

# Training loss vs. prediction loss

Training loss:  $J_{\hat{\theta}}(\underline{\theta}) = \frac{1}{n} \sum_{i=1}^n L(h_{\hat{\theta}}(x_i; \underline{\theta}), y_i)$

hyperparameters      tuning parameters

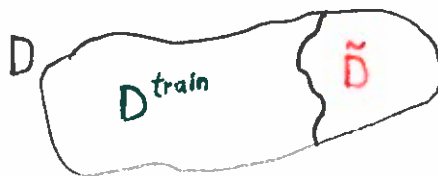
Prediction loss:  $J(h_{\hat{\theta}}(.; \hat{\underline{\theta}})) = \mathbb{E}_{x,y}[L(h_{\hat{\theta}}(x; \hat{\underline{\theta}}), y)]$

Expected prediction loss:  $\bar{J}(\mathcal{A}) = \mathbb{E}_{D^{\text{train}}}[J(\mathcal{A}(D^{\text{train}}))]$

Goal: Avoid overfitting and underfitting

Estimate prediction loss / expected prediction loss

Prediction loss: held-out set



Expected prediction loss: cross validation

$D_1^{\text{val}} \mid D_1^{\text{train}} \rightarrow h_{\hat{\theta}}(.; \hat{\underline{\theta}}(D_1^{\text{train}})) \quad \hat{J}_1$

$D_2^{\text{tr}} \mid D_2^{\text{val}} \mid D_2^{\text{ain}} \rightarrow h_{\hat{\theta}}(.; \hat{\underline{\theta}}(D_2^{\text{train}})) \quad \hat{J}_2$

$D_3^{\text{train}} \mid D_3^{\text{val}} \rightarrow h_{\hat{\theta}}(.; \hat{\underline{\theta}}(D_3^{\text{train}})) \quad \hat{J}_3$

$\hat{J}(\mathcal{A}) = \text{CV}$

2020 L14 (1)

## Estimating generalisation performance with hyperparameter selection

Idea: - Nest method for estimating generalisation performance to avoid overfitting

- Use held-out or CV in hyperparameter optimisation:

For given hyperparameters  $\underline{\lambda}$ , estimate generalisation performance, then we select best  $\underline{\lambda}$

$\rightarrow$  now  $\underline{\lambda}$  can overfit, so yet another held-out set necessary

Algorithm:

1. Split  $D$  into  $D^{hyp}$  and  $D^{test}$  (e.g. 80:20)

Lock away  $D^{test}$  until step 4.

2. Optimise hyperparameters  $\underline{\lambda}$  on  $D^{hyp}$ :

held-out

Split  $D^{hyp}$  into  $D^{train}$   
and  $D^{val}$

$$\hat{\lambda} = \underset{\lambda}{\operatorname{argmin}} PL(\lambda)$$

where

$$PL(\lambda) = \hat{J}(\hat{h}_{\lambda}, D^{val})$$

$$\text{and } \hat{h}_{\lambda} = \mathcal{A}(D^{train})$$

cross validation

Split  $D^{hyp}$  into  $k$  folds  
and compute  $CV$

for each  $\lambda$

$$\hat{\lambda} = \underset{\lambda}{\operatorname{argmin}} EPL(\lambda)$$

where

$$EPL(\lambda) = CV_{\lambda} = \hat{J}(\mathcal{A}_{\lambda})$$

3. Reestimate  $\Theta$  using  $\hat{\lambda}$ :  $\hat{h} = \mathcal{A}_{\hat{\lambda}}(D^{hyp})$

4. Estimate prediction loss:  $\hat{J} = \hat{J}(\hat{h}; D^{test})$

5. Reestimate  $\hat{h}$  on all data:  $\hat{h} = \mathcal{A}_{\hat{\lambda}}(D)$

Do not estimate prediction loss on all data!

