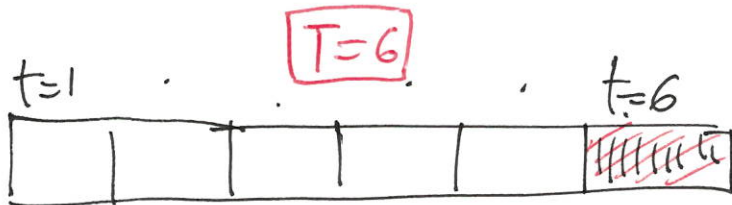


k -NN : k -nearest neighbours
as a predictor.

T -fold cross validation



Split into
 T -folds

All samples $S = \{ \vec{x}^{(i)}, y^{(i)} \} \quad i=1, \dots, N$

$$S = \bigcup_{t=1}^T S_t \quad |S_t| = \frac{|S|}{T}$$

$$\bar{S}_t = S - S_t$$

Equal size, properly
sampled

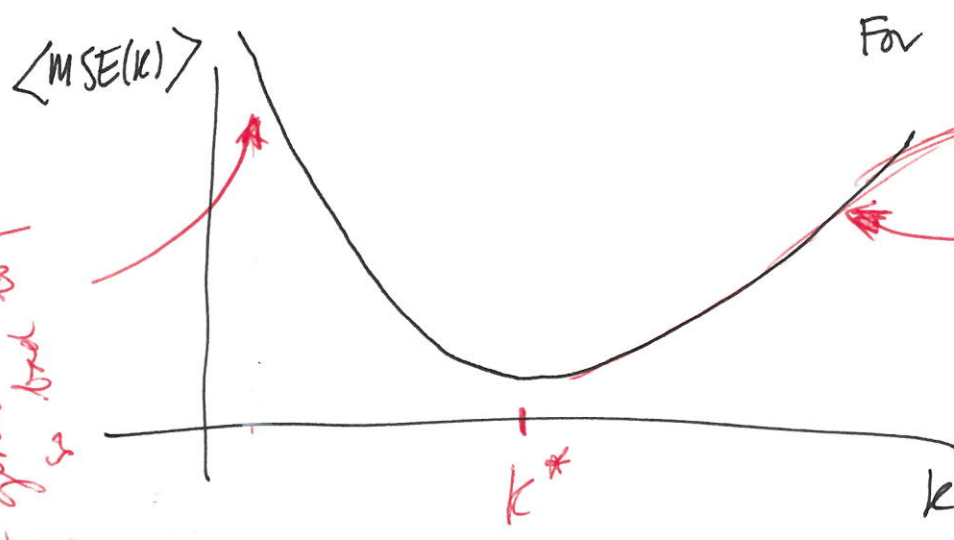
Different model for each \bar{S}_t \rightarrow $\hat{f}_{\bar{S}_t; k}$ $t=1, \dots, T$

Predict \bar{S}_t with the model obtained from \bar{S}_t \Rightarrow
$$\sum_{i \in \bar{S}_t} \left[\hat{f}_{\bar{S}_t}(\vec{x}^{(i)}) - y^{(i)} \right]^2 = \text{MSE}_t$$

$$\langle \text{MSE}(k) \rangle = \frac{1}{T} \sum_{t=1}^T \text{MSE}_t$$

This is a characteristic of
how well the model
predicts out-of-sample

k too small:
Not generalizable to off-sample
Not generalizable to test set



For kNN

k too large:
Predictor too generic

Typical: $T \sim 5, 10$ Machine Learning

If $T = N$ LOO
leave - one - out

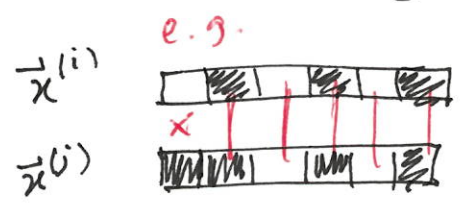
- Good estimate
- computationally expensive
- Overkill.

