

INTRODUCTION TO MACHINE LEARNING

David Kauchak
CS 51A – Fall 2025

1

Admin

Assignment 6

Assignment 7


Mentor hour change:

- Carmen: Wednesday 8-9pm (used to be 7-8pm)

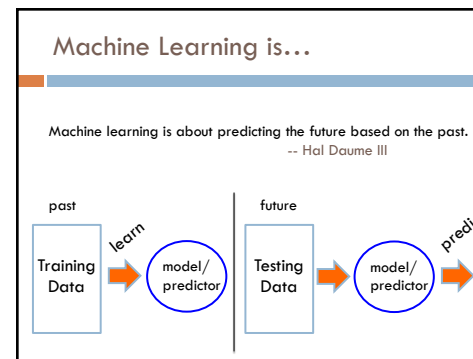
2

Machine Learning is...

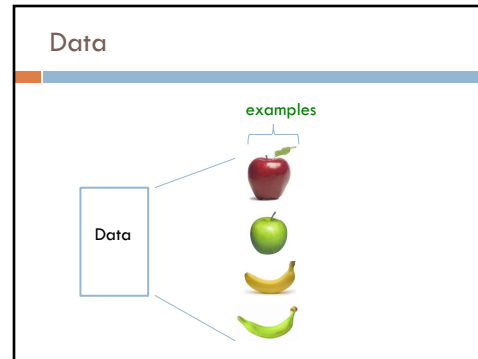
Machine learning is about predicting the future based on the past.
-- Hal Daume III



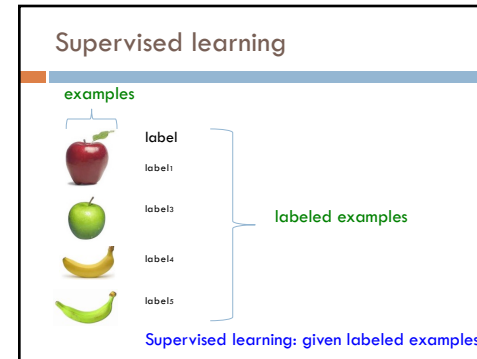
3



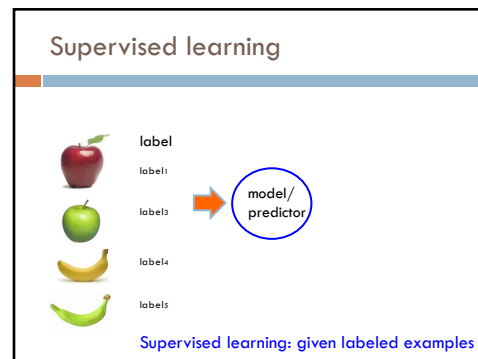
4



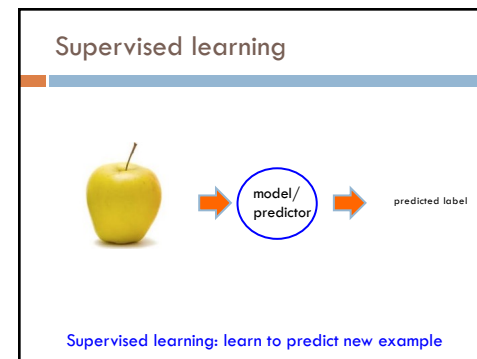
5



6







7



8

Supervised learning: classification

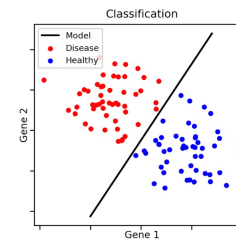
	label	
	apple	
	apple	
	banana	
	banana	

Classification: a finite set of labels

Supervised learning: given labeled examples

9

Classification Example



10

Classification Applications

Face recognition

Character recognition

Spam detection

Medical diagnosis: From symptoms to illnesses

Biometrics: Recognition/authentication using physical and/or behavioral characteristics: Face, iris, signature, etc

...

11

Supervised learning: regression

	label	
	-4.5	
	10.1	
	3.2	
	4.3	

Regression: label is real-valued

Supervised learning: given labeled examples

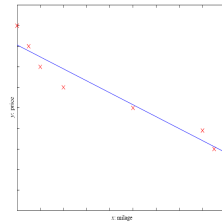
12

Regression Example

Price of a used car

x : car attributes
(e.g. mileage)

y : price



13

Regression Applications

Economics/Finance: predict the value of a stock

Epidemiology

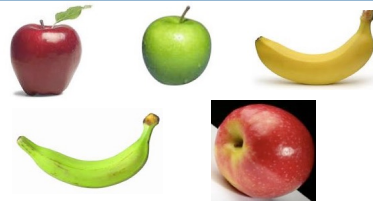
Car/plane navigation: angle of the steering wheel,
acceleration, ...

Temporal trends: weather over time

...