"Two-times held out"

"CV and held-out"

## Loss functions

- Regression:

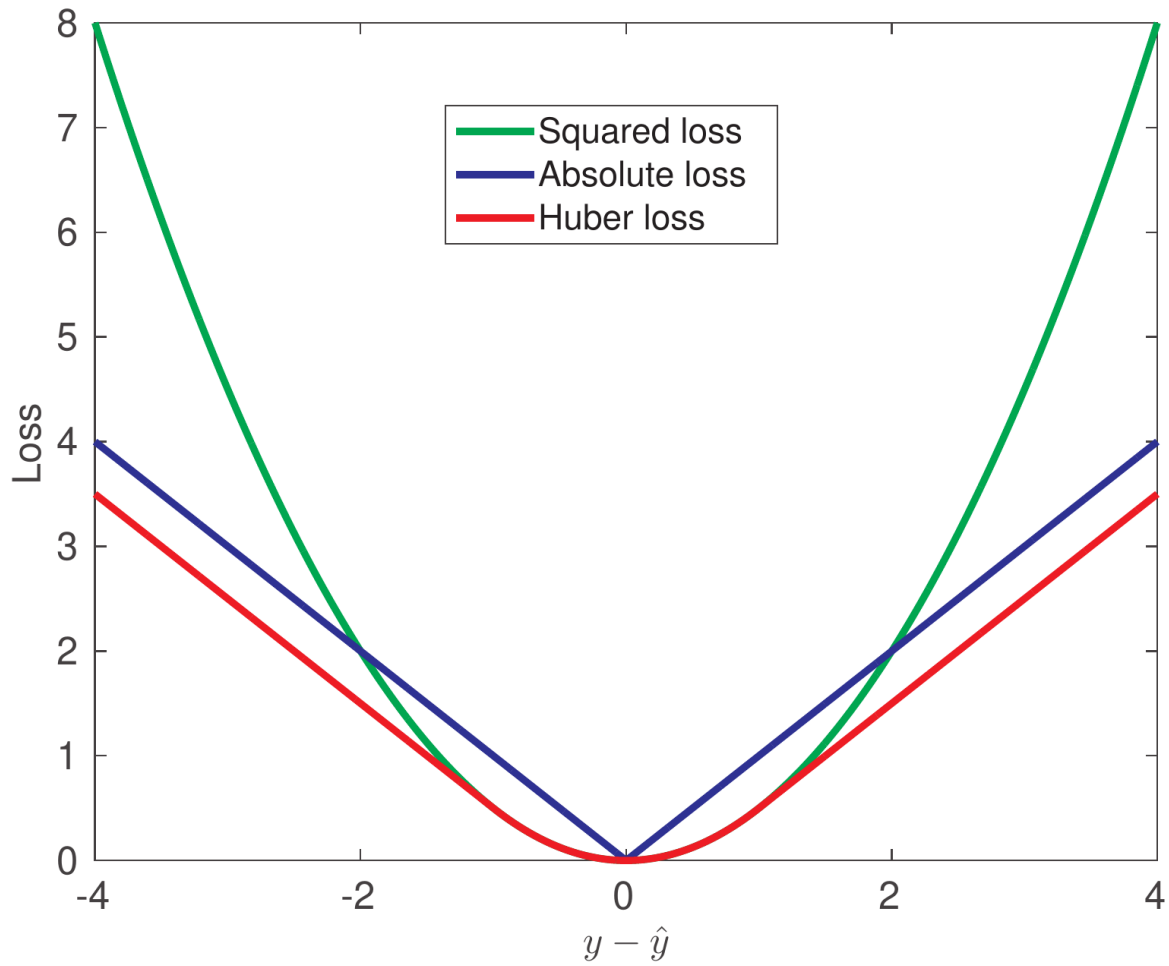
$$L(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2 \quad \text{"square loss"}$$

$$L(\hat{y}, y) = |\hat{y} - y| \quad \text{"absolute loss"}$$

$$L(\hat{y}, y) = \begin{cases} \frac{1}{2}(\hat{y} - y)^2 & ; \text{ if } |\hat{y} - y| < \delta \\ \delta|\hat{y} - y| - \frac{1}{2}\delta^2 & ; \text{ otherwise} \end{cases}$$

"Huber loss"

# Loss Functions



## - Classification loss functions

For  $k$  different classes  $\{1, \dots, k\}$

loss function  $L(\hat{y}, y) \leadsto$  matrix

$$L = \begin{pmatrix} L(1,1) & L(1,2) & \dots & L(1,k) \\ L(2,1) & \dots & \dots & \vdots \\ \vdots & \dots & \dots & \vdots \\ L(k,1) & \dots & \dots & L(k,k) \end{pmatrix}$$

