

PROBABILISTIC MODELS

David Kauchak
CS158 – Fall 2013

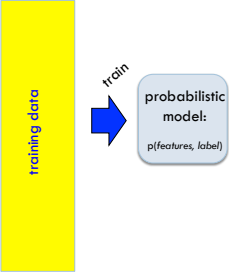
Admin

Assignment 6

Office hours today:

- ▣ 2:30-3:15
- ▣ 4-4:45

Probabilistic Modeling



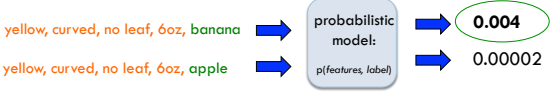
Model the data with a probabilistic model

specifically, learn $p(\text{features}, \text{label})$

$p(\text{features}, \text{label})$ tells us how likely these features and this example are

Probabilistic models

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



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

Basic steps for probabilistic modeling

	Probabilistic models
Step 1: pick a model	Which model do we use, i.e. how do we calculate $p(\text{feature}, \text{label})$?
Step 2: figure out how to estimate the probabilities for the model	How do we train the model, i.e. how do we estimate the probabilities for the model?
Step 3: (optional): deal with overfitting	How do we deal with overfitting?

Basic steps for probabilistic modeling

	Probabilistic models
Step 1: pick a model	Which model do we use, i.e. how do we calculate $p(\text{feature}, \text{label})$?
Step 2: figure out how to estimate the probabilities for the model	How do we train the model, i.e. how do we estimate the probabilities for the model?
Step 3 (optional): deal with overfitting	How do we deal with overfitting?

Some maths

$$\begin{aligned}
 p(\text{features}, \text{label}) &= p(x_1, x_2, \dots, x_m, y) \\
 &= p(y)p(x_1, x_2, \dots, x_m | y) \\
 &= p(y)p(x_1 | y)p(x_2, \dots, x_m | y, x_1) \\
 &= p(y)p(x_1 | y)p(x_2 | y, x_1)p(x_3, \dots, x_m | y, x_1, x_2) \\
 &= p(y) \prod_{j=1}^m p(x_j | y, x_1, \dots, x_{j-1})
 \end{aligned}$$

Step 1: pick a model

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

So, far we have made NO assumptions about the data

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

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

Full distribution tables

x_1	x_2	x_3	...	y	$p(\cdot)$
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: $2^{7000} = ?$

2^{7000}

```

162189758662020264466650854783770951911134303637432562359820841515270231627023529870800237879
4460046519960190995309845386525578925465132041070221102535646586474315852207076599373340842842
722420012281878260072931082617043194484266392077841250999968601694360066600112098175792966787
819625523770655294757256678055809293844627218440216108862600816097132874749204352087401101862
69084232750172460523112939523305905442145547725090096507889478094683592939574112569473438
619121529684847434440647120417402088754037186942170155022073539838122429925874355753611041593
435945576665617017909041725970253365266626820218084938928126997095285708906637557541434487608
824836994199380241519751451012512704382908728091953847630285781185402409995895964192277601255
3604911562403499947144160905700842429213962119953679373012944795600248332570738998392029910222
3465980389530690429801740098017325210491307971242016963397230218353007589784519525848553710885
819563173700074380516741118913461750148452176798429678284228737312742212202251759535994839257
02987907706355347902449354353866405125910795672914312162977887848185522928196541764009803989
9799168140474938621574515802603811510682864067897304838222034604277576550737657654750702714
4662263487685709621261074762705203049488907208978593689047064328548531668665657327174660658185
6090664849508012761754614572161769555751992117507514067775104496728590822558547771447242334900
76402632176089211355256124119453870268029904400183838505719369689759366121356888838680023840
93256738077501891470304962150969698385397500715493963392320287592041517949020709977853625108
320092836048072379548870695466216880446521124930762900919907177423550391351174415329737479300
89955830518884135334798464113680004994037372456003542881123263282186611310645507728992299646
91560185808398207417046068321243881530260995846965881613758263829210295473438888321636712302
91220793384848354835371040340778917741702636365620222695637517780741313455101810009468804
078112205738033371124632958916237089580476224595091825301636909236240671411644331656159828058
3720783439888562390892028440902553829376
    
```

Any problems with this?