14

Unsupervised learning



Unsupervised learning: given data, i.e. examples, but no labels

18

Unsupervised learning applications

learn clusters/groups without any label

customer segmentation (i.e. grouping)

image compression

bioinformatics: learn motifs

...

19

Reinforcement learning

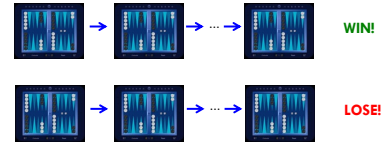
left, right, straight, left, left, left, straight	GOOD
left, straight, straight, left, right, straight, straight	BAD
left, right, straight, left, left, left, straight	18.5
left, straight, straight, left, right, straight, straight	-3

Given a **sequence** of examples/states and a **reward** after completing that sequence, learn to predict the action to take in for an individual example/state

20

Reinforcement learning example

Backgammon



Given sequences of moves and whether or not the player won at the end, learn to make good moves

21

Other learning variations

What data is available:

- Supervised, unsupervised, reinforcement learning
- semi-supervised, active learning, ...

How are we getting the data:

- online vs. offline learning

Type of model:

- generative vs. discriminative
- parametric vs. non-parametric

22

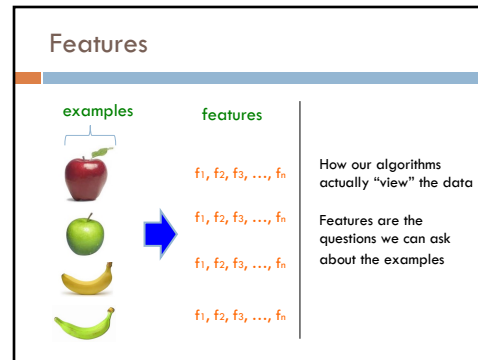
Representing examples

examples

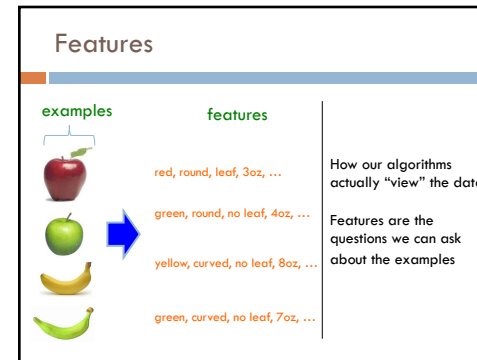


What is an example?
How is it represented?

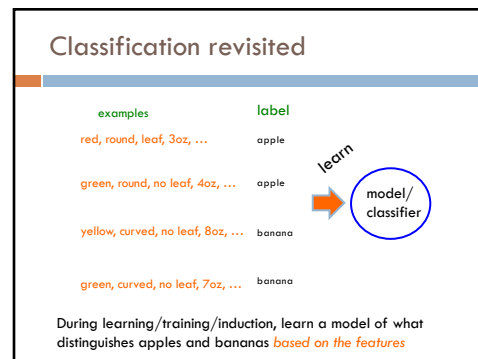
23



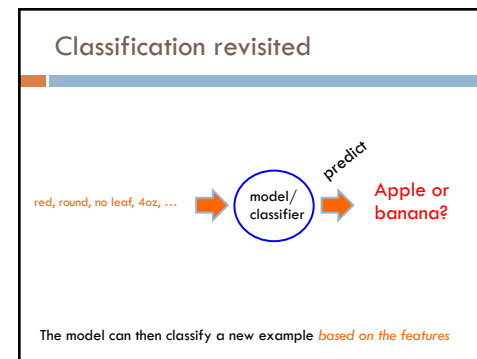
24



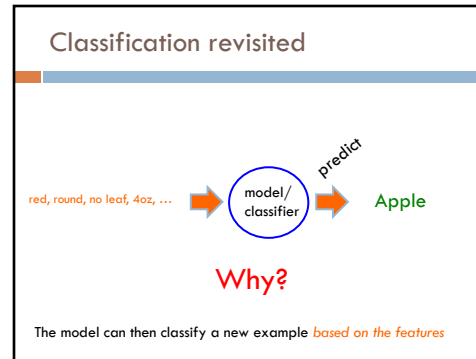
25



26



27



28

Classification revisited

Training data		Test set
examples	label	
red, round, leaf, 3oz, ...	apple	
green, round, no leaf, 4oz, ...	apple	red, round, no leaf, 4oz, ... ?
yellow, curved, no leaf, 4oz, ...	banana	
green, curved, no leaf, 5oz, ...	banana	