$L(i,i)$  typically zero, other entries  $> 0$

Special case: "Zero-one loss"

$$L(i,j) = \begin{cases} 0 & ; i=j \\ 1 & ; i \neq j \end{cases}$$

Expectation:

$$\begin{aligned} J(h) &= \mathbb{E}_{\underline{x}, y} [L(h(\underline{x}), y)] = \mathbb{E}_{\hat{y}, y} [L(\hat{y}, y)] \\ &= \sum_{i,j} p(i,j) L(i,j) = \sum_{i \neq j} p(i,j) = P(\hat{y} \neq y) \\ &= P(h(\underline{x}) \neq y) \end{aligned}$$

"misclassification" or "error rate" 2020 L14 (5)

For two classes (i.e.  $K=2$ ): let  $\hat{y}, y \in \{-1, 1\}$

"true positive rate":  $p(\hat{y}=1|y=1) = \frac{p(\hat{y}=1, y=1)}{p(y=1)}$

("sensitivity", "hit rate", "recall")

"true negative rate":  $p(\hat{y}=-1|y=-1) = \frac{p(\hat{y}=-1, y=-1)}{p(y=-1)}$

("specificity")

"false positive rate":  $p(\hat{y}=1|y=-1) = \dots$

("type 1 error")

"false negative rate":  $p(\hat{y}=-1|y=1) = \dots$

("type 2 error")

Loss function that penalises false positive and false negative rates:

$$L = \begin{pmatrix} 0 & \frac{1}{p(y=-1)} \\ \frac{1}{p(y=1)} & 0 \end{pmatrix}$$

Then:

$$\begin{aligned} J(h) &= \sum_{i,j} p(i,j) L(i,j) \\ &= \frac{p(1,-1)}{p(y=-1)} + \frac{p(-1,1)}{p(y=1)} \\ &= p(\hat{y}=1 | y=-1) + p(\hat{y}=-1 | y=1) \end{aligned}$$

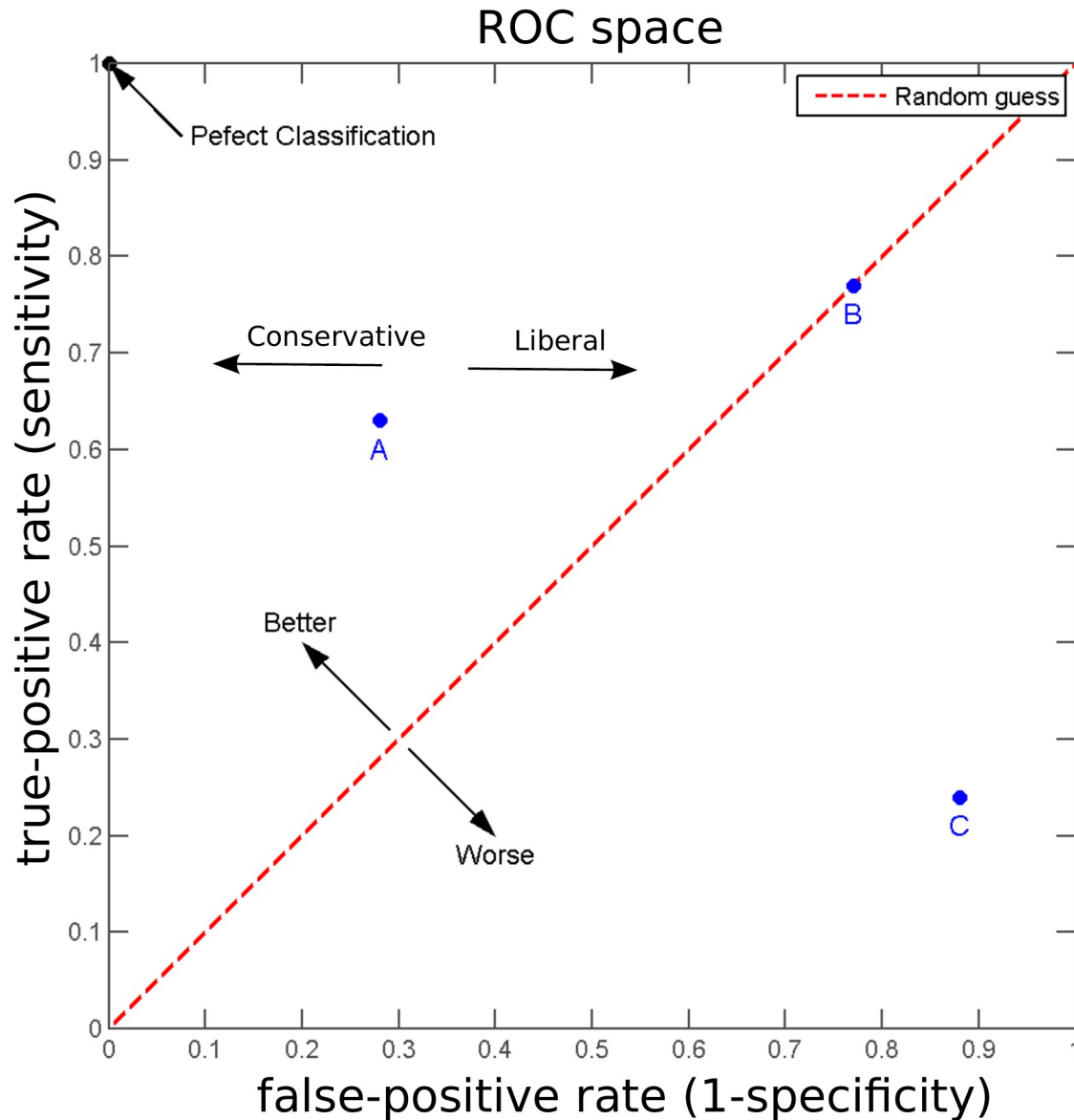
Useful for imbalances in  $y=1$ ,  $y=-1$

