# PROBABILISTIC MODELS

David Kauchak
CS158 – Fall 2025

1

## Admin

Assignment 6

Midterm

No class Thursday

No office hours Thursday

2

## Probabilistic Modeling

training data → *train* → probabilistic model: p(*features, label*)

Model the data with a probabilistic model

specifically, learn p(*features, label*)

p(*features, label*) tells us how likely these features and this example are

3

## Probabilistic models

Probabilistic models define a *probability distribution* over features and labels:

yellow, curved, no leaf, 6oz, banana → probabilistic model: p(*features, label*) → **0.004**

yellow, curved, no leaf, 6oz, apple → probabilistic model: p(*features, label*) → 0.00002

For each label, ask for the probability under the model
Pick the label with the highest probability

4

## Basic steps for probabilistic modeling

**Probabilistic models**

Step 1: pick a model

Step 2: figure out how to estimate the probabilities for the model

Step 3: (optional): deal with overfitting

Which model do we use, i.e. how do we calculate p(*feature, label*)?

How do train the model, i.e. how to we we estimate the probabilities for the model?

How do we deal with overfitting?

5

## Basic steps for probabilistic modeling

**Probabilistic models**

Step 1: pick a model

Step 2: figure out how to estimate the probabilities for the model

Step 3 (optional): deal with overfitting

Which model do we use, i.e. how do we calculate p(*feature, label*)?

How do train the model, i.e. how to we we estimate the probabilities for the model?

How do we deal with overfitting?

6

## Some math

$$p(features, label) = p(x_1, x_2, ..., x_m, y)$$

$$= p(y)p(x_1, x_2, ..., x_m \mid y)$$

$$= p(y)p(x_1 \mid y)p(x_2, ..., x_m \mid y, x_1)$$

$$= p(y)p(x_1 \mid y)p(x_2 \mid y, x_1)p(x_3, ..., x_m \mid y, x_1, x_2)$$

$$= p(y)\prod_{j=1}^{m} p(x_i \mid y, x_1, ..., x_{i-1})$$

7

## Step 1: pick a model

$$p(features, label) = p(y)\prod_{j=1}^{m} p(x_i \mid y, x_1, ..., x_{i-1})$$

So, far we have made NO assumptions about the data

$$p(x_m \mid y, x_1, x_2, ..., x_{m-1})$$

How many entries would the probability distribution table have if we tried to represent all possible values (e.g. for the wine data set)?

8

## Full distribution tables

| x₁ | x₂ | x₃ | ... | y | p( ) |
|---|---|---|---|---|---|
| 0 | 0 | 0 | ... | 0 | * |
| 0 | 0 | 0 | ... | 1 | * |
| 1 | 0 | 0 | ... | 0 | * |
| 1 | 0 | 0 | ... | 1 | * |
| 0 | 1 | 0 | ... | 0 | * |
| 0 | 1 | 0 | ... | 1 | * |
| | | | ... | | |

Wine problem:

- all possible combination of features
- ~7000 binary features
- Sample space size: $2^{7000}$ = ?

9

---

## $2^{7000}$

16216967556622020264666650B547837709519111243036374325623598208415152702316270235298708023787944600004651996019099530984538652557892546513204107022110253564658647431585227076599373340842842722420012281878260072931082617043194484266392077784125099999686016943600666001120981757929667871896255237700655294757256678055809293844627218640216108862600816097132874749204350870101862690842327501724605231129395523505905454421453447250950909650788947809468359293937411256947343861912152968484743444067412004170208875403718694217015502207353983812242992587435373561610415934359455766656170179090417259702533652666268202180849389281269970952857089069637557541434487608824836994199380241519751451012512704382908728091953847630285781185402409995889596419227760125536049911562403499947144160905730842429313962119953679373012994479560024833570738998392029910322346598038953069042980174009801732521069130797124201696339723021833530075897845195258485537108858195631737000743805167411189134617501484521767984296782842287373127422122022517597535994839257029877907706355334790244935435386660512591079567291431216297788784818552292819654176600980398997991681404749384215743518026038115106828640678973048382920234604277576530737765675475070271446622634876857096212610747627052030494889072089785936890470634285485316686656573271746606581856090664849508012761754614572161769555751992117507514067775104496728590822558547771447242334900764026321760892113552561241194538702680299044000183858505767193696897593661213568888838680023840932567380777501891470304962150996983853975207154939633923720287592041517294937079097785362510832009283960480723795488706954662168804465211249307629009199071774235503913517441532973747930089955830518884135334798464113680040999403737245600354288112326328218661131064550772899229969469156018580839820741704606832124388152026099584696588161375826382921029547343888832163627122302921229795384868355483535710603407789177417026363636202726955437517780741313455101810009476880940781122057380335371124632958916237089580476224595091825301636909236240671411644331656159828058372078343988856239089202844090255382937 6