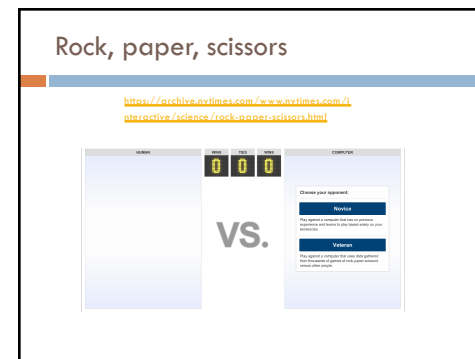
29

Classification revisited

Training data		Test set
examples	label	
red, round, leaf, 3oz, ...	apple	
green, round, no leaf, 4oz, ...	apple	red, round, no leaf, 4oz, ... ?
yellow, curved, no leaf, 4oz, ...	banana	
green, curved, no leaf, 5oz, ...	banana	

Learning is about **generalizing** from the training data

30



31

models

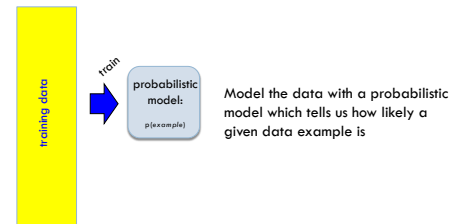


We have many, many different options for the model

They have different characteristics and perform differently (accuracy, speed, etc.)

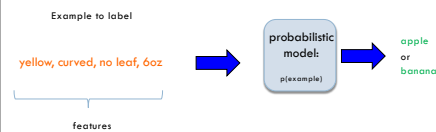
33

Probabilistic modeling



34

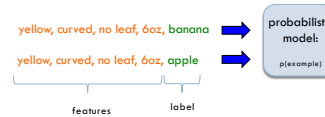
Probabilistic models



35

Probabilistic models

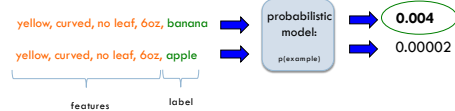
For each label, ask for the probability



36

Probabilistic models

Pick the label with the highest probability



37

Probability basics

A **probability distribution** gives the probabilities of all possible values of an event

For example, say we flip a coin three times. We can define the probability of the number of time the coin came up heads.

P(num heads)
P(3) = ?
P(2) = ?
P(1) = ?
P(0) = ?

38

Probability distributions

What are the possible outcomes of three flips (hint, there are eight of them)?

TTT
TTH
THT
THH
HTT
HTH
HHT
HHH

39

Probability distributions

Assuming the coin is fair, what are our probabilities?

$$\text{probability} = \frac{\text{number of times it happens}}{\text{total number of cases}}$$

TTT
TTH
THT
THH
HTT
HTH
HHT
HHH

P(num heads)
P(3) = ?
P(2) = ?
P(1) = ?
P(0) = ?

40

Probability distributions

Assuming the coin is fair, what are our probabilities?

$$\text{probability} = \frac{\text{number of times it happens}}{\text{total number of cases}}$$

TTT
TTH
THT
THH
HTT
HTH
HHT
HHH

P(num heads)
P(3) = ?
P(2) = ?
P(1) = ?
P(0) = ?

41

Probability distributions

Assuming the coin is fair, what are our probabilities?

$$\text{probability} = \frac{\text{number of times it happens}}{\text{total number of cases}}$$

TTT
TTH
THT
THH
HTT
HTH
HHT
HHH

P(num heads)
P(3) = 1/8
P(2) = ?
P(1) = ?
P(0) = ?

42

Probability distributions

Assuming the coin is fair, what are our probabilities?

$$\text{probability} = \frac{\text{number of times it happens}}{\text{total number of cases}}$$

TTT
TTH
THT
THH
HTT
HTH
HHT
HHH

P(num heads)
P(3) = 1/8
P(2) = ?
P(1) = ?
P(0) = ?

43

Probability distributions

Assuming the coin is fair, what are our probabilities?

$$\text{probability} = \frac{\text{number of times it happens}}{\text{total number of cases}}$$

TTT
TTH
THT
THH
HTT
HTH
HHT
HHH

P(num heads)
P(3) = 1/8
P(2) = 3/8
P(1) = ?
P(0) = ?

44

Probability distributions

Assuming the coin is fair, what are our probabilities?

$$\text{probability} = \frac{\text{number of times it happens}}{\text{total number of cases}}$$

TTT
TTH
THT
THH
HTT
HTH
HHT
HHH

P(num heads)
P(3) = 1/8
P(2) = 3/8
P(1) = 3/8
P(0) = 1/8

45

Probability distribution

A probability distribution assigns probability values to *all possible* values

Probabilities are between 0 and 1, inclusive

The sum of all probabilities in a distribution must be 1

P(num heads)
P(3) = 1/8
P(2) = 3/8
P(1) = 3/8
P(0) = 1/8

46

Probability distribution

A probability distribution assigns probability values to *all possible* values

Probabilities are between 0 and 1, inclusive

The sum of all probabilities in a distribution must be 1

P
P(3) = 1/2
P(2) = 1/2
P(1) = 1/2
P(0) = 1/2

P
P(3) = -1
P(2) = 2
P(1) = 0
P(0) = 0

47

Some example probability distributions

probability of heads
(distribution options: heads, tails)

probability of passing class
(distribution options: pass, fail)

probability of rain today
(distribution options: rain or no rain)

probability of getting an 'A'
(distribution options: A, B, C, D, F)

