# MACHINE LEARNING BASICS

David Kauchak
CS159 Fall 2024

1

---

## Admin

Assignment 6

2

---

## Quiz #3

45 minutes (plus 20 to scan/upload)

Open book and notes

Text Similarity (10/3) through Machine Translation (10/24)

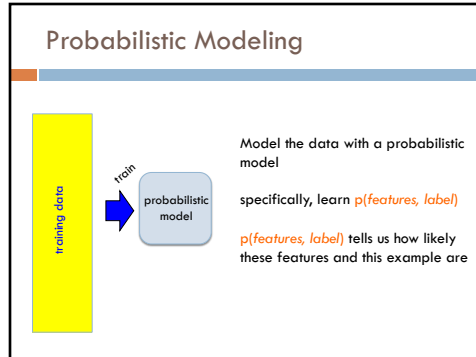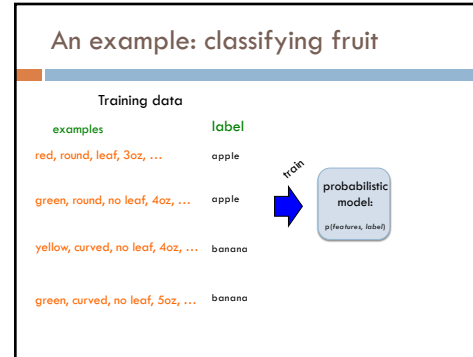Class will be 30 min discussion about final projects

3

---

## Machine Learning is...

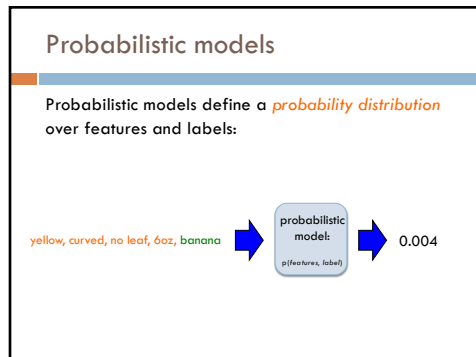Machine learning is about predicting the future based on the past.
-- Hal Daume III

4

---

## Probabilistic Modeling

training data → train → probabilistic model

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

5

## An example: classifying fruit

Training data

| examples | label |
|---|---|
| red, round, leaf, 3oz, … | apple |
| green, round, no leaf, 4oz, … | apple |
| yellow, curved, no leaf, 4oz, … | banana |
| green, curved, no leaf, 5oz, … | banana |

train → probabilistic model: p(features, label)

6

## 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

7

## 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 → 0.00002

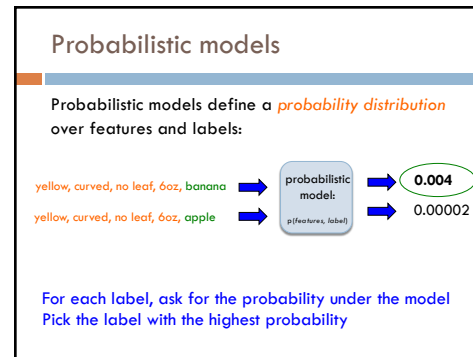For each label, ask for the probability under the model
Pick the label with the highest probability

8

## Probabilistic models: big questions

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

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

3. How do we deal with overfitting (i.e. smoothing)?

9

## 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?

10

## Some math

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

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

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

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

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

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

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

How many entries would the probability distribution table have if we tried to represent all possible values and we had 7000 binary features?

12

## 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 | * |
| | | | … | | |

All possible combination of features!