Generally, trade-off between true-positive and false-negative rates

Visualised by "receiver-operating characteristic" (ROC)



# Receiver Operating Characteristic



## Differentiable loss functions for classification

Here : binary classification only

$$\hat{y}(\underline{x}) = \text{sign}(h(\underline{x})) \in \{-1, 1\}$$

where  $h$  real-valued

$\rightarrow$  correct classification of  $\underline{x} \Leftrightarrow y h(\underline{x}) > 0$

$y h(\underline{x})$  "margin", similar to  $y - h(\underline{x})$  in regression

Zero-one loss recovered for  $L(h(\underline{x}), y) = \begin{cases} 1 & ; y h(\underline{x}) \leq 0 \\ 0 & ; \text{otherwise} \end{cases}$

Logistic regression :

$$p(y=1 | \underline{x}; h) = \frac{1}{1 + \exp(-h(\underline{x}))} \quad p(y=-1 | \underline{x}; h) = \frac{1}{1 + \exp(h(\underline{x}))}$$

Relation to loss :

maximise log-likelihood  $\Leftrightarrow$  minimisation log-loss

# Classification Loss Functions that Operate on Margin

$$L(h(\mathbf{x}), y) = (h(\mathbf{x}) - y)^2 = (1 - yh(\mathbf{x}))^2 \quad (\text{square loss})$$

$$L(h(\mathbf{x}), y) = \log(1 + \exp(-yh(\mathbf{x}))) \quad (\text{log loss})$$

$$L(h(\mathbf{x}), y) = \exp(-yh(\mathbf{x})) \quad (\text{exponential loss})$$

$$L(h(\mathbf{x}), y) = \max(0, 1 - yh(\mathbf{x})) \quad (\text{hinge loss})$$

$$L(h(\mathbf{x}), y) = \max(0, 1 - yh(\mathbf{x}))^2 \quad (\text{square hinge loss})$$

$$L(h(\mathbf{x}), y) = \begin{cases} -4yh(\mathbf{x}) & \text{if } yh(\mathbf{x}) < -1 \\ \max(0, 1 - yh(\mathbf{x}))^2 & \text{otherwise} \end{cases} \quad (\text{Huberised square hinge loss})$$