48

Conditional probability distributions

Sometimes we may know extra information about the world that may change our probability distribution

$P(X|Y)$ captures this (read "probability of X given Y")

- Given some information (Y) what does our probability distribution look like
- Note that this is still just a normal probability distribution

49

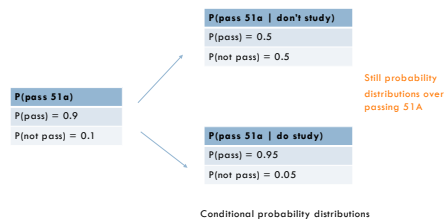
Conditional probability example

P(pass S1a)
 $P(\text{pass}) = 0.9$
 $P(\text{not pass}) = 0.1$

Unconditional probability distribution

50

Conditional probability example



51

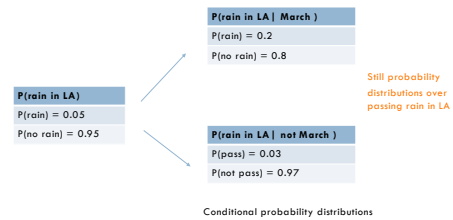
Conditional probability example

P(rain in LA)
 $P(\text{rain}) = 0.05$
 $P(\text{no rain}) = 0.95$

Unconditional probability distribution

52

Conditional probability example



53

Joint distribution

Probability over two events: $P(X,Y)$

Has probabilities for all possible combinations over the two events

51 Pass, EngPass	P(51 Pass, EngPass)
true, true	.88
true, false	.01
false, true	.04
false, false	.07

54

Joint distribution

Still a probability distribution

All questions/probabilities that we might want to ask about these two things can be calculated from the joint distribution

51 Pass, EngPass	P(51 Pass, EngPass)
true, true	.88
true, false	.01
false, true	.04
false, false	.07

What is $P(51 \text{ pass} = \text{true})$?

55

Joint distribution

51 Pass, EngPass	P(51 Pass, EngPass)
true, true	.88
true, false	.01
false, true	.04
false, false	.07

There are two ways that a person can pass 51: they can do it while passing or not passing English

$$P(51 \text{ Pass} = \text{true}) = P(\text{true, true}) + P(\text{true, false}) = 0.89$$

56

Relationship between distributions

$$P(X, Y) = P(Y) * P(X|Y)$$

joint distribution
unconditional distribution
conditional distribution

Can think of it as describing the two events happening in two steps:

The likelihood of X and Y happening:

1. How likely it is that Y happened?
2. Given that Y happened, how likely is it that X happened?

57

Relationship between distributions

$$P(51Pass, EngPass) = P(EngPass) * P(51Pass|EngPass)$$

The probability of passing CS51 and English is:

1. Probability of passing English *
2. Probability of passing CS51 **given** that you passed English

58

Relationship between distributions

$$P(51Pass, EngPass) = P(51Pass) * P(EngPass|51Pass)$$

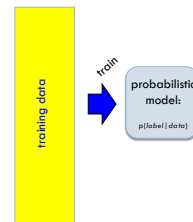
The probability of passing CS51 and English is:

1. Probability of passing **CS51** *
2. Probability of passing **English** **given** that you passed **CS51**

Can also view it with the other event happening first

59

Back to probabilistic modeling



Build a model of the conditional distribution:

$$P(\text{label} | \text{data})$$

How likely is a label given the data

60

Back to probabilistic models

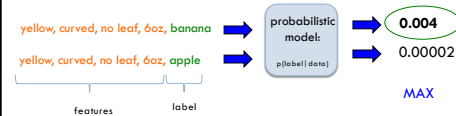
For each label, calculate the probability of the label given the data



61

Back to probabilistic models

Pick the label with the highest probability



62

Naïve Bayes model

Two parallel ways of breaking down the joint distribution

$$P(\text{data}, \text{label}) = P(\text{label}) * P(\text{data}|\text{label})$$

$$P(\text{data}, \text{label}) = P(\text{data}) * P(\text{label}|\text{data})$$

$$P(\text{label}) * P(\text{data}|\text{label}) = P(\text{data}) * P(\text{label}|\text{data})$$

What is $P(\text{label}|\text{data})$?

63

Naïve Bayes

$$P(\text{label}) * P(\text{data}|\text{label}) = P(\text{data}) * P(\text{label}|\text{data})$$



$$P(\text{label}|\text{data}) = \frac{P(\text{label}) * P(\text{data}|\text{label})}{P(\text{data})}$$

(This is called Bayes' rule!)