Full distribution tables

x_1	x_2	x_3	...	y	$p(\cdot)$
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?

Step 1: pick a model

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

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

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

independent or dependent?

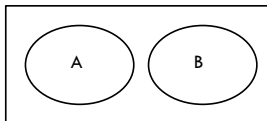
Catching a cold and having a cat-allergy

Miles per gallon and driving habits

Height and longevity of life

Independent variables

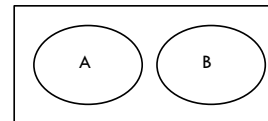
How does independence affect our probability equations/properties?



If A and B are independent (written ...)

- ▣ $P(A,B) = ?$
- ▣ $P(A|B) = ?$
- ▣ $P(B|A) = ?$

Independent variables



If A and B are independent (written ...)

- ▣ $P(A,B) = P(A)P(B)$
- ▣ $P(A|B) = P(A)$
- ▣ $P(B|A) = P(B)$

How does independence help us?

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

Conditional Independence

Dependent events can become independent given certain other events

Examples,

- height and length of life
- "correlation" studies
 - size of your lawn and length of life

If A, B are **conditionally independent** given 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 $P(A,B) \neq P(A)P(B)$

Naïve Bayes assumption

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

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

What does this assume?

Naïve Bayes assumption

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

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

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

For the wine problem?

Naïve Bayes assumption

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

Assumes feature i is independent of 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?

Naïve Bayes assumption

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

For most applications, this is not true!

For example, the fact that “pinot” occurs will probably make it *more likely* that “noir” occurs (or take a compound phrase like “San Francisco”)

However, this is often a reasonable approximation:

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

Naïve Bayes model

$$\begin{aligned} p(\text{features}, \text{label}) &= p(y) \prod_{j=1}^m p(x_j | y, x_1, \dots, x_{j-1}) \\ &= p(y) \prod_{j=1}^m p(x_j | y) \quad \text{naïve bayes assumption} \end{aligned}$$

$p(x_i | y)$ 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

$p(x | y)$

Binary features:

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

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)
- ...

Basic steps for probabilistic modeling

Step 1: pick a model	Probabilistic models Which model do we use, i.e. how do we calculate $p(\text{feature}, \text{label})$?
Step 2: figure out how to estimate the probabilities for the model	How do train the model, i.e. how do we estimate the probabilities for the model?
Step 3 (optional): deal with overfitting	How do we deal with overfitting?

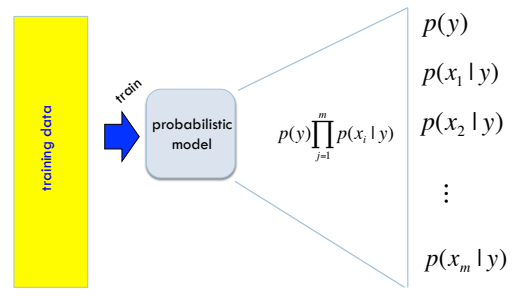
Obtaining probabilities



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 that any review is about Pinot Noir?
- What is the probability that a particular review is about Pinot Noir?

Obtaining probabilities



training data

train

probabilistic model

$p(y)$
 $p(x_1 | y)$
 $p(x_2 | y)$
 \vdots
 $p(x_m | y)$

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

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?

Maximum Likelihood Estimation (MLE)

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

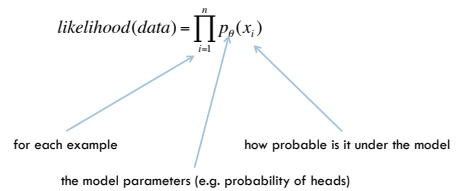
You flip a coin 100 times. 60 times you get heads and 40 times you get tails.

