

Nearest Neighbor classification

• example: $Y \in \{\text{small, medium, tall}\}$
threshold classifier

$$Y = \begin{cases} \text{small} & \text{if } x < 1.6 \text{ m} \\ \text{medium} & 1.6 \leq x < 2.0 \text{ m} \\ \text{tall} & x \geq 2.0 \text{ m} \end{cases}$$

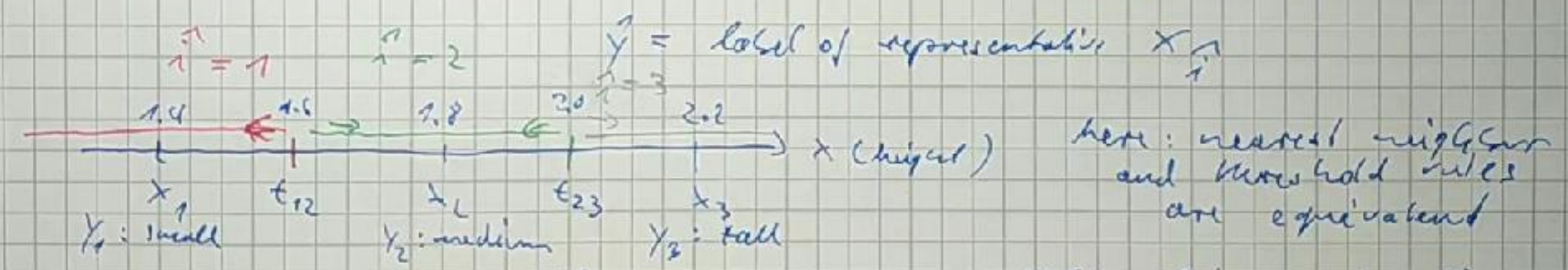
t_{12}
 t_{23}

alternatively: define typical representatives for each class:

$$x_{\text{small}} = 1.4 \text{ m} = x_1, \quad x_{\text{medium}} = 1.8 \text{ m} = x_2, \quad x_{\text{tall}} = 2.2 \text{ m} = x_3$$

\hat{Y} = label of the closest representative for a given x

two stages: $\hat{i} = \arg \min_{i \in \text{representatives}} \|x - x_i\|$ nearest-neighbor



• nearest neighbors work in any dimensions, when a suitable distance function is available: $d(x, x')$ ($d(x, x) = 0$, $d(x, x') > 0$ if $x \neq x'$)
 $d(x, x') + d(x', x'') \geq d(x, x'')$ triangle inequality

\Rightarrow given a set of representatives $TS = \{(x_i, Y_i)\}_{i=1}^N$ Y_i = true label of instance i
nearest neighbor rule $\hat{i} = \arg \min_i d(x, x_i)$ x = test feature return $\hat{Y} = Y_{\hat{i}}$

- advantage: Standard distance functions often work reasonably well / most popular:

Euclidean distance $d(x, x') = \|x - x'\|_2 = \sqrt{\sum_{j=1}^D (x_j - x'_j)^2}$

or weighted Euclidean distance

$$d(x, x') = \sqrt{\sum_{j=1}^D (x_j - x'_j)^2 \cdot w_j}$$

↑ weight of feature j

intuitive meaning: measure features in units such that the magnitudes of different features are comparable

example: Body-mass-index: if height is measured in cm or km instead of m, it doesn't work

standard weight: is inverse variance $w_j = \frac{1}{\sigma_j^2}$ ← variance of feature j

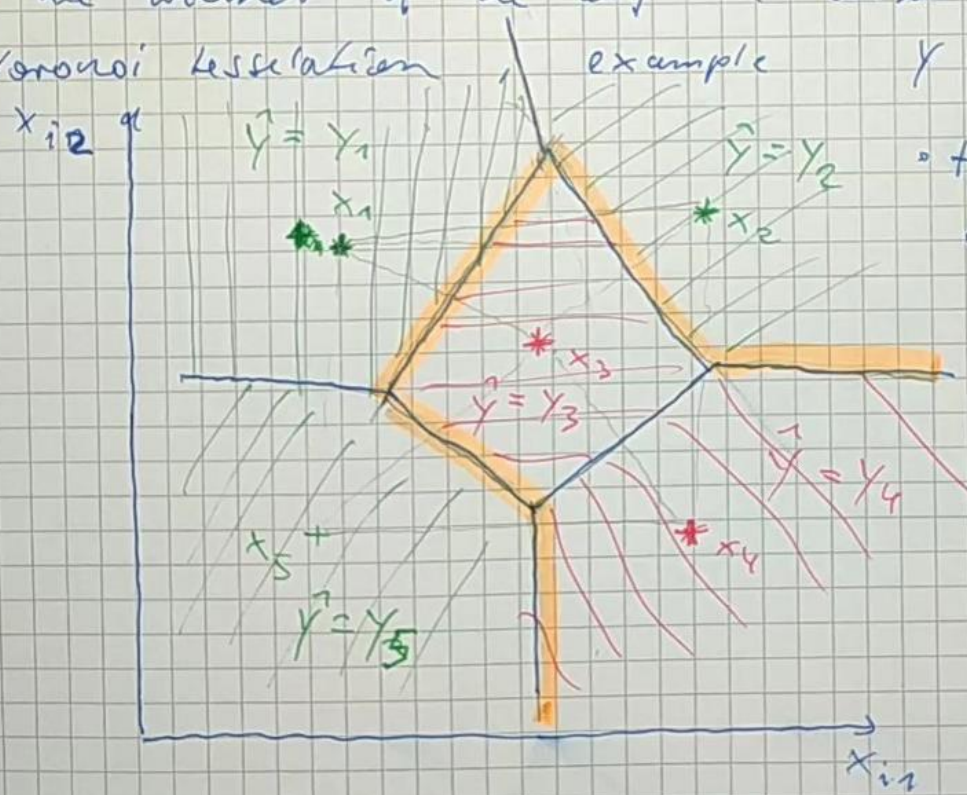
- but if standard distance do ~~not~~ not work, defining a good distance $d(x, x')$ is hard (as hard as defining a loss $\ell = g(x)$) \Rightarrow metric learning

