

## error analysis of a classifier

- exact analytical results are usually impossible for real world problems
- asymptotic results for  $N \rightarrow \infty$  can be derived [may be unrealistic for small  $N$ ]
- worst case bounds for fixed  $N$  ("how bad can it be, when we get an unfortunate training set")
- [often too pessimistic, because worst case rarely happens]

- empirical error estimation: make experiments with the trained algorithm (i.e. algorithm is now fixed)

generalization error for new data ( $\equiv$  test error)  $\rightarrow$  average over all possible data

$$Err = \mathbb{E}_{x, y \sim p(x, y)} \left[ \underbrace{p(\text{error}(x, y) \mid \text{fixed algorithm})}_{\hat{f}(x) \neq y} \right]$$

$\hat{f}(x)$  after training

training tries to minimize the training error

$$err = \mathbb{E}_{x, y \sim TS} \left[ \underbrace{p(\text{error} \mid \text{current alg})}_{\text{change the alg as long as err reduces}} \right]$$

$\Rightarrow$  in general  $err < Err$ , sometimes  $err \ll Err$

because we made  $err$  as small as possible ("training")

$Err - err =$  "optimism" of our predictor

big optimism  $\hat{=}$  over-fitting (extreme case: memorize training data, but has no idea about new data)

we need to ensure that  $Err$  is small enough for our predictor to be applicable in practice



four levels of empirical error analysis:

- (1) (luxury solution, for very critical applications): defines test procedures independently of training data and algorithms, example: certification for autonomous cars  
VDA Leibnizinitiative ~~www.vda1.de~~ ~~www.vda1.de~~

- (2) (very good, for scientific benchmarks): benchmark website, which offers TS for some problem, and ~~the~~ a set of test  $X$ 's ( $Y$ 's secret)  
- people train their algorithm on TS and apply to test set, upload  $\hat{Y}$   
- benchmark managers compare  $\hat{Y}$  with secret  $Y^*$   $\Rightarrow$  leader board  
example: CREMI Neurosegmentation Challenge

very popular in the life sciences (medicine, biology, neuroscience, ...)

- (3) (standard) provide TS and test sets with revealed ground truth  
 $\Rightarrow$  use to design better algorithms (people could cheat, but rarely do)  
example: MNIST dataset (handwritten digits by US Postal service for automated ZIP-code detection/reading)

- (4) (fall back: if you do not have a test set): split TS into a training and a test subset, clever version m-fold cross-validation

alg: (0) split the TS into  $M$  subsets of size  $\frac{N}{M}$  (roughly)  
 $\propto$  "folds"

- in every iteration, train from scratch on TS \ fold  $m$
- (1) for each  $m = 1 \dots M$   
(a) use subset  $m$  as test set, train on remaining data  
 $\Rightarrow$  test error  $\hat{Err}_m = \frac{\# \text{errors}}{N_m}$   $N_m \leftarrow$  size of subset  $m$
- (2) compute test error  $\hat{Err} = \frac{1}{M} \sum_{m=1}^M \hat{Err}_m$



## Improvements for the Nearest Neighbor classifier

- for large data sets ( $N$  sig) and/or high-dimensional features ( $D$  sig), naive slow complexity of "naive" implementation:
  - distance comp. takes  $D$  steps (look at every feature)
  - finding the neighbor  $N$  steps

$O(N \cdot D)$  steps

reduce  $N$  and  $D$

- ① reduce  $D$ : - feature selection: drop all features that are not very influential on the outcome

backward: - start with all features

- try, how much results worsen, when each feature is dropped in turn ~~and~~

$\Rightarrow D$ -times cross-validation

- remove feature with lowest performance drop

- repeat (drop as many features as possible)

forward: - start with the "best" single feature

- add features one at a time, until good performance

- feature computation:  $\tilde{X}_i = f(X_i)$   $f: \mathbb{R}^D \rightarrow \mathbb{R}^{D'}$   $D' < D$

- by hand: BMI (1-dimensional) =  $f(\underbrace{\text{weight, height}}_{2\text{-dimensional}})$

- learn informative features  $\Leftrightarrow$  later (chapter "unsupervised learning")

- ② reduce  $N$ : only keep very informative / typical & training influences:

- exact algorithm: compute the Voronoi diagram of all training data,  $O(N^{D/2})$

- prediction (position of decision boundary) unchanged, if drop instance ~~will~~ not at boundary



- greedy approximation:
  - maintain active set of instances, initially empty
  - for each  $i = 1, \dots, N$ : if instance  $i$  is incorrectly classified by the current active set  $\Rightarrow$  add  $(x_i, y_i)$  to active set otherwise, drop it

- clustering algorithm:
  - group similar instances together ("cluster")
  - choose one representative per cluster

later in chapter "unsupervised l."

- relax criterion: choose an approximate nearest-neighbour



$$d' \leq d'' \leq \begin{cases} d' + \epsilon & \text{bounded additive error} \\ d'(1 + \epsilon) & \text{multiplicative error} \end{cases}$$

large field of research: how to guarantee that  $x''$  is not too far off and the search is much faster than  $O(N)$

- clever data structure(?): in 1-D, a binary search tree can find an instance in  $O(\log N)$  time

generalisation to higher dimensions: k-D tree (or variants thereof)

search complexity:  $O(D \cdot \min(N, 2^D \log N))$

in practice: only workable faster for  $D \leq 20$

improved variants

$$D \leq 30 \dots 50$$

↑ tree search  
↑ sequential search

- two-stage procedure:
  - use approximate alg. to filter out far-neighbors
  - use exact alg. on remaining near neighbor candidate set

e.g. locality-preserving hashing

$$O(N') \quad N' \ll N$$



- nearest neighbor alg. is not consistent, ~~not~~ i.e. does not converge to the optimal Bayes classifier even when  $N \rightarrow \infty$

$\Rightarrow$  make better use of the TL via k-nearest neighbor alg.  
M-nearest neighbor alg.

- instead of looking for nearest neighbor, find the M nearest neighbors

- count fraction of class k instances in this set  $\hat{p}_k = \frac{M_k}{M} = \hat{p}(Y=k|X)$

$M_k$ : # of  $Y=k$  in M-neighbor set

- return the vector of posteriors  $(\hat{p}_1, \dots, \hat{p}_C)$  or  $\hat{Y} = \arg \max_k \hat{p}_k$

theorem: M-nearest neighbor rule is consistent if  $\underline{M \rightarrow \infty}$ , as  $\underline{N \rightarrow \infty}$ ,

and  $\underline{\frac{M}{N} \rightarrow 0}$  as  $N \rightarrow \infty$  (ex:  $M = \log N$ )  $\swarrow$  ensures that  $\hat{p}_k(X) \rightarrow p^*(Y|X)$

ensures that the radius of the M-neighbor set in feature space  $\rightarrow 0$ , i.e. condenses around query point X

- nearest neighbor alg. only makes sense, if  $d(X, X')$  distance for neighbor search represents the semantics of our application: if  $d(X, X') \leq 0 \Rightarrow$  events X, X' should be similar ~~for~~ from application perspective. This is difficult for more realistic applications.

$\Rightarrow$  metric learning, similarity learning  $\Rightarrow$  own research field