64

Naïve Bayes

$$P(\text{label}|\text{data}) = \frac{P(\text{label}) * P(\text{data}|\text{label})}{P(\text{data})}$$

probabilistic model: $p(\text{label}|\text{data})$

$$\frac{P(\text{positive}) * P(\text{data}|\text{positive})}{P(\text{data})} \quad \text{MAX}$$

$$\frac{P(\text{negative}) * P(\text{data}|\text{negative})}{P(\text{data})}$$

65

One observation

$$\frac{P(\text{positive}) * P(\text{data}|\text{positive})}{P(\text{data})} \quad \text{MAX}$$

$$\frac{P(\text{negative}) * P(\text{data}|\text{negative})}{P(\text{data})}$$

For picking the largest $P(\text{data})$ doesn't matter!

66

One observation

$$\frac{P(\text{positive}) * P(\text{data}|\text{positive})}{P(\text{data})} \quad \text{MAX}$$

$$\frac{P(\text{negative}) * P(\text{data}|\text{negative})}{P(\text{data})}$$

For picking the largest $P(\text{data})$ doesn't matter!

67

A simplifying assumption (for this class)

$$\frac{P(\text{positive}) * P(\text{data}|\text{positive})}{P(\text{data})} \quad \text{MAX}$$

$$\frac{P(\text{negative}) * P(\text{data}|\text{negative})}{P(\text{data})}$$

If we assume $P(\text{positive}) = P(\text{negative})$ then:

$$\frac{P(\text{data}|\text{positive})}{P(\text{data}|\text{negative})} \quad \text{MAX}$$

68

Naïve Bayes Assumption

$$P(\text{data}|\text{label}) = P(f_1, f_2, \dots, f_n|\text{label})$$

$$\approx P(f_1|\text{label}) * P(f_2|\text{label}) * \dots * P(f_n|\text{label})$$

This is generally not true!

However..., it makes our life easier.

This is why the model is called **Naïve** Bayes

69

Naïve Bayes

$$P(f_1|\text{positive}) * P(f_2|\text{positive}) * \dots * P(f_n|\text{positive})$$

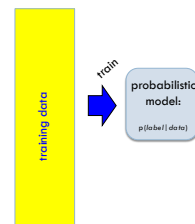
MAX

$$P(f_1|\text{negative}) * P(f_2|\text{negative}) * \dots * P(f_n|\text{negative})$$

Where do these come from?

70

Training Naïve Bayes



71

An aside: P(heads)

What is the $P(\text{heads})$ on a fair coin?

0.5

What if you didn't know that, but had a coin to experiment with?

Flip it a bunch of times and count how many times it comes up heads

$$P(\text{heads}) = \frac{\text{number of times heads came up}}{\text{total number of coin tosses}}$$

72

Try it out...

73

$P(\text{feature} | \text{label})$

$$P(\text{heads}) = \frac{\text{number of times heads came up}}{\text{total number of coin tosses}}$$

Can we do the same thing here? What is the probability of a feature given positive, i.e. the probability of a feature occurring in the positive label?

$$P(\text{feature} | \text{positive}) = ?$$

74

$P(\text{feature} | \text{label})$

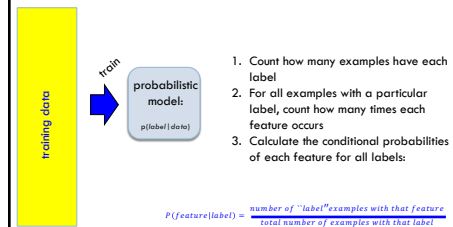
$$P(\text{heads}) = \frac{\text{number of times heads came up}}{\text{total number of coin tosses}}$$

Can we do the same thing here? What is the probability of a feature given positive, i.e. the probability of a feature occurring in the positive label?

$$P(\text{feature} | \text{positive}) = \frac{\text{number of positive examples with that feature}}{\text{total number of positive examples}}$$

75

Training Naïve Bayes



76

Classifying with Naïve Bayes

For each label, calculate the product of $p(\text{feature} | \text{label})$ for each label

yellow, curved, no leaf, 6oz \Rightarrow $P(\text{yellow} | \text{banana}) * \dots * P(6\text{oz} | \text{banana})$ **MAX**
 \Rightarrow $P(\text{yellow} | \text{apple}) * \dots * P(6\text{oz} | \text{apple})$

77

Naïve Bayes Text Classification

Positive

I loved it
 I loved that movie
 I hated that I loved it

Negative

I hated it
 I hated that movie
 I loved that I hated it

Given examples of text in different categories, learn to predict the category of new examples

Sentiment classification: given positive/negative examples of text (sentences), learn to predict whether new text is positive/negative