**Any problems with this?**

10

---

## Full distribution tables

| x₁ | x₂ | x₃ | ... | y | p( ) |
|---|---|---|---|---|---|
| 0 | 0 | 0 | ... | 0 | * |
| 0 | 0 | 0 | ... | 1 | * |
| 1 | 0 | 0 | ... | 0 | * |
| 1 | 0 | 0 | ... | 1 | * |
| 0 | 1 | 0 | ... | 0 | * |
| 0 | 1 | 0 | ... | 1 | * |
| | | | ... | | |

- Storing a table of that size is impossible
- How are we supposed to learn/estimate each entry in the table?

11

---

## Step 1: pick a model

$$p(features, label) = p(y)\prod_{j=1}^{m} p(x_i \mid y, x_1, ..., x_{i-1})$$

So, far we have made **NO** assumptions about the data

Model selection involves making assumptions about the data

We did this before, e.g. assume the data is linearly separable

These assumptions allow us to represent the data *more compactly* and to estimate the parameters of the model

12

## An aside: independence

Two variables are independent if one has nothing to do with the other

For two independent variables, knowing the value of one does not change the probability distribution of the other variable (or the probability of any individual event)
- the result of the toss of a coin is independent of a roll of a die
- the price of tea in England is independent of the whether or not you pass ML

13

## independent or dependent?

Catching a cold and whether it's raining currently in NY

Miles per gallon and driving habits

Height and longevity of life

Ice cream sales and shark attacks

14

## Independent variables

How does independence affect our probability equations/properties?

If A and B are independent (written $A \perp\!\!\!\perp B$ )
- $P(A,B) = $ ?
- $P(A|B) = $ ?
- $P(B|A) = $ ?

15

## Independent variables

How does independence affect our probability equations/properties?

If A and B are independent (written $A \perp\!\!\!\perp B$ )
- $P(A,B) = P(A)P(B)$
- $P(A|B) = P(A)$
- $P(B|A) = P(B)$

How does independence help us?

16

4

## Independent variables

If A and B are independent
- P(A,B) = P(A)P(B)
- P(A|B) = P(A)
- P(B|A) = P(B)

Reduces the storage requirement for the distributions

Reduces the complexity of the distribution

Reduces the number of probabilities we need to estimate

17

## Conditional Independence

Dependent events can become independent given certain other events

Examples,
- height and length of life (or ice cream and shark attacks)
- "correlation" studies
  - size of your lawn and length of life

If A, B are conditionally independent given C (written $A \perp\!\!\!\perp B \,|C$)
- P(A,B|C) = P(A|C)P(B|C)
- P(A|B,C) = P(A|C)
- P(B|A,C) = P(B|C)
- but P(A,B) ≠ P(A)P(B)

18

## Naïve Bayes assumption

$$p(features, label) = p(y)\prod_{j=1}^{m} p(x_i \,|\, y, x_1, ..., x_{i-1})$$

$$p(x_i \,|\, y, x_1, x_2, ..., x_{i-1}) = p(x_i \,|\, y)$$

What does this assume?

19

## Naïve Bayes assumption

$$p(features, label) = p(y)\prod_{j=1}^{m} p(x_i \,|\, y, x_1, ..., x_{i-1})$$

$$p(x_i \,|\, y, x_1, x_2, ..., x_{i-1}) = p(x_i \,|\, y)$$

Assumes feature i is independent of the the other features **given the label** (i.e. is conditionally independent given the label)

For the wine problem?

