2} < \frac{1}{2} e^{\hat{\beta}_0 + \hat{\beta}^T x}$$

$$\Longleftrightarrow \quad 0 < \hat{\beta}_0 + \hat{\beta}^T x$$

The classification rule is :

$$\hat{y} = \begin{cases} 1 & \hat{\beta}_0 + \hat{\beta}^T \underline{x} > 0 \\ 0 & \hat{\beta}_0 + \hat{\beta}^T \underline{x} < 0 \end{cases}$$

The line $\hat{\beta}_0 + \hat{\beta}^T \underline{x} = 0$ is the <u>decision boundary</u>.