Table size: $2^{7000} =$ **?**

13

---

## $2^{7000}$

16216967556622020264666650854783770951911124303637432562359820841515270231627023529870802378794460004651996019099530984338652557892546513204107022110253564658647431585227076599373340842842722420012281878260072931082617043194484266392077784125099999686016943600666001120981757929667878196255237700655294757256678055809293844627218640216108862600816097132874749204350287401101862690842327501724605231129395523505905454421453477250950909650788947809468359293937411256947343861912152968484743444067412004170208875403718694217015502207353983812240992587435373536161041593435945576665617017909041725970253365266626820218084938928126997095285708906963755754143448760882483699419938024151975145101251270438290872809195384763028578118540240999588959641922776012553604911562403499947144160905730842429313962119953679373012944795600248333570738998392029910322346598038593306904298017400980173252106913079712420169633972302183530075897845195258485537108858195631737000743805167411189136175014845216798429678284228737313124221220225175975335994839257029877907706355334790244935435386660512591079567291431216297788784818552292819654176600980398997991681404749384215743515802603811510682864067897304838292203460427757653073776567547507027144662263487685709621261074762705203049488907208978593689040706342854853166866565732717466065818560906648495080127617546145721617695557519921175075140677751044967285908225585477714472423349007640263217608921135525612411945387026802990440018385850576719369689759366121356888838680023840932567380777501891470304962150996983853975207154939633923720287592041517294937079097785362510832009283960480723795488706954662168804465211249307629009199071774235503913517441532973747930089955830518884135334798464113680004999940373724560035428811232632828618661131064550772899229969469156018580839820741704606832124388152026099584696588161375826382921029547343888883216362712230292122979538486835548353571060340778917747170263636560202726955437517780743134551018100094680940781122057380335371124632958916237089580476224595091825301636909236240671411164433165615982805837207834398885623908920284409025538293376

**Any problems with this?**

14

---

## 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?

15

---

## 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've done this before, n-gram language model, parsing, etc.

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

16

## Naïve Bayes assumption

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

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

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

17

## Naïve Bayes model

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

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

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

How do we model this?
- for binary features (e.g., "banana" occurs in the text)
- for discrete features (e.g., "banana" occurs $x_i$ times)
- for real valued features (e.g, the text contains $x_i$ proportion of verbs)

18

## p(x | y)

Binary features (aka, Bernoulli Naïve Bayes) :

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

19

## 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?

20

## Obtaining probabilities

training data → (train) probabilistic model

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

$p(y)$

$p(x_1\mid y)$

$p(x_2\mid y)$

$\vdots$

$p(x_m\mid y)$

(m = number of features)

21

## MLE estimation for Bernoulli NB

training data → (train) probabilistic model

$$p(y)\prod_{i=1}^{m}p(x_j\mid y)$$

$p(y)$          $p(x_j\mid y)$

What are the MLE estimates for these?

22

## Maximum likelihood estimates

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

number of examples with label

total number of examples

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

number of examples with the label with feature

number of examples with label

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

23

## Text classification

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

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

Unigram features:
$w_i$, whether or not word $w_i$ occurs in the text

What are these counts for text classification with unigram features?

24

## Text classification

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

number of texts with label
—————————————
total number of texts

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

number of texts with the label with word $w_j$
—————————————
number of texts with label

25

## Naïve Bayes classification

yellow, curved, no leaf, 6oz, banana → NB Model $p(features, label)$ → 0.004

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

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

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

26

## NB classification

yellow, curved, no leaf, 6oz, banana → probabilistic model: $p(features, label?)$ $p(y=1)\prod_{j=1}^{m} p(x_j \mid y=1)$ →

yellow, curved, no leaf, 6oz, apple → $p(y=2)\prod_{j=1}^{m} p(x_j \mid y=2)$ →

pick largest

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

27

## NB classification

yellow, curved, no leaf, 6oz, banana → probabilistic model: $p(features, label?)$ $p(y=1)\prod_{j=1}^{m} p(x_j \mid y=1)$ →

yellow, curved, no leaf, 6oz, apple → $p(y=2)\prod_{j=1}^{m} p(x_j \mid y=2)$ →

pick largest

Notice that each label has its own separate set of parameters, i.e. $p(x_i \mid y)$

28

## Bernoulli NB for text classification

probabilistic
model:

p(features, label?)

$(1, 1, 1, 0, 0, 1, 0, 0, ...)$

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

$p(y=2)\prod_{j=1}^{m} p(w_j \mid y=2)$

pick largest

How good is this model for *text* classification?

29

## Bernoulli NB for text classification

$(1, 1, 1, 0, 0, 1, 0, 0, ...)$

$p(y=1)\prod_{j=1}^{m} p(w_j \mid y=1)$       $p(y=2)\prod_{j=1}^{m} p(w_j \mid y=2)$

pick largest

For text classification, what is this computation?
Does it make sense?

30

## Bernoulli NB for text classification

$(1, 1, 1, 0, 0, 1, 0, 0, ...)$

$p(y=1)\prod_{j=1}^{m} p(w_j \mid y=1)$       $p(y=2)\prod_{j=1}^{m} p(w_j \mid y=2)$

pick largest

Each word that occurs, contributes $p(w_i \mid y)$
Each word that does NOT occur, contributes $1-p(w_i \mid y)$

31

## 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* an example

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

32

8

## Bernoulli NB generative story

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

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 does this mean for text classification, assuming unigram features?

