

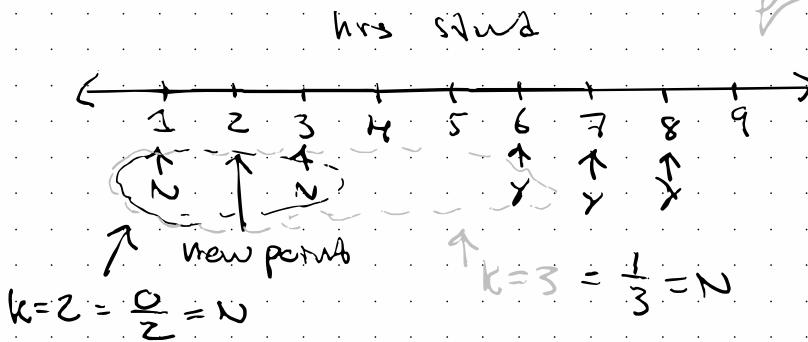
Supervised Learning in Python

- Classification Algorithms in this course:
 - ↳ CNN, N. Bayes, Decision Tree, Perceptron
 - ↳ Pros vs Cons
- Practical topics:
 - hyperparameters, CV, feature extraction & selection

K-nearest neighbors

example: will you pass an exam?

↳ find students
who studied for some
time & estimate outcome based on their
based on distance.



→ Prediction pseudocode

predict(x_0)

closest-dist = inf

" -class = -1

for x, y in train-dat

$d = \text{dist}(x, x_0)$

if $d < \text{closest-dist}$

closest-dist = d ,

closest-class = y

return closest class

↳ Notice how calculation gets more brutal w/ large datasets.

↳ Optimize performance w/ sorting algorithms

Break ties?

↳ averages or pick at random

How to choose k ?

↳ Cross validation

→ When KNN can fail

↳ data, grid of alternating dots

Always will classify the other class (due to alternation)

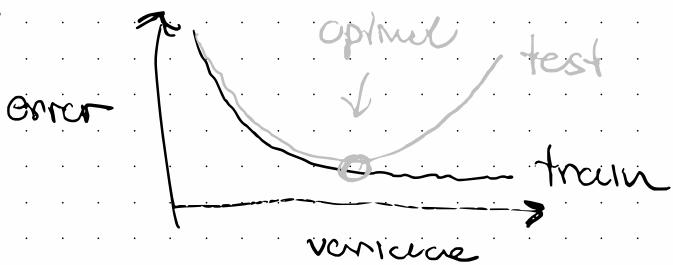


→ Effect of K

↳ Small values of k create smooth boundaries. (high bias)

↳ Large vals give flexible boundaries
(high var)

Need to balance the two so you don't overfit.



→ knn pseudocode using matrix algebra

X has shape $m \times n_x$ (100×2)

\hat{y} has shape $m_y \times n_y$ (100×1)

$\hat{X}_{new} \in (1, 2)$

as many examples as wanted

$X_{newz} \in \hat{X}_{new} . repeat(100, 2)$

$$D \triangleq = X^2 - X_{new}^2$$

$\hat{d} = \hat{d}$ add along rows

argmin of \hat{d} gives closest point

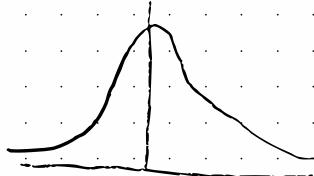
(most frequent of remaining vals gives classification)

→ Naive Bayes

↳ from conditional probability

$$P(A|B) \leftarrow \text{prob. } A \text{ given } B$$

i.e. given a person's height, what is prob of m/f



→ Probabilities (continuous)

↳ $f(x) = \text{pdf}(x)$ ↪ not a probability
but prob. DENSITY

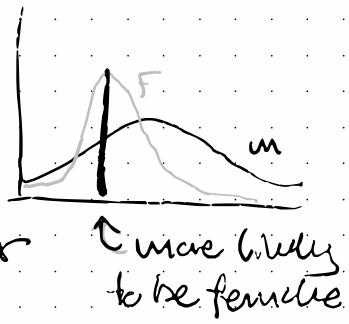
$$P(g) \sim \text{B}(0.5)$$

$$g = \begin{cases} 0: \text{male}, \\ 1: \text{female} \end{cases}$$

$$P(h|g) \sim N(70\text{in}, 16\text{in}^2) I_{g=0}$$

$$+ N(65\text{in}, 12.5\text{in}^2) I_{g=1}$$

↳ given our assumption of the probability of each gender & the data over heights we can predict someone's gender given their height



- ↪ Note that Normalisation is not necessary
- ↪ just choose argmax for classification
- ↪ comparing between the same $P(x)$?