78

Text classification training

Positive

I loved it
 I loved that movie
 I hated that I loved it

Negative

I hated it
 I hated that movie
 I loved that I hated it

We'll assume words just occur once in any given sentence

79

Text classification training

Positive

I loved it
 I loved that movie
 I hated that loved it

Negative

I hated it
 I hated that movie
 I loved that hated it

We'll assume words just occur once in any given sentence

80

Training the model

Positive	Negative
I loved it	I hated it
I loved that movie	I hated that movie
I hated that loved it	I loved that hated it

For each word and each label, learn:

$$p(\text{word} \mid \text{label})$$

81

Training the model

Positive	Negative
I loved it	I hated it
I loved that movie	I hated that movie
I hated that loved it	I loved that hated it

$P(I \mid \text{positive}) = ?$

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

82

Training the model

Positive	Negative
I loved it	I hated it
I loved that movie	I hated that movie
I hated that loved it	I loved that hated it

$P(I \mid \text{positive}) = 3/3 = 1.0$

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

83

Training the model

Positive	Negative
I loved it	I hated it
I loved that movie	I hated that movie
I hated that loved it	I loved that hated it

$P(I \mid \text{positive}) = 1.0$
 $P(\text{loved} \mid \text{positive}) = ?$

$$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

84

Training the model

Positive	Negative
I loved it	I hated it
I loved that movie	I hated that movie
I hated that loved it	I loved that hated it

$P(I \mid \text{positive}) = 1.0$
 $P(\text{loved} \mid \text{positive}) = 3/3$

$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$

85

Training the model

Positive	Negative
I loved it	I hated it
I loved that movie	I hated that movie
I hated that loved it	I loved that hated it

$P(I \mid \text{positive}) = 1.0$
 $P(\text{loved} \mid \text{positive}) = 3/3$
 $P(\text{hated} \mid \text{positive}) = ?$

$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$

86

Training the model

Positive	Negative
I loved it	I hated it
I loved that movie	I hated that movie
I hated that loved it	I loved that hated it

$P(I \mid \text{positive}) = 1.0$
 $P(\text{loved} \mid \text{positive}) = 2/3$
 $P(\text{hated} \mid \text{positive}) = 1/3$
 $P(I \mid \text{negative}) = ?$

$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$

87

Training the model

Positive	Negative
I loved it	I hated it
I loved that movie	I hated that movie
I hated that loved it	I loved that hated it

$P(I \mid \text{positive}) = 1.0$
 $P(\text{loved} \mid \text{positive}) = 2/3$
 $P(\text{hated} \mid \text{positive}) = 1/3$
 $P(I \mid \text{negative}) = 1.0$

$P(\text{word} \mid \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$

88

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

$$\begin{aligned} P(I | \text{positive}) &= 1.0 \\ P(\text{loved} | \text{positive}) &= 2/3 \\ P(\text{hated} | \text{positive}) &= 1/3 \\ \dots \end{aligned}$$

$$P(\text{word} | \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

Negative

I hated it
I hated that movie
I loved that hated it

$$\begin{aligned} P(I | \text{negative}) &= 1.0 \\ P(\text{movie} | \text{negative}) &= ? \end{aligned}$$

89

Training the model

Positive

I loved it
I loved that movie
I hated that loved it

$$\begin{aligned} P(I | \text{positive}) &= 1.0 \\ P(\text{loved} | \text{positive}) &= 2/3 \\ P(\text{hated} | \text{positive}) &= 1/3 \\ \dots \end{aligned}$$

$$P(\text{word} | \text{label}) = \frac{\text{number of times word occurred in "label" examples}}{\text{total number of examples with that label}}$$

Negative

I hated it
I hated that movie
I loved that hated it

$$\begin{aligned} P(I | \text{negative}) &= 1.0 \\ P(\text{movie} | \text{negative}) &= 1/3 \\ \dots \end{aligned}$$

90

Classifying

$P(I \text{positive}) = 1.0$	$P(I \text{negative}) = 1.0$
$P(\text{loved} \text{positive}) = 1.0$	$p(\text{hated} \text{negative}) = 1.0$
$p(\text{it} \text{positive}) = 2/3$	$p(\text{that} \text{negative}) = 2/3$
$p(\text{that} \text{positive}) = 2/3$	$P(\text{movie} \text{negative}) = 1/3$
$p(\text{movie} \text{positive}) = 1/3$	$p(\text{it} \text{negative}) = 2/3$
$P(\text{hated} \text{positive}) = 1/3$	$p(\text{loved} \text{negative}) = 1/3$

Notice that each label has its own probability distribution

$P(\text{loved} \text{positive})$
$P(\text{loved} \text{positive}) = 2/3$
$P(\text{no loved} \text{positive}) = 1/3$

91

Trained model

$P(I \text{positive}) = 1.0$	$P(I \text{negative}) = 1.0$
$P(\text{loved} \text{positive}) = 2/3$	$p(\text{hated} \text{negative}) = 1.0$
$p(\text{it} \text{positive}) = 2/3$	$p(\text{that} \text{negative}) = 2/3$
$p(\text{that} \text{positive}) = 2/3$	$P(\text{movie} \text{negative}) = 1/3$
$p(\text{movie} \text{positive}) = 1/3$	$p(\text{it} \text{negative}) = 2/3$
$P(\text{hated} \text{positive}) = 1/3$	$p(\text{loved} \text{negative}) = 1/3$

How would we classify: "I hated movie"?