33

## Bernoulli NB generative story

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

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

34

## Bernoulli NB

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

Pros/cons?

35

## Bernoulli NB

Pros
  - Easy to implement
  - Fast!
  - Can be done on large data sets

Cons
  - Naïve Bayes assumption is generally not true
  - Performance isn't as good as other models
  - For text classification (and other sparse feature domains) the $p(x_i=0 \mid y)$ can be problematic

36

## Another generative story

Randomly draw words from a "bag of words" until document length is reached

37

## Draw words from a fixed distribution

Selected: $w_1$

38

## Draw words from a fixed distribution

Selected: $w_1$

Put a copy of $w_1$ back

sampling with replacement

39

## Draw words from a fixed distribution

Selected: $w_1$  $w_1$

40

## Draw words from a fixed distribution

Selected: W1 W1

Put a copy of w1 back

sampling with replacement

41

## Draw words from a fixed distribution

Selected: W1 W1 W2

42

## Draw words from a fixed distribution

Selected: W1 W1 W2

Put a copy of w2 back

sampling with replacement

43

## Draw words from a fixed distribution

Selected: W1 W1 W2 …

44

## Draw words from a fixed distribution

Is this a NB model, i.e. does it assume each individual word occurrence is independent?



45

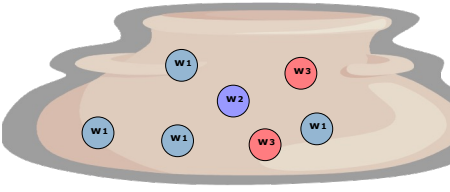## Draw words from a fixed distribution

Yes! Doesn't matter what words were drawn previously, still the same probability of getting any particular word



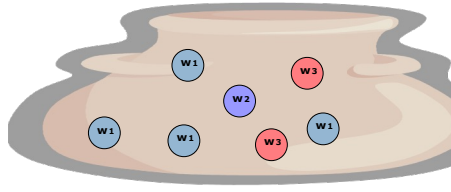46

## Draw words from a fixed distribution

Does this model handle multiple word occurrences?



47

## Draw words from a fixed distribution

Selected: W1 W1 W2 …



48

## NB generative story

### Bernoulli NB

1. Pick a label according to p(y)
   - roll a biased, num_labels-sided die

2. For each word in your vocabulary:
   - Flip a biased coin:
     - if heads, include the word in the text
     - if tails, don't include the word

### Multinomial NB

1. Pick a label according to p(y)
   - roll a biased, num_labels-sided die

2. Keep drawing words from *p(words|y)* until text length has been reached.

49

## Probabilities

### Bernoulli NB

1. Pick a label according to p(y)
   - roll a biased, num_labels-sided die

2. For each word in your vocabulary:
   - Flip a biased coin:
     - if heads, include the word in the text
     - if tails, don't include the word

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

(1, 1, 1, 0, 0, 1, 0, 0, ...)

### Multinomial NB

1. Pick a label according to p(y)
   - roll a biased, num_labels-sided die

2. Keep drawing words from *p(words|y)* until document length has been reached

**?**

(4, 1, 2, 0, 0, 7, 0, 0, ...)

50

## A digression: rolling dice

What's the probability of getting a 3 for a single roll of this dice?

1/6

51

## A digression: rolling dice

What is the probability distribution over possible single rolls?

| 1/6 | 1/6 | 1/6 | 1/6 | 1/6 | 1/6 |
|-----|-----|-----|-----|-----|-----|
| 1   | 2   | 3   | 4   | 5   | 6   |

52

## A digression: rolling dice

What if I told you 1 was twice as likely as the others?

| 2/7 | 1/7 | 1/7 | 1/7 | 1/7 | 1/7 |
|-----|-----|-----|-----|-----|-----|
| 1 | 2 | 3 | 4 | 5 | 6 |

53

## A digression: rolling dice

What if I rolled 400 times and got the following number?

