## Logistic regression

- Classification setup
- Logistic regression is the classification counterpart of linear regression.
  - we focus on binary classification, i.e. we assume
  $$X \in \mathbb{R}^d, \quad Y = \{yes, no\}$$

For example,
- the two classes can be associated with whether an individual will default on their credit card payments.
  - The attributes may be the person's annual income & monthly credit card balance
    - ($d = 2$ in this case).

## Probability prediction

- The default labels fall into one of the two categories, i.e. $Y \in \{yes, no\}$

- Idea: rather than modelling $Y$ directly, logistic regression models
  $$p(X) = Prob(Y = yes | X).$$
  i.e. the probability that $Y$ is yes (given $X$)

- Let $\tilde{X} = [1, X]^T$ with $X \in \mathbb{R}^d$, then

$$F(X, \lambda) = p_\lambda(X) = \sigma(\tilde{X}^T \lambda) = \sigma\left(\lambda_0 + \sum_{i=1}^{d} \lambda_i X_i\right)$$

where; $\sigma: \mathbb{R} \to [0,1]$ is the logistic function.

# The sigmoid function

- The logistic function

$$\sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{1+e^x}$$

- is nonnegative and bounded.
- is monotone
    - i.e. its derivative $\sigma'(x) = \sigma(x)(1-\sigma(x))$
        - is positive everywhere
- maps from the real line to the interval $[0,1]$ & it can be naturally interpreted as a probability

- $\lim\limits_{x\to\infty} \sigma(x) = 1$ and $\lim\limits_{x\to-\infty} \sigma(x) = 0$

## Odds

- The odds of the conditional probability
$$\text{Prob}(Y = yes|x) = p_\lambda(x)$$
are defined as
$$\frac{p_\lambda(x)}{1 - p_\lambda(x)} = e^{\tilde{x}^T \lambda}$$
and their logarithms,
  i.e. the log-odds of $p_\lambda(x)$, are linear in the attribute
  i.e.
$$\log \frac{P(x)}{1 - P(x)} = \tilde{x}^T \lambda$$

- Note that: $1 - p_\lambda(x) = \text{Prob}(Y = no|x)$

## Maximum likelihood estimation

- Let $D = \{(z_n, y_n) \in \mathbb{R}^d \times \{yes, no\}\}_{n=1}^N$ be a training data set.

- The interpretation of $F(x, \lambda) = \sigma(\tilde{x}^T \lambda) \in [0, 1]$
  - as the conditional probability of observing $Y = yes$ given $x$
    - allows you to estimate $\lambda$ by
    - maximizing the likelihood of $D$.

The ML estimate on $D$ is the parameter that maximises the probability of observing $D$.

## Conditional likelihood maximisation

- As the attributes are assumed to be always known,
  - we focus on maximizing the probability of observing the labels $y_1, \ldots, y_N \in D$
  → More precisely, we let
$$\hat{\lambda} = \arg \max_\lambda \mathcal{L}(D, \lambda)$$
where;
$$\mathcal{L}(D, \lambda) = \text{Prob}(y_1, \ldots, y_N | x_1, \ldots, x_N, \lambda).$$
$$= \prod_{i=1}^N \mathbb{1}[y_i = yes] F(x_i, \lambda) + \mathbb{1}[y_i = no](1 - F(x_i, \lambda))$$

## Log-likelihood minimization

- An numerically easier but equivalent problem
  - is to minimize the negative of the logarithm
  of $\mathcal{L}(D,\lambda)$,
  i.e. to let

$$\hat{\lambda} = \arg \min_{\lambda \in \mathbb{R}^{d+1}} J(D,\lambda) = -\log(\mathcal{L}(D,\lambda))$$

$$\hat{\lambda} = -\log\left(\prod_{i=1}^{N} \mathbb{1}[y_i = yes] F(x_i, \lambda) + \mathbb{1}[y_i = no](1 - F(x_i, \lambda))\right)$$

- In particular, using standard property of the logarithm
  we have
  $$J(D,\lambda) = \sum_{i=1}^{N} \log(y_i F(x_i, \lambda) + (1 - y_i)(1 - F(x_i, \lambda))).$$

  where:
  - the representation numerical of the nominal variables
    i.e. {yes, no} $\to$ {1,0}
  - allows us to rewrite the $i^{th}$ factor in
    $\mathcal{L}(D,\lambda)$ as
    $$y_i F(x_i, \lambda) + (1 - y_i)(1 - F(x_i, \lambda))$$

  - which is **differentiable** in $\lambda$.

**Example :**