92

Trained model

$P(I \text{positive})$	$= 1.0$	$P(I \text{negative})$	$= 1.0$
$P(\text{loved} \text{positive})$	$= 2/3$	$p(\text{hated} \text{negative})$	$= 1.0$
$p(\text{it} \text{positive})$	$= 2/3$	$p(\text{that} \text{negative})$	$= 2/3$
$p(\text{that} \text{positive})$	$= 2/3$	$P(\text{movie} \text{negative})$	$= 1/3$
$p(\text{movie} \text{positive})$	$= 1/3$	$p(\text{it} \text{negative})$	$= 2/3$
$P(\text{hated} \text{positive})$	$= 1/3$	$p(\text{loved} \text{negative})$	$= 1/3$

$$P(I | \text{positive}) * P(\text{hated} | \text{positive}) * P(\text{movie} | \text{positive}) = 1.0 * 1/3 * 1/3 = 1/9$$

$$P(I | \text{negative}) * P(\text{hated} | \text{negative}) * P(\text{movie} | \text{negative}) = 1.0 * 1.0 * 1/3 = 1/3$$

93

Trained model

$P(I \text{positive})$	$= 1.0$	$P(I \text{negative})$	$= 1.0$
$P(\text{loved} \text{positive})$	$= 2/3$	$p(\text{hated} \text{negative})$	$= 1.0$
$p(\text{it} \text{positive})$	$= 2/3$	$p(\text{that} \text{negative})$	$= 2/3$
$p(\text{that} \text{positive})$	$= 2/3$	$P(\text{movie} \text{negative})$	$= 1/3$
$p(\text{movie} \text{positive})$	$= 1/3$	$p(\text{it} \text{negative})$	$= 2/3$
$P(\text{hated} \text{positive})$	$= 1/3$	$p(\text{loved} \text{negative})$	$= 1/3$

How would we classify: "I hated the movie"?

94

Trained model

$P(I \text{positive})$	$= 1.0$	$P(I \text{negative})$	$= 1.0$
$P(\text{loved} \text{positive})$	$= 2/3$	$p(\text{hated} \text{negative})$	$= 1.0$
$p(\text{it} \text{positive})$	$= 2/3$	$p(\text{that} \text{negative})$	$= 2/3$
$p(\text{that} \text{positive})$	$= 2/3$	$P(\text{movie} \text{negative})$	$= 1/3$
$p(\text{movie} \text{positive})$	$= 1/3$	$p(\text{it} \text{negative})$	$= 2/3$
$P(\text{hated} \text{positive})$	$= 1/3$	$p(\text{loved} \text{negative})$	$= 1/3$

$$P(I | \text{positive}) * P(\text{hated} | \text{positive}) * P(\text{the} | \text{positive}) * P(\text{movie} | \text{positive}) =$$

$$P(I | \text{negative}) * P(\text{hated} | \text{negative}) * P(\text{the} | \text{negative}) * P(\text{movie} | \text{negative}) =$$

95

Trained model

$P(I \text{positive})$	$= 1.0$	$P(I \text{negative})$	$= 1.0$
$P(\text{loved} \text{positive})$	$= 2/3$	$p(\text{hated} \text{negative})$	$= 1.0$
$p(\text{it} \text{positive})$	$= 2/3$	$p(\text{that} \text{negative})$	$= 2/3$
$p(\text{that} \text{positive})$	$= 2/3$	$P(\text{movie} \text{negative})$	$= 1/3$
$p(\text{movie} \text{positive})$	$= 1/3$	$p(\text{it} \text{negative})$	$= 2/3$
$P(\text{hated} \text{positive})$	$= 1/3$	$p(\text{loved} \text{negative})$	$= 1/3$

$$P(I | \text{positive}) * P(\text{hated} | \text{positive}) * P(\text{the} | \text{positive}) * P(\text{movie} | \text{positive}) =$$

$$P(I | \text{negative}) * P(\text{hated} | \text{negative}) * P(\text{the} | \text{negative}) * P(\text{movie} | \text{negative}) =$$

What are these?