20

## Naïve Bayes assumption

$$p(x_i \mid y, x_1, x_2, ..., x_{i-1}) = p(x_i \mid y)$$

Assumes feature *i* is independent of the the other features *given the label*

Assumes the probability of a word occurring in a review is independent of the other words *given the label*

For example, the probability of "pinot" occurring is independent of whether or not "wine" occurs given that the review is about "chardonnay"

Is this assumption true?

21

## Naïve Bayes assumption

$$p(x_i \mid y, x_1, x_2, ..., x_{i-1}) = p(x_i \mid y)$$

For most applications, this is not true!

For example, the fact that "pinot" occurs will probably make it *more likely* that "noir" occurs (or other compound phrases like "San Francisco")

However, this is often a reasonable approximation:

$$p(x_i \mid y, x_1, x_2, ..., x_{i-1}) \approx p(x_i \mid y)$$

22

## Naïve Bayes model

$$p(features, label) = p(y) \prod_{j=1}^{m} p(x_i \mid y, x_1, ..., x_{i-1})$$

$$= p(y) \prod_{j=1}^{m} p(x_i \mid y) \quad \text{naïve bayes assumption}$$

$p(x_i \mid y)$ is the probability of a particular feature value given the label

How do we model this?
- for binary features
- for discrete features, i.e. counts
- for real valued features

23

## p(x|y)

Binary features:

$$p(x_i \mid y) = \begin{cases} \theta_i & if \ x_i = 1 \\ 1 - \theta_i & otherwise \end{cases} \quad \text{biased coin toss!}$$

Other features:

Could use a lookup table for each value, but doesn't generalize well

Better, model as a distribution:
- gaussian (i.e. normal) distribution
- poisson distribution
- multinomial distribution (more on this later)
- …

24

## Basic steps for probabilistic modeling

Step 1: pick a model

Step 2: figure out how to estimate the probabilities for the model

Step 3 (optional): deal with overfitting

**Probabilistic models**

Which model do we use, i.e. how do we calculate p(*feature, label*)?

How do train the model, i.e. how to we we estimate the probabilities for the model?

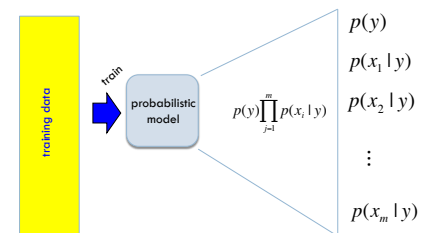How do we deal with overfitting?

25

## Obtaining probabilities



We've talked a lot about probabilities, but not where they come from

- How do we calculate $p(x_i | y)$ from training data?
- What is the probability of surviving the titanic?
- What is the probability that a review is about Pinot Noir?
- What is the probability that a particular review is about Pinot Noir?

26

## Obtaining probabilities



$$p(y) \prod_{j=1}^{m} p(x_i | y)$$

$p(y)$

$p(x_1 | y)$

$p(x_2 | y)$

$\vdots$

$p(x_m | y)$

27

## Estimating probabilities

**What is the probability of a pinot noir review?**

**We don't know!**

**We can *estimate* it based on data, though:**

$$\frac{\text{number of reviews labeled pinot noir}}{\text{total number of reviews}}$$

This is called the maximum likelihood estimation.  Why?

28

## Maximum Likelihood Estimation (MLE)

Maximum likelihood estimation picks the values for the model parameters that *maximize the likelihood* of the training data

You flip a coin 100 times. 60 times you get heads and 40 times you get tails.

What is the MLE estimate for heads?

p(head) = 0.60        why?

29

## Likelihood

The *likelihood* of a data set is the probability that a particular model (i.e. a model and estimated probabilities) assigns to the data

$$likelihood(data) = \prod_{i=1}^{n} p_\theta(x_i)$$

for each example                    how probable is it under the model

the model parameters (e.g. probability of heads)

30

## Likelihood

You flip a coin 100 times. 60 times you get heads and 40 times you get tails.

What is the likelihood of this data with Θ=p(head) = 0.6 ?

$$likelihood(data) = \prod_{i=1}^{n} p_\theta(x_i)$$

for each example                    how probable is it under the model

the model parameters (e.g. probability of heads)

31

## Likelihood

You flip a coin 100 times. 60 times you get heads and 40 times you get tails.