1: 100
2: 50
3: 50
4: 100
5: 50
6: 50

| 1/4 | 1/8 | 1/8 | 1/4 | 1/8 | 1/8 |
|-----|-----|-----|-----|-----|-----|
| 1 | 2 | 3 | 4 | 5 | 6 |

54

## A digression: rolling dice

1. What is the probability of rolling a 1 and a 5 (in any order)?

2. Two 1s and a 5 (in any order)?

3. Five 1s and two 5s (in any order)?

| 1/4 | 1/8 | 1/8 | 1/4 | 1/8 | 1/8 |
|-----|-----|-----|-----|-----|-----|
| 1 | 2 | 3 | 4 | 5 | 6 |

55

## A digression: rolling dice

1. What is the probability of rolling a 1 and a 5 (in any order)?

$(1/4 * 1/8) * 2 = 1/16$

prob. of those two rolls          number of ways that can happen
(1,5 and 5,1)

1. Two 1s and a 5 (in any order)?
$((1/4)^2 * 1/8) * 3 = 3/128$

2. Five 1s and two 5s (in any order)?
$((1/4)^5 * (1/8)^3) * 21 = 21/524,288 = 0.00004$          General formula?

| 1/4 | 1/8 | 1/8 | 1/4 | 1/8 | 1/8 |
|-----|-----|-----|-----|-----|-----|
| 1 | 2 | 3 | 4 | 5 | 6 |

56

14

## Multinomial distribution

Multinomial distribution: independent draws over *m* possible categories

If we have frequency counts x1, x2, ..., xm over each of the categories, the probability is:

$$p(x_1, x_2, ..., x_m \mid \theta_1, \theta_2, ..., \theta_m) = \frac{n!}{\prod_{j=1}^{m} x_j!} \prod_{j=1}^{m} \theta_j^{x_j}$$

number of different ways to get those counts     probability of particular counts

| $\theta_1$ | $\theta_2$ | $\theta_3$ | $\theta_4$ | $\theta_5$ | $\theta_6$ | ... |
|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | |

57

## Multinomial distribution

$$p(x_1, x_2, ..., x_m \mid \theta_1, \theta_2, ..., \theta_m) = \frac{n!}{\prod_{j=1}^{m} x_j!} \prod_{j=1}^{m} \theta_j^{x_j}$$

What are $\theta_i$?

Are there any constraints on the values that they can take?

| $\theta_1$ | $\theta_2$ | $\theta_3$ | $\theta_4$ | $\theta_5$ | $\theta_6$ | ... |
|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | |

58

## Multinomial distribution

$$p(x_1, x_2, ..., x_m \mid \theta_1, \theta_2, ..., \theta_m) = \frac{n!}{\prod_{j=1}^{m} x_j!} \prod_{j=1}^{m} \theta_j^{x_j}$$

$\theta_j$: probability of rolling "j"

$$\theta_j \geq 0$$

$$\sum_{j=1}^{m} \theta_j = 1$$

| $\theta_1$ | $\theta_2$ | $\theta_3$ | $\theta_4$ | $\theta_5$ | $\theta_6$ | ... |
|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | |

59

## Back to words…

Why the digression?

$$p(x_1, x_2, ..., x_m \mid \theta_1, \theta_2, ..., \theta_m) = \frac{n!}{\prod_{j=1}^{m} x_j!} \prod_{j=1}^{m} \theta_j^{x_j}$$

Drawing words from a bag is the same as rolling a die!

number of sides = number of words in the vocabulary

60

15

## Back to words…

**Why the digression?**

$$p(x_1, x_2, ..., x_m | \theta_1, \theta_2, ..., \theta_m) = \frac{n!}{\prod_{j=1}^{m} x_j!} \prod_{j=1}^{m} \theta_j^{x_j}$$

$$p(features, label) = p(y) \frac{n!}{\prod_{j=1}^{m} x_j!} \prod_{j=1}^{m} (\theta_y)_j^{x_j}$$

$\theta_j$ for class y

61

## Basic steps for probabilistic modeling

Model each class as a multinomial:

$$p(features, label) = p(y) \frac{n!}{\prod_{j=1}^{m} x_j!} \prod_{j=1}^{m} (\theta_y)_j^{x_j}$$

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

How do we train the model, i.e. estimate $\theta_j$ for each class?

62

## A digression: rolling dice

What if I rolled 400 times and got the following number?