96

Trained model

$P(I \text{positive})$	$= 1.0$	$P(I \text{negative})$	$= 1.0$
$P(\text{loved} \text{positive})$	$= 2/3$	$p(\text{hated} \text{negative})$	$= 1.0$
$p(\text{it} \text{positive})$	$= 2/3$	$p(\text{that} \text{negative})$	$= 2/3$
$p(\text{that} \text{positive})$	$= 2/3$	$P(\text{movie} \text{negative})$	$= 1/3$
$p(\text{movie} \text{positive})$	$= 1/3$	$p(\text{it} \text{negative})$	$= 2/3$
$P(\text{hated} \text{positive})$	$= 1/3$	$p(\text{loved} \text{negative})$	$= 1/3$

$$P(I | \text{positive}) * P(\text{hated} | \text{positive}) * P(\text{the} | \text{positive}) * P(\text{movie} | \text{positive}) =$$

$$P(I | \text{negative}) * P(\text{hated} | \text{negative}) * P(\text{the} | \text{negative}) * P(\text{movie} | \text{negative}) =$$

0. Is this a problem?

97

Trained model

$P(I \text{positive})$	$= 1.0$	$P(I \text{negative})$	$= 1.0$
$P(\text{loved} \text{positive})$	$= 2/3$	$p(\text{hated} \text{negative})$	$= 1.0$
$p(\text{it} \text{positive})$	$= 2/3$	$p(\text{that} \text{negative})$	$= 2/3$
$p(\text{that} \text{positive})$	$= 2/3$	$P(\text{movie} \text{negative})$	$= 1/3$
$p(\text{movie} \text{positive})$	$= 1/3$	$p(\text{it} \text{negative})$	$= 2/3$
$P(\text{hated} \text{positive})$	$= 1/3$	$p(\text{loved} \text{negative})$	$= 1/3$

$$P(I | \text{positive}) * P(\text{hated} | \text{positive}) * P(\text{the} | \text{positive}) * P(\text{movie} | \text{positive}) =$$

$$P(I | \text{negative}) * P(\text{hated} | \text{negative}) * P(\text{the} | \text{negative}) * P(\text{movie} | \text{negative}) =$$

Yes. They make the entire product go to 0 !

98

Trained model

$P(I \text{positive})$	$= 1.0$	$P(I \text{negative})$	$= 1.0$
$P(\text{loved} \text{positive})$	$= 2/3$	$p(\text{hated} \text{negative})$	$= 1.0$
$p(\text{it} \text{positive})$	$= 2/3$	$p(\text{that} \text{negative})$	$= 2/3$
$p(\text{that} \text{positive})$	$= 2/3$	$P(\text{movie} \text{negative})$	$= 1/3$
$p(\text{movie} \text{positive})$	$= 1/3$	$p(\text{it} \text{negative})$	$= 2/3$
$P(\text{hated} \text{positive})$	$= 1/3$	$p(\text{loved} \text{negative})$	$= 1/3$

$$P(I | \text{positive}) * P(\text{hated} | \text{positive}) * P(\text{the} | \text{positive}) * P(\text{movie} | \text{positive}) =$$

$$P(I | \text{negative}) * P(\text{hated} | \text{negative}) * P(\text{the} | \text{negative}) * P(\text{movie} | \text{negative}) =$$

Our solution: assume any unseen word has a small, fixed probability, e.g. in this example $1/10$

99

Trained model

$P(I \text{positive})$	$= 1.0$	$P(I \text{negative})$	$= 1.0$
$P(\text{loved} \text{positive})$	$= 2/3$	$p(\text{hated} \text{negative})$	$= 1.0$
$p(\text{it} \text{positive})$	$= 2/3$	$p(\text{that} \text{negative})$	$= 2/3$
$p(\text{that} \text{positive})$	$= 2/3$	$P(\text{movie} \text{negative})$	$= 1/3$
$p(\text{movie} \text{positive})$	$= 1/3$	$p(\text{it} \text{negative})$	$= 2/3$
$P(\text{hated} \text{positive})$	$= 1/3$	$p(\text{loved} \text{negative})$	$= 1/3$

$$P(I | \text{positive}) * P(\text{hated} | \text{positive}) * P(\text{the} | \text{positive}) * P(\text{movie} | \text{positive}) = 1/90$$

$$P(I | \text{negative}) * P(\text{hated} | \text{negative}) * P(\text{the} | \text{negative}) * P(\text{movie} | \text{negative}) = 1/30$$

Our solution: assume any unseen word has a small, fixed probability, e.g. in this example $1/10$

100

Full disclaimer

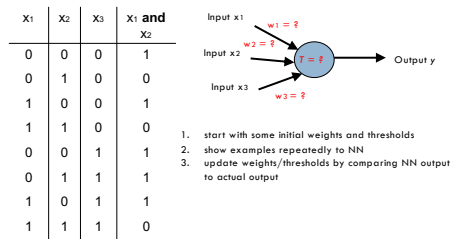
I've fudged a few things on the Naïve Bayes model for simplicity

Our approach is very close, but it takes a few liberties that aren't technically correct, but it will work just fine 😊

If you're curious, I'd be happy to talk to you offline

101

Training neural networks



102