Also at high dimensions prior & likelihoods start to reach 0.

↪ "Space" between points gets larger

$$k^* = \arg\max_k (\log P(x|y) + \log P(y))$$

↑
use log to make calculations robust to small probabilities
(esp w/ high dim.)

Probabilities (Discrete)

↪ example: spam classification

email → Clean up → feature extraction
wrapping

↪ n word appearance in corpus

↪ word is same or not

↪ series of Bernoulli trials \rightarrow Binomial

If $n = 100$ & $\theta = \text{prob of } \text{word} = 8\%$
^{specific}

& it appears 5 times (k)

$$f(5) = \binom{100}{5} 0.08^5 (1 - 0.08)^{100-5} = 99$$

i.e. we picked 100 words & get 5
of interest w/ $p(\theta) = 0.08$

(\hookrightarrow) there are $100c5$ ways we could
have done this

(\hookrightarrow) thus we can calculate which is more

likely θ { spam = 8%

{ !spam = 1%

(\hookrightarrow) get $p(\times | \theta)p(\theta)$ for all θ

& compare w/o norm w/ $p(x)$

(\hookrightarrow) increase the number of features

w/ multinomial θ

(lazy programmer kind of sucks at
teaching this)

Naive Bayes

(\hookrightarrow) sometimes features are not independent
independence is assumed!

Pseudocode

def fit(X, y)

gauss-dict = dict

priors = dict

for c in classes:

$X_c = X \text{ [corresponding } y = c]$

mu, var = mean & diagonal covar of X_c

dict of gauss(x) mu } var

priors [C] = len(X_c) / len(X)

def predict(X):

predictions = []

max post = -inf

best class = none

for x in X:

 for c in classes:

 mu, var = dict of gaussians[c]

 post = log-pdf(x, mu, var)
 + log(priors[c])

 if post > max post:

 max post = post

 best class = c

 pred.append(best-class)

return pred

Naïve Bayes

↳ not possible $p(X|C) = \prod_i p(X_i|C)$

→ Must find other routes

- ↳ full cov matrix (Σ)
 - ↳ Hidden markov model
 - ↳ Cust. Bayes Net
- can be arbitrarily complex

LDA

↳ Assumptions: - data for each feature is gauss distributed.

- (only consider bin class rel/0)

If we are checking

$$\frac{p(X|C=1)p(C=1)}{p(X)} > \frac{p(X|C=0)p(C=0)}{p(X)}$$

then

$$p(X|C=1)p(C=1) > p(X|C=0)p(C=0)$$

↳ plug in gauss (well)boxed for $p(x|c=)$
 $c \in \{0, 1\}$

Then note that $P(C=1) = 1 - P(C=0)$

→ & take the log derivative

↳ move to left hand side

(c) simplify & combine

(c) result will have 3 parts & depends on x

$$\left. \begin{array}{l} \textcircled{1} x^T \frac{1}{2} (\Sigma_0^{-1} - \Sigma_1^{-1}) x + \mu_1^T \Sigma_1^{-1} - \mu_0^T \Sigma_0^{-1} \\ \textcircled{2} \frac{1}{2} \mu_1^T \Sigma_1^{-1} \mu_1 + \frac{1}{2} \mu_0^T \Sigma_0^{-1} \mu_0 \\ \textcircled{3} \log \lambda - \log(1-\alpha) + u_1 - u_0 \end{array} \right\} 20$$

① Quadratic terms

② Concret ferrum

(3) const. terms

Quadratic discriminant analysis

$$x^T A x + w^T c + b > 0$$

↑ way simpler form

where

$$A = \frac{1}{2} (\Sigma_0^{-1} - \Sigma_1^{-1})$$

$$w^T x = (\mu_1^T \Sigma_1^{-1} \mu_1 - \mu_0^T \Sigma_0^{-1} \mu_0)$$

$$b = -\frac{1}{2} \mu_1^T \Sigma_1^{-1} \mu_1 + \frac{1}{2} \mu_0^T \Sigma_0^{-1} \mu_0 + k_1 + k_0 + \log \alpha - \log(1-\alpha)$$

if Σ is diagonal then this simplifies w/
quad term vanishing

(*) if assumptions about data is
true then bayes clf is optimal.

