

Artificial Intelligence/Machine Learning/Deep Learning: 'Bridging the Skills Gap'

Lesson 6: 'Naïve' Bayes Classifier

We will look at a bunch of classifiers in this course including:

Classification Algorithm	Coding (Yes/No)	Exercise	Comment
K-Nearest Neighbors (KNN)	✓		
Perceptron - Multi-Layer Perceptron	X		Perceptron is the simplest form of a Neural Network.
Naïve Bayes	✓		
Logistics Regression	✓		
Support Vector Machines and Kernel Functions	✓		
Decision Trees	✓		
Neural Networks	✓		

Naive Bayes models are **fast and simple** classification algorithms that are **often suitable for very high-dimensional datasets**. One application of naïve Bayes is text classification, aiming to assign documents (emails, tweets, posts, news) to one or many categories. One example is **email spam/not-spam** classification. The idea of naïve Bayes is that spam and not-spam emails have a different probability distributions.

For naïve Bayes the data doesn't need to be linearly separable. Naïve Bayes finds the hyperplane that best differentiates one distribution from another...it DOES NOT find hyperplane that separates + from -
A perceptron separates the data and will not converge when data is not linearly separable

Naive Bayes Classifier

(1)

↳ typically used for spam filtering < 'spam' 'ham' >

↳ but also used for other classification applications

- recognize handwritten digits on envelopes
- AUTOMATIC ESSAY grading
- medical diagnosis
- Fraud detection

Prior = Probability of A before B occurred

Bayes Rule

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

probability of A when B has occurred
→ POSTERIOR

> links probability before an EVENT B has occurred to probability before B has occurred!

GIVEN: A test point x_{test} (email) with features x_1, x_2, \dots, x_d ^{words of email}

goal → maximize $P(Y=y | x_1 = \text{word}_1, x_2 = \text{word}_2, \dots, x_d = \text{word}_d)$
'spam' 'ham' exact email never seen before

→ maximize $\frac{P(x_1 = \text{word}_1, \dots, x_d = \text{word}_d | Y=y) \cdot P(Y=y)}{P(x_1 = \text{word}_1, \dots, x_d = \text{word}_d)}$

Bayes Rule

given the email is spam
→ what would be the email?

Prior belief spam or ham before we take EVIDENCE INTO ACCOUNT

(2)

$$\rightarrow \text{MAXIMIZE } P(x_1 = \text{word}_1 | Y=y) \dots P(x_d = \text{word}_d | Y=y) \cdot P(Y=y)$$

Naïve Bayes

All words in the email received
ARE INDEPENDENT of each other
provided that you know the label y
(spam)

Naïve.

but computing is
cheap and results are
not bad!

ex: GIVEN that email is SPAM,
what is the probability that it
contains the word 'VIAGRA'?

Does Naïve Bayes overestimate or underestimate the
probability of spam?

assume $x_1 = \text{VIAGRA}$
 $x_2 = \text{KENIA}$ } email with any of these 2 words
will be likely SPAM!

BUT. unlikely for both words to be in same
email.

$$\textcircled{1} \quad P(x_1 = \text{VIAGRA}, x_2 = \text{KENIA} | Y = \text{spam})$$

Naïve Bayes $\rightarrow \textcircled{2} \quad P(x_1 = \text{VIAGRA} | Y = \text{spam}) \cdot P(x_2 = \text{KENIA} | Y = \text{spam})$

