# Logistic regression as a "neural network"[1]

Reminder: $\quad (x_1^{(i)}, \dots x_p^{(i)}) \qquad y^{(i)} \in \{0,1\} \qquad i = 1, \dots, N$
logreg

Assume observable is Bernoulli: $\quad P(Y=y \mid \vec{x}) = P(y=1)^y (1 - P(y=1))^{1-y}$
$\underset{y}{\phantom{xxxx}}$
$$= Ber(y \mid P(Y=1))$$

If we express probability in terms of 'log-odds':

$$\eta = \log \frac{P(Y=1)}{P(Y=0)} = \log \frac{P(Y=1)}{1-P(Y=1)} \Rightarrow P(Y=1) = \frac{1}{1+e^{-\eta}}$$

$$\text{sigm}(\eta) = h(\eta)$$



Model: log odds is linear function of descriptors:

$$X = \begin{bmatrix} 1 & x_1^{(1)} \cdots & x_p^{(1)} \\ \vdots & & \\ 1 & x_1^{(N)} \cdots & x_p^{(N)} \end{bmatrix} \qquad \eta = \vec{x}^T \cdot \vec{\beta} \qquad \vec{x} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{pmatrix} \qquad \vec{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$

Infer $\vec{\beta}$ from data by maximising likelihood

$$P(\vec{y} \mid X, \vec{\beta}) = \prod_{i=1}^{N} Ber\left(y^{(i)} \mid h\left(\vec{x}^{(i)T} \cdot \vec{\beta}\right)\right) \qquad \boxed{\begin{array}{l} Ber(\ ) \equiv \\ Bernoulli \end{array}}$$

$$L = \sum_{i=1}^{N} \left[ y^{(i)} \log\left[h\left(\vec{x}^{(i)T} \cdot \vec{\beta}\right)\right] + (1-y^{(i)}) \log\left(1 - h\left(\vec{x}^{(i)T} \cdot \vec{\beta}\right)\right) \right]$$

Optimisation: $\quad \nabla_{\vec{\beta}} L \big|_{\vec{\beta}^*_{log}} = 0 \qquad \overset{\text{Normal}}{\text{equations}} \quad \boxed{X^T \left[\vec{y} - \vec{h}\left(X \vec{\beta}^*_{log}\right)\right] = \vec{0}}$

$$h_i(X\vec{\beta}^*_{log}) = h\left(\vec{x}^{(i)T} \cdot \vec{\beta}^*_{log}\right)$$

and the problem is $\underline{convex}$ with $\underline{Hessian}$:

$$H = \nabla_{\vec{\beta}}\left[\nabla_{\vec{\beta}} L\right] = -X^T\left[\text{diag}(\vec{h}) \cdot \left[I - \text{diag}(\vec{h})\right]\right] X$$
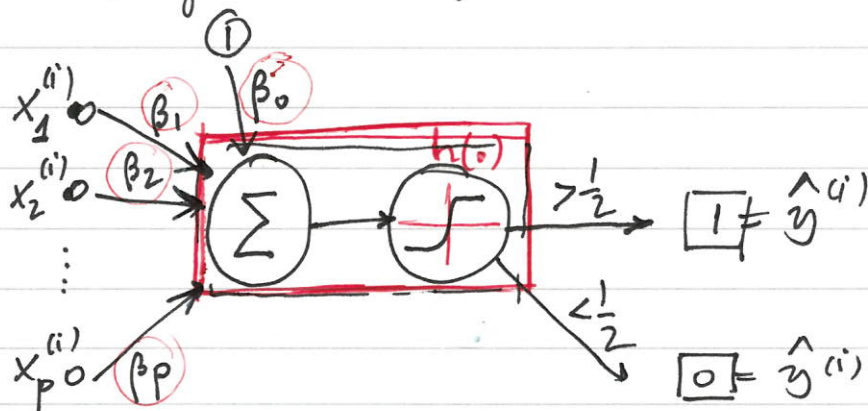
optimise using Newton or gradient methods, etc.
to obtain $\vec{\beta}^*_{log}$

Classifier:

Given $\vec{x}^{in}$, compute

$$P(y \mid \vec{x}^{in}) = \frac{1}{1 + e^{-\vec{x}^{in T} \cdot \vec{\beta}^*_{log}}}$$

$> \frac{1}{2} \quad \hat{y} = 1$

$< \frac{1}{2} \quad \hat{y} = 0$

---

Diagrammatically:



where $\vec{\beta}$ have to be optimised:

$$\max_{\vec{\beta}} L = \sum_{i=1}^{N}\left[y^{(i)} \log\left[h\left(\vec{x}^{(i)T} \cdot \vec{\beta}\right)\right] + \left(1 - y^{(i)}\right) \log\left(1 - h\left(\vec{x}^{(i)T} \cdot \vec{\beta}\right)\right)\right]$$

log likelihood is equal
to minus cross-entropy between
$\{y^{(i)}\}$ and $\{h(\vec{x}^{(i)T} \cdot \vec{\beta})\}$

Max of $L$ is minimising cross-entropy!

More generally:



$$\Pi_A + \Pi_B = 1$$

$$\begin{pmatrix} \Pi_A \\ \Pi_B \end{pmatrix} = \vec{\pi} = \begin{pmatrix} \Pi_1 \\ \vdots \\ \Pi_J \end{pmatrix}$$

$$\hat{y} = \arg\max_{J} \vec{\pi}$$

J classes

In general for J classes

J = 2 in this case

Infer $\underset{\approx}{\beta} = \begin{bmatrix} \vec{\beta}_A & \vec{\beta}_B \end{bmatrix}$  p×J

J is number of classes.

In this case $\vec{\beta}_A$ and $\vec{\beta}_B$ are related:

Original formulation
$$P(y=1) = \frac{e^{\vec{x}^T \cdot \vec{\beta}_A}}{e^{+\vec{x}^T \cdot \beta_A} + 1}$$

$$1 - P(y=1) = \frac{1}{e^{+x^T \beta_A} + 1}$$

$$\vec{\beta}'_A = \frac{\vec{\beta}_A}{2} = -\vec{\beta}'_B$$

Rewrite as

$$P(y=1) \boxed{A} = \frac{e^{\vec{x}^T \cdot \vec{\beta}'_A}}{e^{\vec{x}^T \cdot \vec{\beta}'_A} + e^{\vec{x}^T \cdot \vec{\beta}'_B}}$$

Rewritten in this form

$$P(y=0) \boxed{B} = \frac{e^{\vec{x}^T \cdot \vec{\beta}'_B}}{e^{\vec{x}^T \cdot \beta'_A} + e^{\vec{x}^T \beta'_B}}$$