Generative vs Discriminative Models

- ↳ probability, predict $P(C|X)$
 - ↳ logit does this directly
also KNN
- ↳ Generative
 - ↳ calculates $P(X|C)$ & reverses prob w/ Bayes rule
for each C

Generative

- ↳ each var modeled explicitly w/ high level of explainability

Discriminative

- ↳ classifies better

→ Decision Trees

↳ makes predictions based on decisions



Challenges: how to make splits

how complex to make tree?

Splits are only possible orthogonal to axes

↳ previous could generate boundaries to arbitrary angle to axes

How to make splits?

↳ max entropy (split which maximizes information gain)

$$H(p) = - \sum_x p(x) \log(p(x)) \quad \begin{matrix} \text{relates} \\ \text{to variance} \end{matrix}$$

for binary clif

$$H(p) = -p \log(p) - (1-p) \log(1-p) \approx$$

$$\frac{\partial H(p)}{\partial p} = 0 \leftarrow \text{gives optimal value}$$

"entropy is a measure of how much info we get from making a split"

Algorithms for Maximising entropy gain

ID3:

for each C in Cols:

↳ find C that best splits data

↳ repeat for each split until all C is used

Others:

↳ use C multiple times

Pseudo code

def fit(X, y)

best Ig = 0

Best attr = None

for c in Cols:

cond = find_splc(X, y, c)

$x_{left} = X[X < c]$ # meets cond

$x_{right} = X[X \geq c]$ # !meets cond

l-gain = $H(y) - p(left)H(y_{left}) - p(right)H(y_{right})$

If best

Save to final vars

→ Perceptron

Setup

↳ Binary Classification
w/ targets -1 / +1

Predict

calc $w^T x + b$ } if ≥ 0 df = ?
 } if > 0 df = 1
 } if < 0 df = -1

train/fit

w = random vector

b = 0

eta = 0.1 # or 1 or .001 etc

for i in n-epochs:

get misclassifications (try, except break)

$x, y \in \text{get 1 mis. df}$

$w = w + \text{eta} \cdot y \cdot x$

→ Practical Machine Learning

- ↳ for each algorithm, what is best for hyperparameter tuning?
 - ↳ use CV w/ gridsearch (CV) (or variant of it)



↑ train test w/ grid
validate w/ left out

→ "Generalization"

→ Parameter space for gs?

- ↳ depends on compute power
- ↳ experience/trial & error

→ Feature extraction & selection

- ↳ vital for decent problem
- ↳ 1 feature (outputs) makes it super easy to do classification

→ requires knowledge of data
(domain knowledge)

- ↳ Not all features are useful & can lead to overfitting

various algorithms to choose

- ↳ select all & elim
- ↳ select 1 & build
 - ↳ try all combns
 - ↳ keep only best

etc

Other methods:

↳ PCA

↳ Select features w/ no corr

↳ Problem: only a lin. transform

→ Classic algos vs dL

↳ Some algos train & predict faster than others

↳ Depends on use case

↳ Decision boundaries

↳ Depends on variance/bias tradeoff
(also interpretable)

→ Multiclass Classifications

→ 1 vs rest

→ many 1 vs 1

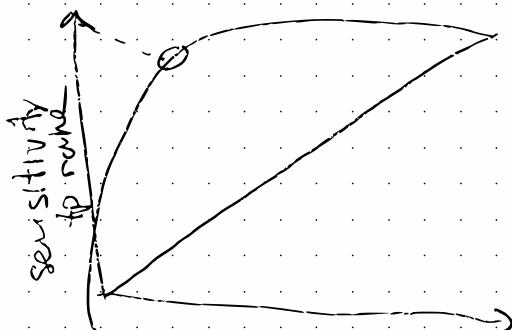
↳ cumbersome & more difficult

\rightarrow ROC/AUC

↳ if output is a probability, then
make classification based on threshold

AUC = performance

point closest to $(0, 1)$
is best threshold to
use.



$1 - \text{specificity}$
 $= \text{fp rate}$