

$$P(\{y^{(i)}\} | \{\vec{x}^{(i)}\}, \vec{\beta}) = \prod_{i=1}^N h_{\vec{\beta}}(\vec{x}^{(i)})^{y^{(i)}} (1 - h_{\vec{\beta}}(\vec{x}^{(i)}))^{1-y^{(i)}}$$

log-likelihood:

$$\mathcal{L} = \sum_{i=1}^N y^{(i)} \log h_{\vec{\beta}}(\vec{x}^{(i)}) + (1-y^{(i)}) \log (1 - h_{\vec{\beta}}(\vec{x}^{(i)}))$$

$\Rightarrow \nabla_{\vec{\beta}} \mathcal{L}$ to be maximised.

On model:

$$\begin{cases} P(y=1) = h_{\vec{\beta}}(\vec{x}^{(i)}) = \frac{1}{1 + e^{-\vec{x}^{(i)T} \vec{\beta}}} = h(\vec{x}^{(i)T} \vec{\beta}) \\ P(y=0) = 1 - h_{\vec{\beta}}(\vec{x}^{(i)}) = \frac{e^{-\vec{x}^{(i)T} \vec{\beta}}}{1 + e^{-\vec{x}^{(i)T} \vec{\beta}}} = e^{-\vec{x}^{(i)T} \vec{\beta}} \cdot h(\vec{x}^{(i)T} \vec{\beta}) \end{cases}$$

Note that:

$$\log(1 - h_{\vec{\beta}}(\vec{x}^{(i)})) = -\vec{x}^{(i)T} \vec{\beta} + \log(h(\vec{x}^{(i)T} \vec{\beta}))$$

$$\mathcal{L} = \sum_{i=1}^N \underbrace{\log h_{\vec{\beta}}(\vec{x}^{(i)})}_{\log h(\vec{x}^{(i)T} \vec{\beta})} - (1-y^{(i)}) (\vec{x}^{(i)T} \vec{\beta})$$

$$\begin{aligned} \nabla_{\vec{\beta}} [\log h_{\vec{\beta}}(\vec{x}^{(i)})] &= \frac{e^{-\vec{x}^{(i)T} \vec{\beta}}}{1 + e^{-\vec{x}^{(i)T} \vec{\beta}}} \vec{x}^{(i)} \\ &= (1 - h_{\vec{\beta}}(\vec{x}^{(i)})) \vec{x}^{(i)} \quad (p+1) \times 1 \end{aligned}$$

$$\begin{aligned}\nabla_{\vec{\beta}} L &= \sum_{i=1}^N \left[1 - h_{\vec{\beta}}(\vec{x}^{(i)}) \right] \vec{x}^{(i)} - (1 - y^{(i)}) \vec{x}^{(i)} \\ &= \sum_{i=1}^N \left[y^{(i)} - h_{\vec{\beta}}(\vec{x}^{(i)}) \right] \vec{x}^{(i)}\end{aligned}$$

Maximum:

$$\nabla_{\vec{\beta}} L \Big|_{\vec{\beta}^*} = 0$$

Equivalent of
normal equations
in LS.

$$\sum_{i=1}^N \left[y^{(i)} - h(\vec{x}^{(i)T} \cdot \vec{\beta}_{\log}^*) \right] \vec{x}^{(i)} = 0_v$$

(p+1) equations

L is concave \Rightarrow global maximum is $\vec{\beta}_{\log}^*$

which can be found by
standard optimisation.

$$P(\hat{y}=1) = h(\vec{x}^{inT} \cdot \vec{\beta}_{\log}^*) \in [0, 1]$$

Classifier:

$$P(\hat{y}=1 | \vec{x}^{in}) > \mathcal{T}$$

$$\vec{x}^{in} \in \{1\}$$

\mathcal{T} is the threshold for the
classifier: usually

$$\mathcal{T} = \frac{1}{2}$$

[But it can be
tuned a posteriori]

Summary of our solution for logistic regression:

$$X_{N \times (p+1)}$$

$$X\vec{\beta}^* = \vec{\hat{y}}_{N \times 1}$$

Normal equations

$$X^T [\vec{y} - \vec{h}(X\vec{\beta}^*)] = \vec{0}$$

$$\vec{h}_{N \times 1}(X\vec{\beta}^*) \Rightarrow h_i = \left(\frac{1}{1 + e^{-\vec{x}^{(i)T} \vec{\beta}^*}} \right)$$

log

Remember that for LS:

Normal equations for LS

$$X^T [\vec{y} - X\vec{\beta}_{LS}^*] = \vec{0}$$

Quality of the classifier

confusion matrix
or

contingency table

Two-classes

For a $\{0,1\}$ classifier

		True	
		1	0
Predicted	1	True positive (TP)	False positive (FP) Type I error
	0	False negative (FN) Type II error	True negative (TN)

We want most cases on the diagonal

$$\left\{ \begin{array}{l} \text{True positive rate} = \frac{TP}{TP + FN} = \text{TPR} \\ \text{Sensitivity or recall} \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{True negative rate} = \frac{TN}{TN + FP} = \text{TNR} \\ \text{Specificity} \end{array} \right.$$

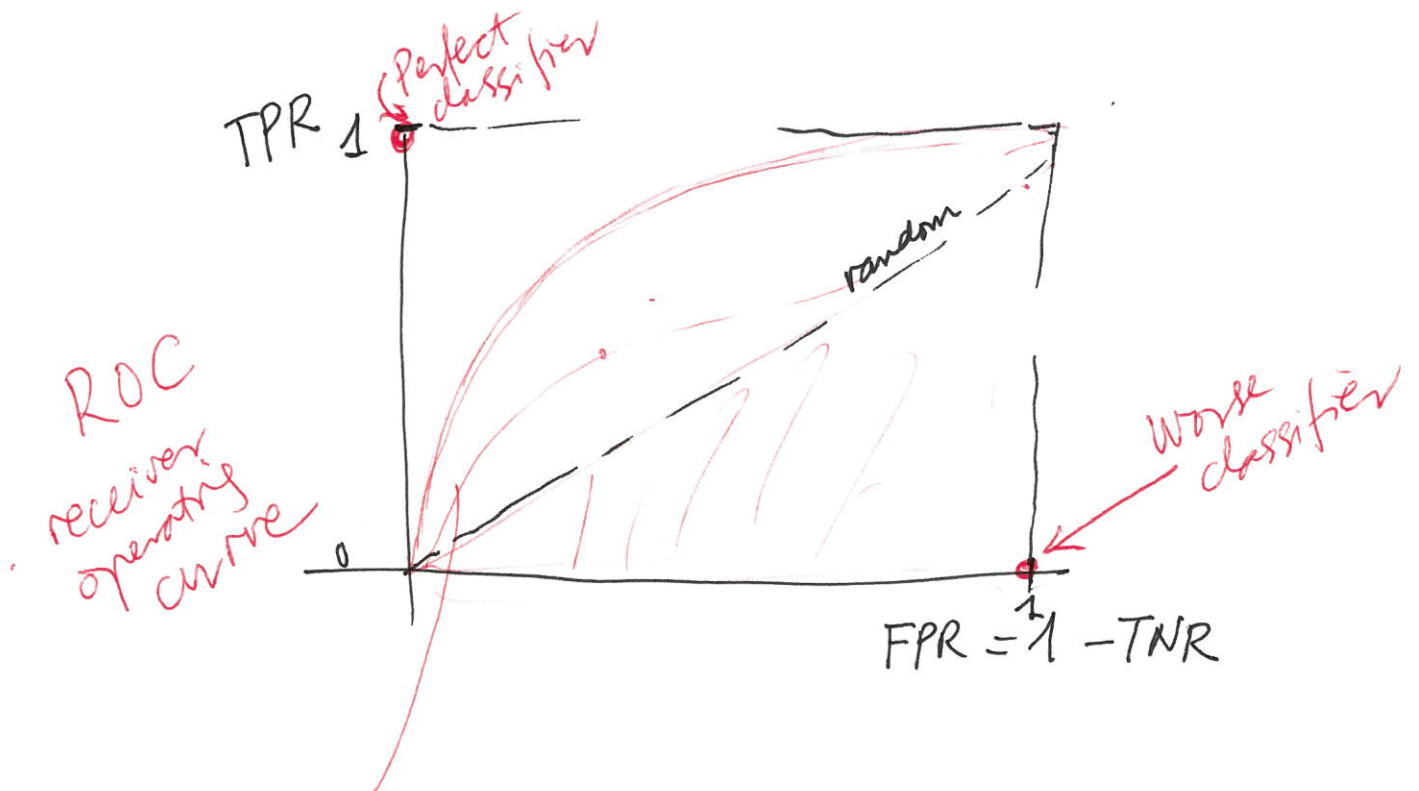
$$\left\{ \begin{array}{l} \text{Accuracy} = \frac{TP + TN}{N_{\text{validation}}} \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{Precision} = \frac{TP}{TP + FP} \end{array} \right.$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$F = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$



AUC : Area Under the Curve
is an overall
measure of quality

$$AUC > \frac{1}{2}$$

Better than random
over a range of
parameters