Nearest Neighbor Classifier Recap:

- Training is trivial: just store ("memorize") the $TS = \{(X_i, Y_i)\}_{i=1}^N$
- Testing / Inference:
 - need a suitable distance function for features $d(X, X')$
 - for test features X : find index of the nearest (= most similar) training instance: $\hat{i} = \arg \min_i d(X, X_i)$

- where the label of instance \hat{i} : $\hat{Y} = Y_{\hat{i}}$
- the training set partitions the feature space into regions, where each i is the winner of the $\arg \min$ distance: Voronoi regions

\Rightarrow Voronoi tessellation, example $Y \in \{0, 1\}$ $X_i \in \mathbb{R}^D$



- to construct Voronoi regions: draw bisector between neighboring points
- merge all regions of the same label into one big ~~single~~ region per label "decision regions"
- the boundary between decision regions: "decision boundary"

• How good can the Nearest neighbor rule perform?

~~three ways to analyse the error~~

• the training error of a NN cl. is meaningless: whenever query point x equals a training point x_i : $x = x_i$ for some i , then $\hat{y} = y_i$

\Rightarrow the label of the response is always correct $\hat{=}$ training error $= 0$

• this doesn't imply that the test error / generalization error is also small

• three ways to analyse the generalization error

- compute an exact analytical formula - usually only possible for toy problems
- compute an analytic asymptotic formula, i.e. for infinitely large TS, $N \rightarrow \infty$
- estimate it empirically, either with separate test data set or cross validation (2 fold test set)

Analytic error analysis for a toy problem

priors: $p(Y=0) = p(Y=1) = 1/2$ $x \in \{0, 1\}$

likelihoods $p(x|Y=0) = 1-2x$, $x \in [0, 1]$

$p(x|Y=1) = 2x$

evidence: $p(x) = 1$

posteriors: $p(Y=0|x) = 1-x$

$p(Y=1|x) = x$

Reduce the NN cl. for $N=2$

to threshold classification

one training instance with either label

e.g. $TS = \{(x_1, y_1=0), (x_2, y_2=1)\}$

cases: - if $x_1 < x_2$: \Rightarrow type A classifier

- if $x_1 > x_2$: \Rightarrow type B anti-classifier

with threshold $t = \frac{x_1 + x_2}{2}$

\Rightarrow calculate prob. of different TS \Rightarrow expected NN error

• How good can the Nearest neighbor rule perform?

~~three ways to analyse the error~~

• the training error of a NN cl. is meaningless: whenever query point x equals a training point x_i : $x = x_i$ for some i , then $\hat{y} = y_i$

\Rightarrow the label of the response is always correct $\hat{=}$ training error $= 0$

• this doesn't imply that the test error / generalization error is also small

• three ways to analyse the generalization error

- compute an exact analytical formula - usually only possible for toy problems
- compute an analytic asymptotic formula, i.e. for infinitely large TS, $N \rightarrow \infty$
- estimate it empirically, either with separate test data set or cross validation (2 fake test set)

Analytic error analysis for a toy problem

priors: $p(Y=0) = p(Y=1) = 1/2$ $Y \in \{0, 1\}$

likelihoods $p(X|Y=0) = 2 - 2x$, $x \in [0, 1]$

$p(X|Y=1) = 2x$

evidence: $p(x) = 1$

posteriors: $p(Y=0|x) = 1 - x$

$p(Y=1|x) = x$

Reduce the NN cl. for $N=2$

to threshold classification

one training instance with either label

e.g. $TL = \{(X_1, Y_1=0), (X_2, Y_2=1)\}$

cases: - if $x_1 < x_2$: \Rightarrow type A classifier

- if $x_1 > x_2$: \Rightarrow type B anti-classifier

with threshold $t = \frac{x_1 + x_2}{2}$

\Rightarrow calculate prob. of different TS \Rightarrow expected NN error

probability of TS: pair $(x_1, y_1=0)$ has prob $p(x_1 | Y=0)$
 pair $(x_2, y_2=1)$ has prob. $p(x_2 | Y=1)$