What is the likelihood of this data with Θ=p(head) = 0.6 ?

$$likelihood(data) = \prod_{i=1}^{n} p_\theta(x_i)$$

$0.60^{60} * 0.40^{40} = 5.908465121038621e\text{-}30$

60 heads with p(head) = 0.6            40 tails with p(tail) = 0.4

32

## MLE example

Can we do any better?

$$likelihood(data) = \prod_i p(x_i)$$

$0.60^{60} * 0.40^{40} = 5.908465121038621e\text{-}30$

60 heads with p(head) = 0.6      40 tails with p(tail) = 0.4

What about p(head) = 0.5?

33

## MLE example

Can we do any better?

$$likelihood(data) = \prod_i p(x_i)$$

$0.60^{60} * 0.40^{40} = 5.908465121038621e\text{-}30$

60 heads with p(head) = 0.6      40 tails with p(tail) = 0.4

$0.50^{60} * 0.50^{40} = 7.888609052210118e\text{-}31$

60 heads with p(head) = 0.5      40 tails with p(tail) = 0.5

34

## MLE example

Can we do any better?

$$likelihood(data) = \prod_i p(x_i)$$

$0.60^{60} * 0.40^{40} = 5.908465121038621e\text{-}30$

60 heads with p(head) = 0.6      40 tails with p(tail) = 0.4

What about p(head) = 0.7?

35

## MLE example

Can we do any better?

$$likelihood(data) = \prod_i p(x_i)$$
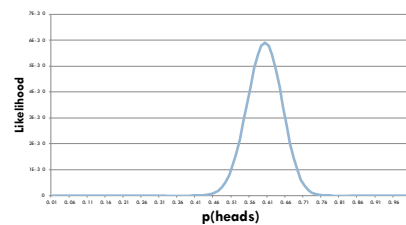
$0.60^{60} * 0.40^{40} = 5.908465121038621e\text{-}30$

60 heads with p(head) = 0.6      40 tails with p(tail) = 0.4

$0.70^{60} * 0.30^{40} = 6.176359828759916e\text{-}31$

60 heads with p(head) = 0.7      40 tails with p(tail) = 0.3

36

## MLE Example



## Maximum Likelihood Estimation (MLE)

The *maximum likelihood* estimate for a model parameter is the one that maximizes the likelihood of the training data

$$MLE = \arg\max_{\theta} \prod_{i=1}^{n} p_{\theta}(x_i)$$

Often easier to work with log-likelihood:

$$MLE = \operatorname{argmax}_{\theta} \log(\prod_{i=1}^{n} p_{\theta}(x_i))$$

Why is this ok?

$$= \operatorname{argmax}_{\theta} \sum_{i=1}^{n} \log(p(x_i))$$

## Calculating MLE

The *maximum likelihood* estimate for a model parameter is the one that maximize the likelihood of the training data

$$MLE = \operatorname{argmax}_{\theta} \sum_{i=1}^{n} \log(p(x_i))$$

Given some training data, how do we calculate the MLE?

You flip a coin 100 times.  60 times you get heads and 40 times you get tails.

## Calculating MLE

You flip a coin 100 times.  60 times you get heads and 40 times you get tails.

$$\log-likelihood = \sum_{i=1}^{n} \log(p(x_i))$$

$$= 60\log(p(heads)) + 40\log(p(tails))$$

$$= 60\log(\theta) + 40\log(1-\theta)$$

$$MLE = \arg\max_{\theta} 60\log(\theta) + 40\log(1-\theta)$$
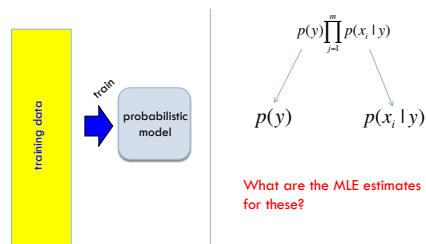
How do we find the max?