OVERESTIMATES
Probability

$$P(\textcircled{2}) \gg P(\textcircled{1})$$

Evidence

Several models could explain the data

\rightarrow Pick model that maximizes $P(\text{model} | \text{data})$

$$\sim P(\text{DATA} | \text{model}) \cdot P(\text{model})$$

Naive Bayes Summary → linear classifier

3

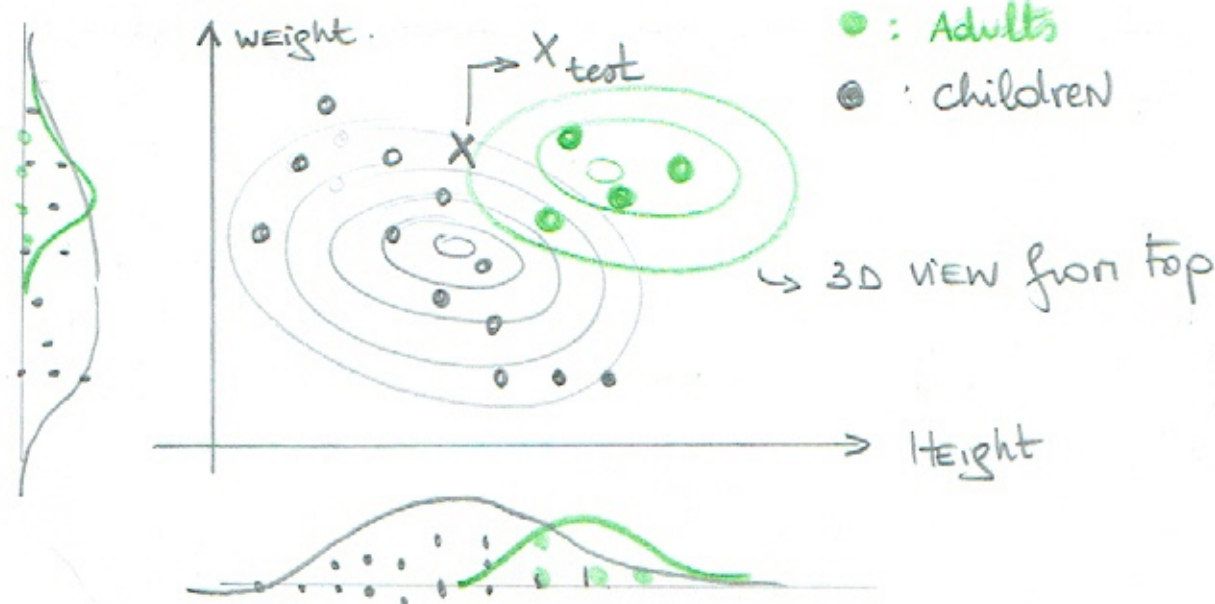
- ① CONTINUOUS Example : 2 classes $\begin{cases} \text{Adults : } a=y \\ \text{children : } c=y \end{cases}$
4 adults, 12 children

training data (h, w, y)
 ↓ height ↓ weight

- ① compute Prior : $P(a) = \frac{4}{16} = 25\%$
 $P(c) = 75\%$

- ② model selection : GAUSSIAN (μ, σ^2)
for distribution of weight, height for adults and for kids
ASSUME iid $\begin{cases} P(x|a) \\ P(x|c) \end{cases}$

DATA Scientist job



$$\text{Distribution} = \underbrace{\text{Distribution}_c}_{\text{GAUSS}} \cdot \underbrace{\text{Distribution}_a}_{\text{GAUSS}}$$

↑
GAUSS!

→ once we have the distributions we can forget about the data.

$$P(x|a) = P(h_x|a) P(w_x|a)$$

$$P(x|c) = P(h_x|c) P(w_x|c)$$

$$\rightarrow P(a|x) = \frac{P(x|a) P(a)}{P(x|a) P(a) + P(x|c) P(c)}$$

② Discrete Example: Multinomial distribution

GIVEN: MOVIE REVIEWS $\begin{cases} P \text{ (Positive)} \\ N \text{ (Negative)} \end{cases}$

① BAD, terrible, BORING N

② GOOD, would recommend P

③ would NOT recommend N

④ BORING, terrible N

(V) = vocabulary
= 7

PRIORS: $P(N) = 3/4 : 75\%$

$P(P) = 1/4 : 25\%$

Conditional Probabilities - likelihoods → SMOOTHING

$$P(w|c) = \frac{\text{count}(w,c) + 1}{\text{count}(c) + |V|}$$

↓ TOTAL words in class.

$$P(\text{BAD}|N) = \frac{1+1}{8+7} = \frac{2}{15}$$

$$P(\text{terrible}|N) = \frac{2+1}{8+7} = \frac{3}{15}$$

$$P(\text{BORING}|N) = \frac{2+1}{8+7} = \frac{3}{15}$$

$$P(\text{Good}|N) = \frac{0+1}{8+7} = \frac{1}{15}$$

$$P(\text{would}|N) = \frac{2}{8+7} = \frac{2}{15}$$

$$P(\text{recommend}|N) = \frac{2}{8+7} = \frac{2}{15}$$

$$P(\text{NOT}|N) = \frac{2}{8+7} = \frac{2}{15}$$

$$P(\text{BAD}|P) = \frac{0+1}{3+7} = \frac{1}{10}$$

$$P(\text{Terrible}|P) = \frac{0+1}{3+7} = \frac{1}{10}$$

$$P(\text{BORING}|P) = \frac{0+1}{3+7} = \frac{1}{10}$$

$$P(\text{Good}|P) = \frac{1+1}{3+7} = \frac{2}{10}$$

$$P(\text{would}|P) = \frac{1+1}{3+7} = \frac{2}{10}$$

$$P(\text{recommend}|P) = \frac{2}{3+7} = \frac{2}{10}$$

$$P(\text{Not}|P) = \frac{1}{10}$$

Now we look at a Test Point (Test Review)

(5)

(5) NOT GOOD, BORING "TR

CHOOSING A class:

$$P(N|TR) \sim \frac{3}{4} \cdot \frac{1}{15} \cdot \frac{2}{15} \cdot \frac{3}{15} = \frac{18}{13500} = 0.0013 = 0.13\%$$

\downarrow PRIOR. \downarrow $P(\text{word}|N)$
↳ Naive Bayes!

$$P(P|TR) \sim \frac{1}{4} \cdot \frac{1}{10} \cdot \frac{2}{10} \cdot \frac{1}{10} = \frac{4}{4000} = \frac{1}{1000} = 0.1\%$$

↓ Machine will classify
Test Review AS NEGATIVE because
0.13 > 0.1