$$\begin{aligned}
 & \mathbb{E}_{x_1, x_2} \left[p(\text{error} | \text{rule} = A, t = \frac{x_1 + x_2}{2}) \right] \\
 & \quad x_1 < x_2 \\
 &= \int_0^1 \underbrace{p(x_1 | Y=0)}_{= 2 - 2x_1} \cdot \int_{x_1}^1 \underbrace{p(x_2 | Y=1)}_{x_1 = 2x_2} \underbrace{p(\text{error} | \text{rule} = A, t = \frac{x_1 + x_2}{2})}_{= \left(\frac{x_1 + x_2}{2} - \frac{1}{2} \right)^2 + \frac{1}{4}} dx_2 dx_1 \\
 &= \frac{83}{360} \quad (\text{as per Wolfram Cloud})
 \end{aligned}$$

$$\begin{aligned}
 & \mathbb{E}_{x_1, x_2} \left[p(\text{error} | \text{rule} = B, t = \frac{x_1 + x_2}{2}) \right] \\
 & \quad x_1 > x_2 \\
 &= \int_0^1 p(x_1 | Y=0) \int_0^{x_1} p(x_2 | Y=1) p(\text{error} | \text{rule} = B, t = \frac{x_1 + x_2}{2}) dx_2 dx_1 \\
 &= \frac{43}{360} \quad \Rightarrow \text{expected error of the NAV-rule with } N=2 : \\
 & \quad \text{error: } \frac{83}{360} + \frac{43}{360} = \frac{7}{20} = 35\%
 \end{aligned}$$

if you repeat the experiment + draw a TS of $N=2$, test the resulting NAV rule
 infinitely often, the average error is 35% (compare Bayes: 25%, guessing: 50%)

Asymptotic analytic error analysis: $N \rightarrow \infty$ (infinitely many training instances)

Definition: a classification algorithm is called consistent, if it converges to the Bayes classifier when $N \rightarrow \infty$

Result: NN rule is not consistent, but only a factor 2 off
 i.e. $P_{N \rightarrow \infty}(\text{error NN}) \leq 2 \cdot p^*(\text{error Bayes})$

[full derivation: Duda, Hart, Stork, Chapter 4.5]

$$p(\text{error} | X, X') \quad X' = \arg \min_{x_i} d(x, x_i)$$

query feature \uparrow nearest training instance \uparrow

define probability $p(X' | X)$ that X' is the nearest training point of X

$$(1) \quad p(\text{error} | X) = \int p(\text{error} | x, x') p(X' | x) dx' \quad \text{dxx}$$

"marginalization over x' " \Rightarrow eliminate the unknown variable x'

[observation: when $N \rightarrow \infty$, a training example is at every x

$$\Rightarrow p(X' | X) \xrightarrow{N \rightarrow \infty} \delta(x - x') \quad \text{if } p(x) \text{ is smooth.}]$$

$$(2) \quad p(\text{error} | x, x') = 1 - p(\text{correct} | x, x') = 1 - \sum_{k=1}^c p(Y=k, Y'=k | x, x')$$

truth prediction

assumption: query point is selected independently of the TS, but from same distribution

$$\Rightarrow p(Y=k, Y'=k | x, x') = p(Y=k | x) \cdot p(Y'=k | x')$$

$$p(\text{error} | x, x') = 1 - \sum_{k=1}^C p(Y=k|x) p(Y'=k|x')$$

$$(2) \quad p_{N \rightarrow \infty}(\text{error NN} | x) = \lim_{N \rightarrow \infty} p(\text{error NN} | x)$$

$$= \int \left(1 - \sum_{k=1}^C p(Y=k|x) \cdot p(Y'=k|x') \right) \underbrace{\left(\lim_{N \rightarrow \infty} p(x'|x) \right)}_{=\delta(x-x')} dx'$$

$$\boxed{p_{N \rightarrow \infty}(\text{error NN} | x) = 1 - \sum_{k=1}^C p(Y=k|x)^2}$$

\Rightarrow Gini impurity of the ^{label} distribution at x

extreme cases: - best case: only one label can occur for feature x

$$p(Y=Y^* | x) = 1, \quad p(Y \neq Y^* | x) = 0$$

$$\Rightarrow p_{N \rightarrow \infty}(\text{error NN} | x) = 0 \Rightarrow \text{perfect classification}$$

- worst case: all labels are equally ~~likely~~ likely at feature x :

$$p(Y=k | x) = \frac{1}{C} \quad \text{for all } k$$

$$\Rightarrow p_{N \rightarrow \infty}(\text{error NN} | x) = 1 - C \cdot \frac{1}{C^2} = \frac{C-1}{C} \geq \frac{1}{2} \quad \text{pure guessing}$$

typical case: for a given feature, the label probabilities are

$$\text{all different} \quad p(Y=k | x) \quad \text{arbitrary (but positive!)} \quad \sum_k p(Y=k | x) = 1$$