k-NN can also be applied to discrete descriptors

\Rightarrow If

$$\vec{x}^{(i)} = \{0, 1\}^P, y^{(i)} \in \mathbb{R}$$

Hamming distance: $d_H(\vec{x}^{(i)}, \vec{x}^{(j)})$



$$\text{distance} = \frac{1}{P} \sum_{m=1}^P I(x_m^{(i)} \neq x_m^{(j)})$$

Distance is the key:

The rest is the same

$$N_k(\vec{x}^{in}) \text{ and } \hat{y} = \frac{1}{K} \sum_{i \in N_k(\vec{x}^{in})} y^{(i)}$$

\rightarrow Other methods use distances \Rightarrow kernels to have local models

Classification problem with KNN

$\vec{x}^{(i)}$
 \uparrow
 $\left\{ \begin{array}{l} \vec{x}^{(i)} \in \mathbb{R}^p \\ \vec{x}^{(i)} \in \{0,1\}^p \end{array} \right\}$

$y^{(i)} \in C_i$
 $y^{(i)}$ discrete.
 $C_q = \{C_1, \dots, C_q\}$
 $i = 1, \dots, N$

$d(\vec{x}^{(i)}, \vec{x}^{(j)})$

As long as we have the distance:

- same ✓ ① Choose k (number of neighbors)
- same ✓ ② For \vec{x}^{in} find $N_k(\vec{x}^{in})$
- Predictor is different for classification: new ③ $\hat{y} \in C_q$ e.g. by majority rule.



$$k=4 \quad \hat{y}(\vec{x}^{in}) \rightarrow \triangle$$

$$k=7 \quad \hat{y}(\vec{x}^{in}) \rightarrow \square$$

Other predictors are used, including those that give a probabilistic output:

$$\forall i \in N_k(\vec{x}^{in}) \text{ compute } \hat{Q}(\vec{x}^{in}) = \frac{1}{k} \sum_{i \in N_k(\vec{x}^{in})} \mathbb{I}(y^{(i)} \in C_q)$$

$$\mathbb{I}(\cdot) = \begin{cases} 1 & \text{if true} \\ 0 & \text{otherwise} \end{cases}$$

$$\vec{Q}(\vec{x}^{in}) = \begin{pmatrix} \# \text{ of points in } C_1 \\ \vdots \\ \# \text{ of points in } C_g \end{pmatrix} \frac{1}{K}$$

$g \times 1$

↓

Each component Q_i : Probability of \vec{x}^{in} belonging in class C_i

Standard KNN applies an argmax to this probability:

$$\hat{y} = \underset{g}{\text{argmax}} \vec{Q}(\vec{x}^{in})$$

↑

Chooses the class with the maximum probability.

Global models for classification.

$$\vec{x}^{(i)} \in \mathbb{R}^p ; \quad y^{(i)} \in \{0, 1\} \quad i=1, \dots, N$$

Binary.

Classic model: Logistic ~~"regression"~~
↳ classification.

Motivation: log odds are linear combination of descriptors

$$\beta_0 + \vec{x}^T \cdot \vec{\beta} = \log \frac{P(y=1)}{P(y=0)} = \log \frac{P(y=1)}{1 - P(y=1)}$$

$\vec{x}^T \cdot \vec{\beta}$ $(p+1) \times 1$

$$\frac{P(y=1)}{1 - P(y=1)} = e^{\vec{x}^T \cdot \vec{\beta}} \Rightarrow P(y=1) = \frac{1}{1 + e^{-\vec{x}^T \cdot \vec{\beta}}} = h_{\vec{\beta}}(\vec{x})$$

Bernoulli variable with probability of success $P(y=1)$

$$\hookrightarrow P(Y=y | \vec{x}, \vec{\beta}) = h_{\vec{\beta}}(\vec{x})^y (1 - h_{\vec{\beta}}(\vec{x}))^{1-y}$$

If we assume independence in our samples; we can factorize probabilities to get:

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

log-likelihood:

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

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