37

38

39

40

## Calculating MLE

You flip a coin 100 times.  60 times you get heads and 40 times you get tails.

$$\frac{d}{d\theta} 60\log(\theta) + 40\log(1-\theta) = 0$$

$$\frac{60}{\theta} - \frac{40}{1-\theta} = 0$$

$$\frac{40}{1-\theta} = \frac{60}{\theta}$$

$$40\theta = 60 - 60\theta$$

$$100\theta = 60$$

$$\theta = \frac{60}{100} \qquad \text{Yay!}$$

41

## Calculating MLE

You flip a coin n times.  **a** times you get heads and **b** times you get tails.

$$\frac{d}{d\theta} a\log(\theta) + b\log(1-\theta) = 0$$

$$...$$

$$\theta = \frac{a}{a+b}$$

42

## MLE estimation for NB

training data → train → probabilistic model

$$p(y)\prod_{j=1}^{m} p(x_i \mid y)$$

$$p(y) \qquad p(x_i \mid y)$$

What are the MLE estimates for these?

43

## Maximum likelihood estimates

$$p(y) = \frac{count(y)}{n}$$

number of examples with label y

total number of examples

$$p(x_i \mid y) = \frac{count(x_i, y)}{count(y)}$$

number of examples with label y with feature $x_i = 1$

number of examples with label

What does training a NB model then involve?
How difficult is this to calculate?

44

11

## Naïve Bayes classification

yellow, curved, no leaf, 6oz, banana → NB Model $p(features, label)$ → 0.004
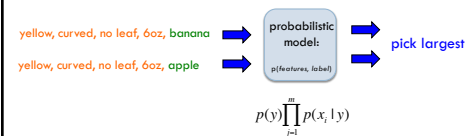
$$p(y)\prod_{j=1}^{m}p(x_i \mid y)$$

Given an unlabeled example: yellow, curved, no leaf, 6oz predict the label

How do we use a probabilistic model for classification/prediction?

45

## Probabilistic models

yellow, curved, no leaf, 6oz, banana → probabilistic model: $p(features, label)$ →

yellow, curved, no leaf, 6oz, apple →

pick largest

$$p(y)\prod_{j=1}^{m}p(x_i \mid y)$$

label $= \arg\max_{y \in labels} p(y)\prod_{j=1}^{m}p(x_i \mid y)$

46

## Generative Story

To classify with a model, we're given an example and we obtain the probability

We can also ask how a given model would *generate* a document

This is the "generative story" for a model

Looking at the generative story can help understand the model

We also can use generative stories to help develop a model

47

## NB generative story

$$p(y)\prod_{j=1}^{m}p(x_i \mid y)$$

What is the generative story for the NB model?

48

12

## NB generative story

$$p(y)\prod_{j=1}^{m} p(x_i \mid y)$$

1. Pick a label according to p(y)
   - roll a biased, num_labels-sided die
2. For each feature:
   - Flip a biased coin:
     - if heads, include the feature
     - if tails, don't include the feature

   What about for modeling wine reviews?

49

## NB decision boundary

$$\text{label} = \text{argmax}_{y \in labels}\, p(y)\prod_{j=1}^{m} p(x_i \mid y)$$

What does the decision boundary for
NB look like if the features are binary?

50

## Some math

$$label = \log(\text{argmax}_{y \in labels}\, p(y)\prod_{j=1}^{m} p(x_i \mid y))$$

$$= \text{argmax}_{y \in labels}\, \log(p(y)) + \sum_{i=1}^{m} \log(p(x_i \mid y))$$

$$= \text{argmax}_{y \in labels}\, \log(p(y)) + \sum_{i=1}^{m} x_i \log(p(x_i \mid y)) + \overline{x}_i \log(1 - p(x_i \mid y))$$

$$p(x_i \mid y) = \begin{cases} \theta_i & if\ x_i = 1 \\ 1 - \theta_i & otherwise \end{cases}$$

51

## Some more math

$$labels = \text{argmax}_{y \in labels}\, \log(p(y)) + \sum_{i=1}^{m} x_i \log(p(x_i \mid y)) + \overline{x}_i \log(1 - p(x_i \mid y))$$

$$= \text{argmax}_{y \in labels}\, \log(p(y)) + \sum_{i=1}^{m} x_i \log(p(x_i \mid y)) + (1 - x_i) \log(1 - p(x_i \mid y))$$