1: 100
2: 50
3: 50
4: 100
5: 50
6: 50

| 1/4 | 1/8 | 1/8 | 1/4 | 1/8 | 1/8 |
|-----|-----|-----|-----|-----|-----|
| 1   | 2   | 3   | 4   | 5   | 6   |

63

## Training a multinomial

label1

label2

| 1/4 | 1/8 | 1/8 | 1/4 | 1/8 | 1/8 |
|-----|-----|-----|-----|-----|-----|
| 1   | 2   | 3   | 4   | 5   | 6   |

64

16

## Training a multinomial

label1

For each label, y:

w1: 100 times
w2: 50 times
w3: 10 times
w4: …

$$\theta_j = \frac{count(w_j, y)}{\sum_{k=1}^{m} count(w_k, y)}$$

$$= \frac{\text{number of times word } w_j \text{ occurs in label y docs}}{\text{total number of words in label y docs}}$$

| 1/4 | 1/8 | 1/8 | 1/4 | 1/8 | 1/8 |
|-----|-----|-----|-----|-----|-----|
| 1   | 2   | 3   | 4   | 5   | 6   |

65

## Classifying with a multinomial

(10, 2, 6, 0, 0, 1, 0, 0, …)

$$p(y=1)\frac{n!}{\prod_{j=1}^{m} x_j!}\prod_{j=1}^{m}(\theta_1)_j^{x_j} \qquad p(y=2)\frac{n!}{\prod_{j=1}^{m} x_j!}\prod_{j=1}^{m}(\theta_2)_j^{x_j}$$

**Any way I can make this simpler?**

pick largest

66

## Classifying with a multinomial

(10, 2, 6, 0, 0, 1, 0, 0, …)

$$p(y=1)\prod_{j=1}^{m}(\theta_1)_j^{x_j} \qquad p(y=2)\prod_{j=1}^{m}(\theta_2)_j^{x_j}$$

$$\frac{n!}{\prod_{j=1}^{m} x_m!} \quad \text{Is a constant!}$$

pick largest

67

## Multinomial finalized

Training:
- Calculate p(label)
- For each label, calculate θs

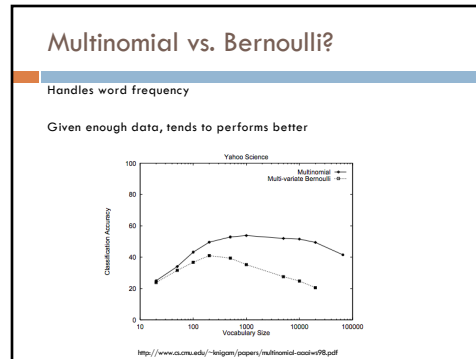$$\theta_j = \frac{count(w_j, y)}{\sum_{k=1}^{m} count(w_k, y)}$$

Classification:
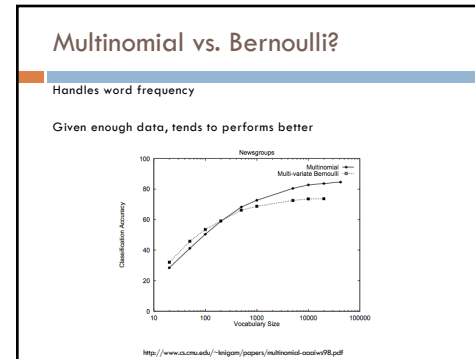- Get word counts
- For each label you had in training, calculate:

$$p(y)\prod_{j=1}^{m}\theta_j^{x_j}$$

and pick the largest

68

17

## Multinomial vs. Bernoulli?

Handles word frequency

Given enough data, tends to performs better



http://www.cs.cmu.edu/~knigam/papers/multinomial-aaaiws98.pdf

69

## Multinomial vs. Bernoulli?

Handles word frequency

Given enough data, tends to performs better



http://www.cs.cmu.edu/~knigam/papers/multinomial-aaaiws98.pdf

70

## Multinomial vs. Bernoulli?

Handles word frequency

Given enough data, tends to performs better



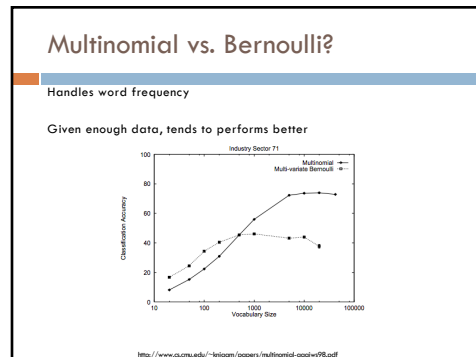http://www.cs.cmu.edu/~knigam/papers/multinomial-aaaiws98.pdf

71

## 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

72

## 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?

73

## Unseen events

training data → positive  banana: 2

negative  banana: 0

$$\theta_j = \frac{count(w_j, y)}{\sum_{k=1}^{m} count(w_k, y)}$$

What will θbanana be for the negative class?