What is the MLE estimate for heads?

$p(\text{head}) = 0.60$ why?

Likelihood

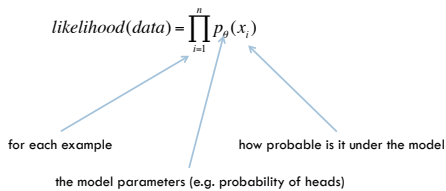
The likelihood of a data set is the probability that a particular model (i.e. a model and estimated probabilities) assigns to the data



Likelihood

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

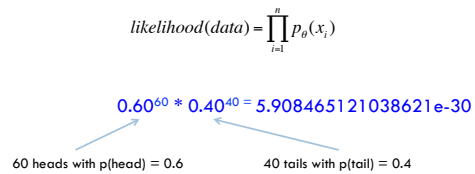
What is the likelihood of this data with $\Theta = p(\text{head}) = 0.6$?



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 $\Theta = p(\text{head}) = 0.6$?



MLE example

Can we do any better? $likelihood(data) = \prod_i p(x_i)$

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

60 heads with $p(head) = 0.6$

40 tails with $p(tail) = 0.4$

What about $p(head) = 0.5$?

MLE example

Can we do any better? $likelihood(data) = \prod_i p(x_i)$

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

60 heads with $p(head) = 0.6$

40 tails with $p(tail) = 0.4$

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

60 heads with $p(head) = 0.5$

40 tails with $p(tail) = 0.5$

MLE example

Can we do any better? $likelihood(data) = \prod_i p(x_i)$

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

60 heads with $p(head) = 0.6$

40 tails with $p(tail) = 0.4$

What about $p(head) = 0.7$?

MLE example

Can we do any better? $likelihood(data) = \prod_i p(x_i)$

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

60 heads with $p(head) = 0.6$

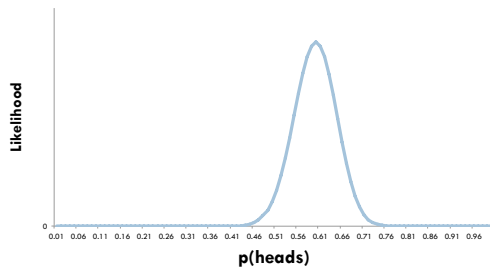
40 tails with $p(tail) = 0.4$

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

60 heads with $p(head) = 0.7$

40 tails with $p(tail) = 0.3$

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:

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

Why is this ok?

$$= \arg \max_{\theta} \sum_{i=1}^n \log(p(x_i))$$

Calculating MLE

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

$$MLE = \arg \max_{\theta} \sum_{i=1}^n \log(p(x_i))$$

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.

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

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

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

$$MLE = \arg \max_{\theta} 60 \log(\theta) + 40 \log(1 - \theta)$$

How do we find the max?

Calculating MLE

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

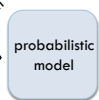
$$\begin{aligned} \frac{d}{d\theta} 60 \log(\theta) + 40 \log(1-\theta) &= 0 \\ \frac{60}{\theta} - \frac{40}{1-\theta} &= 0 \\ \frac{40}{1-\theta} &= \frac{60}{\theta} \\ 40\theta &= 60 - 60\theta \\ 100\theta &= 60 \\ \theta &= \frac{60}{100} \quad \text{Yay!} \end{aligned}$$

Calculating MLE

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

$$\begin{aligned} \frac{d}{d\theta} a \log(\theta) + b \log(1-\theta) &= 0 \\ &\dots \\ \theta &= \frac{a}{a+b} \end{aligned}$$

MLE estimation for NB



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

\swarrow \searrow
 $p(y)$ $p(x_i | y)$

What are the MLE estimates for these?