(because $x_i$ are binary)

$$= \text{argmax}_{y \in labels}\, \log(p(y)) + \sum_{i=1}^{m} x_i \log(p(x_i \mid y)) - x_i \log(1 - p(x_i \mid y)) + \log(1 - p(x_i \mid y))$$

$$= \text{argmax}_{y \in labels}\, \log(p(y)) + \sum_{i=1}^{m} x_i \log\left(\frac{p(x_i \mid y)}{1 - p(x_i \mid y)}\right) + \log(1 - p(x_i \mid y))$$

52

13

## And…

$$labels = \text{argmax}_{y \in labels} \log(p(y)) + \sum_{i=1}^{m} x_i \log\left(\frac{p(x_i \mid y)}{1 - p(x_i \mid y)}\right) + \log(1 - p(x_i \mid y)$$

$$= \text{argmax}_{y \in labels} \log(p(y)) + \sum_{i=1}^{m} \log(1 - p(x_i \mid y)) + \sum_{i=1}^{m} x_i \log\left(\frac{p(x_i \mid y)}{1 - p(x_i \mid y)}\right)$$

**What does this look like?**

53

## And…

$$labels = \text{argmax}_{y \in labels} \log(p(y)) + \sum_{i=1}^{m} x_i \log\left(\frac{p(x_i \mid y)}{1 - p(x_i \mid y)}\right) + \log(1 - p(x_i \mid y)$$

$$= \text{argmax}_{y \in labels} \underbrace{\log(p(y)) + \sum_{i=1}^{m} \log(1 - p(x_i \mid y))} + \sum_{i=1}^{m} x_i \log\left(\frac{p(x_i \mid y)}{1 - p(x_i \mid y)}\right)$$

b        +        $x_i * w_i$

w x + b          **What are the weights?**

**Linear model !!!**

54

## NB as a linear model

$$w_i = \log\left(\frac{p(x_i \mid y)}{1 - p(x_i \mid y)}\right)$$

How likely this feature is to be 1 given the label

How likely this feature is to be 0 given the label

**When is this big/small?**

55

## NB as a linear model

$$w_i = \log\left(\frac{p(x_i \mid y)}{1 - p(x_i \mid y)}\right)$$

How likely this feature is to be 1 given the label

How likely this feature is to be 0 given the label

- low magnitude weights indicate there isn't much difference
- larger weights (positive or negative) indicate feature is important

56

## Maximum likelihood estimation

Intuitive

Sets the probabilities so as to maximize the probability of the training data

**Problems?**
- Overfitting!
- Amount of data
  - particularly problematic for rare events
- Is our training data representative

57

## Basic steps for probabilistic modeling

**Probabilistic models**

Step 1: pick a model

Which model do we use, i.e. how do we calculate p(*feature, label*)?

Step 2: figure out how to estimate the probabilities for the model

How do train the model, i.e. how to we we estimate the probabilities for the model?

Step 3 (optional): deal with overfitting

How do we deal with overfitting?

58

## Coin experiment

59



LAW OF LARGE NUMBERS IN AVERAGE OF DIE ROLLS
AVERAGE CONVERGES TO EXPECTED VALUE OF 3.5

60

## Back to parasitic gaps

Say the actual probability is 1/100,000

We don't know this, though, so we're estimating it from a small data set of 10K sentences

What is the probability that we have a parasitic gap sentence in our sample?

61

## Back to parasitic gaps

p(not_parasitic) = 0.99999

$p(not\_parasitic)^{10000} \approx 0.905$ is the probability of us NOT finding one

Then probability of us finding one is ~10%
- □ 90% of the time we won't find one and won't know anything (or assume p(parasitic) = 0)
- □ 10% of the time we would find one and incorrectly assume the probability is 1/10,000 (10 times too large!)

Solutions?

62

## Priors

Coin1 data: 3 Heads and 1 Tail
Coin2 data: 30 Heads and 10 tails
Coin3 data: 2 Tails
Coin4 data: 497 Heads and 503 tails

If someone asked you what the probability of heads was for each of these coins, what would you say?

63