74

## Unseen events

training data → positive  banana: 2

negative  banana: 0

$$\theta_j = \frac{count(w_j, y)}{\sum_{k=1}^{m} count(w_k, y)}$$

What will θbanana be for the negative class?

0!  Is this a problem?

75

## Unseen events

training data → positive  banana: 2

negative  banana: 0

p("I ate a bad banana", negative) = ?

76

19

## Unseen events



training data

positive — banana: 2

negative — banana: 0

p("I ate a bad banana", negative) = 0
p(".... banana ...", negative) = 0

Solution?

77

## Add lambda smoothing

training data

$$\theta_j = \frac{count(w_j, y)}{\sum_{k=1}^{m} count(w_k, y)}$$

$$\theta_j = \frac{count(w_j, y) + \lambda}{\lambda m + \sum_{k=1}^{m} count(w_k, y)}$$

for each label, pretend like we've seen each feature/word occur in λ additional examples

78

## Different than...

training data

positive — banana: 0

negative — banana: 0

How is this problem different?

79

## Different than...

training data

positive — banana: 0

negative — banana: 0

p("I ate a bad banana", positive)      p("I ate a bad", positive)

p("I ate a bad banana", negative)      p("I ate a bad", negative)

Out of vocabulary. Many ways to solve... for our implementation, we'll just ignore them.

80

20

## 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?

81

## Training revisited

From a probability standpoint, MLE training is selecting the Θ that maximizes:

$$p(\theta \mid data)$$

i.e.

$$\operatorname{argmax}_\theta p(\theta \mid data)$$

We pick the most likely model parameters given the data

82

## Estimating revisited

We can incorporate a prior belief in what the probabilities might be!

To do this, we need to break down our probability

$$p(\theta \mid data) = ?$$

(Hint: Bayes rule)

83

## Estimating revisited

What are each of these probabilities?

$$p(\theta \mid data) = \frac{p(data \mid \theta) p(\theta)}{p(data)}$$

84

21

## Priors

likelihood of the data under the model

probability of different parameters, call the **prior**

$$p(\theta \mid data) = \frac{p(data \mid \theta) p(\theta)}{p(data)}$$

probability of seeing the data (regardless of model)

85

## Priors

$$\theta = \operatorname{argmax}_\theta \frac{p(data \mid \theta) p(\theta)}{p(data)}$$

Does p(data) matter for the argmax?

86

## Priors

likelihood of the data under the model

probability of different parameters, call the **prior**

$$\theta = \operatorname{argmax}_\theta p(data \mid \theta) p(\theta)$$

What does MLE assume for a prior on the model parameters?

87

## Priors

likelihood of the data under the model

probability of different parameters, call the **prior**

$$\theta = \operatorname{argmax}_\theta p(data \mid \theta) p(\theta)$$

- Assumes a uniform prior, i.e. all Θ are equally likely!
- Relies solely on the likelihood

88

22

## A better approach

$$\theta = \text{argmax}_\theta \, p(data \,|\, \theta) p(\theta)$$

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

We can use any distribution we'd like
This allows us to impart addition bias
into the model

89

## Another view on the prior

Remember, the max is the same if we take the log:

$$\theta = \text{argmax}_\theta \, \log(p(data \,|\, \theta)) + \log(p(\theta))$$

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

We can use any distribution we'd like
This allows us to impart addition bias
into the model

90

## What about smoothing?

training data

$$\theta_j = \frac{count(w_j, y)}{\sum_{k=1}^{m} count(w_k, y)}$$

$$\theta_j = \frac{count(w_j, y) + \lambda}{\lambda m + \sum_{k=1}^{m} count(w_k, y)}$$

for each label, pretend like
we've seen each feature/word
occur in λ additional examples

Sometimes this is also called smoothing
because it is seen as smoothing or interpolating
between the MLE and some other distribution

91

## Prior for NB

$$\theta = \text{argmax}_\theta \, \log(p(data \,|\, \theta)) + \log(p(\theta))$$

Uniform prior

Dirichlet prior

λ= 0    increasing

$$p(w_j \,|\, y) = \theta_j = \frac{count(w_j, y)}{\sum_{k=1}^{m} count(w_k, y)}$$

$$\theta_j = \frac{count(w_j, y) + \lambda}{\sum_{k=1}^{m} \left( count(w_k, y) + \lambda \right)} = \frac{count(w_j, y) + \lambda}{\lambda m + \sum_{k=1}^{m} count(w_k, y)}$$

92

23