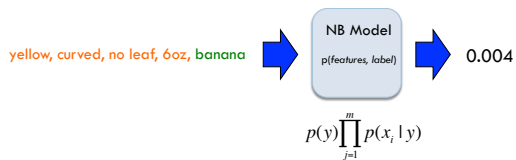
Maximum likelihood estimates

$$p(y) = \frac{\text{count}(y)}{n} \quad \frac{\text{number of examples with label}}{\text{total number of examples}}$$

$$p(x_i | y) = \frac{\text{count}(x_i, y)}{\text{count}(y)} \quad \frac{\text{number of examples with the label with feature}}{\text{number of examples with label}}$$

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

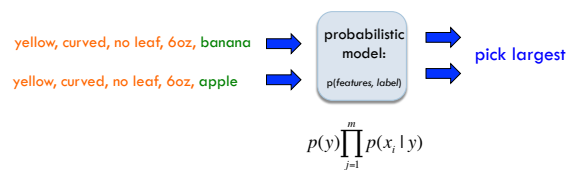
Naïve Bayes classification



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

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

Probabilistic models



$$\text{label} = \arg \max_{y \in \text{labels}} p(y) \prod_{j=1}^m p(x_j | y)$$

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

NB generative story



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

What is the generative story for the NB model?

NB generative story



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

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?

NB decision boundary

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

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

Some maths

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

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

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

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

Some more maths

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

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

(because x_i are binary)

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

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

And...

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

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

What does this look like?

And...

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

$$= \underset{y \in labels}{\operatorname{argmax}} \log(p(y)) + \underbrace{\sum_{i=1}^m \log(1 - p(x_i | y))}_b + \sum_{i=1}^m x_i \log \left(\frac{p(x_i | y)}{1 - p(x_i | y)} \right)$$

$w x + b$ What are the weights?
 Linear model !!!

NB as a linear model

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

How likely this feature is to be 1 given the label

How likely this feature is to be 0 given the label

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

Maximum likelihood estimation

Intuitive

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

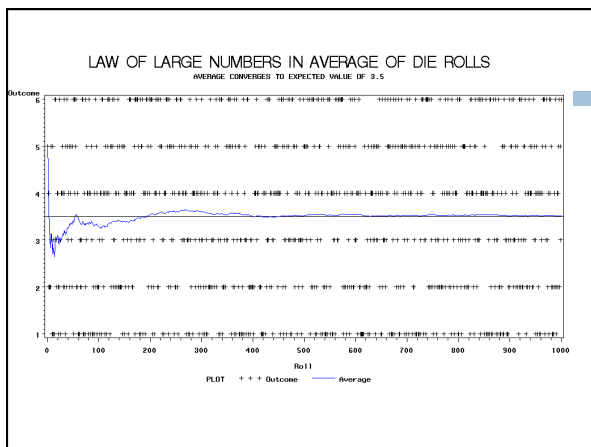
Problems?

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

Basic steps for probabilistic modeling

<p>Step 1: pick a model</p> <p>Step 2: figure out how to estimate the probabilities for the model</p> <div style="border: 1px solid blue; padding: 5px; margin-top: 10px;"> <p>Step 3 (optional): deal with overfitting</p> </div>	<p style="color: blue;">Probabilistic models</p> <p>Which model do we use, i.e. how do we calculate $p(\text{feature}, \text{label})$?</p> <p>How do train the model, i.e. how to we we estimate the probabilities for the model?</p> <p>How do we deal with overfitting?</p>
--	--

Coin experiment



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?

Back to parasitic gaps

$$p(\text{not_parasitic}) = 0.99999$$

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

Then probability of us finding one is $\sim 10\%$

- 90% of the time we won't find one and won't know anything (or assume $p(\text{parasitic}) = 0$)
- 10% of the time we would find one and incorrectly assume the probability is $1/10,000$ (10 times too large!)

Solutions?

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?