

PROBABILISTIC MODELS

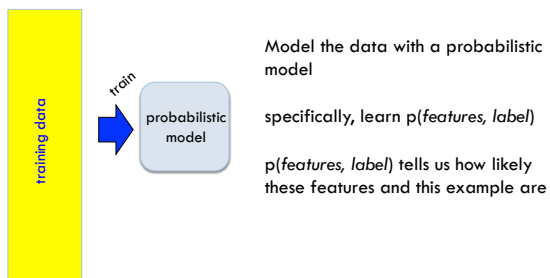
David Kauchak
CS451 – Fall 2013

Admin

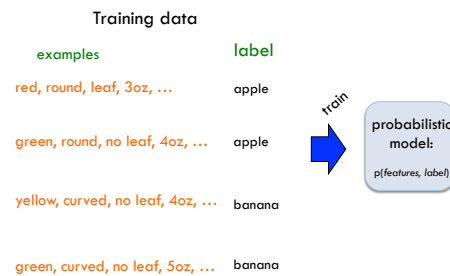
Assignment 6

- L2 normalization constant should be 2 (not 1)
- Just a handful of changes to the Perceptron code!

Probabilistic Modeling



An example: classifying fruit



Probabilistic models

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

A diagram showing the input "yellow, curved, no leaf, 6oz, banana" in multi-colored text. A blue arrow points to a rounded rectangular box labeled "probabilistic model: p(features, label)". A second blue arrow points from the box to the output "0.004".

Probabilistic model vs. classifier

Probabilistic model:

A diagram showing the input "yellow, curved, no leaf, 6oz, banana" in multi-colored text. A blue arrow points to a rounded rectangular box labeled "probabilistic model: p(features, label)". A second blue arrow points from the box to the output "0.004".

Classifier:

A diagram showing the input "yellow, curved, no leaf, 6oz" in multi-colored text. A blue arrow points to a rounded rectangular box labeled "probabilistic model: p(features, label)". A second blue arrow points from the box to the output "banana" in green text.

Probabilistic models: classification

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

A diagram showing the input "yellow, curved, no leaf, 6oz, banana" in multi-colored text. A blue arrow points to a rounded rectangular box labeled "probabilistic model: p(features, label)". A second blue arrow points from the box to the output "0.004".

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

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

Probabilistic models

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

A diagram showing two inputs. The first is "yellow, curved, no leaf, 6oz, banana" in multi-colored text, with a blue arrow pointing to a rounded rectangular box labeled "probabilistic model: p(features, label)". A second blue arrow points from the box to the output "0.004", which is circled in green. The second input is "yellow, curved, no leaf, 6oz, apple" in multi-colored text, with a blue arrow pointing to the same box. A second blue arrow points from the box to the output "0.00002".

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

Probabilistic model vs. classifier

Probabilistic model:

yellow, curved, no leaf, 6oz, banana → probabilistic model: $p(\text{features}, \text{label})$ → 0.004

Classifier:

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

Why probabilistic models?

Probabilistic models

Probabilities are nice to work with

- range between 0 and 1
- can combine them in a well understood way
- lots of mathematical background/theory
- an aside: to get the benefit of probabilistic output you can sometimes **calibrate** the confidence output of a non-probabilistic classifier

Provide a strong, well-founded groundwork

- Allow us to make clear decisions about things like regularization
- Tend to be much less "heuristic" than the models we've seen
- Different models have very clear meanings

Probabilistic models: big questions

Which model do we use, i.e. how do we calculate $p(\text{feature}, \text{label})$?

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

How do we deal with overfitting?

Same problems we've been dealing with so far

Probabilistic models	ML in general
Which model do we use, i.e. how do we calculate $p(\text{feature}, \text{label})$?	Which model do we use (decision tree, linear model, non-parametric)
How do train the model, i.e. how to we we estimate the probabilities for the model?	How do train the model?
How do we deal with overfitting?	How do we deal with overfitting?

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> <p>Step 3 (optional): deal with overfitting</p>	<p style="color: blue; text-align: center;">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>
---	--

Basic steps for probabilistic modeling

<div style="border: 1px solid blue; padding: 2px; margin-bottom: 10px;">Step 1: pick a model</div> <p>Step 2: figure out how to estimate the probabilities for the model</p> <p>Step 3 (optional): deal with overfitting</p>	<p style="color: blue; text-align: center;">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>
--	--

What was the data generating distribution?

The diagram illustrates a data generating distribution represented by a blue oval containing various fruits: an apple, a banana, and a green apple. Two green arrows point upwards from this oval to two separate groups of fruit. The group on the left is labeled 'Training data' and contains a mix of these fruits. The group on the right is labeled 'Test set' and also contains a mix of these fruits, representing a sample from the same underlying distribution.

Step 1: picking a model

What we're really trying to do is model is the data generating distribution, that is how likely the feature/label combinations are

The diagram shows a blue oval containing a single instance of each fruit type: one apple, one banana, and one green apple, representing the data generating distribution.

Some maths

$$p(\text{features}, \text{label}) = p(x_1, x_2, \dots, x_m, y)$$

$$= p(y)p(x_1, x_2, \dots, x_m | y)$$

What rule?

Some maths

$$p(\text{features}, \text{label}) = p(x_1, x_2, \dots, x_m, y)$$

$$= p(y)p(x_1, x_2, \dots, x_m | y)$$

$$= p(y)p(x_1 | y)p(x_2, \dots, x_m | y, x_1)$$

$$= p(y)p(x_1 | y)p(x_2 | y, x_1)p(x_3, \dots, x_m | y, x_1, x_2)$$

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

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: $2^{7000} = ?$

2⁷⁰⁰⁰

```

163169755662202064664588547837799519112430363743286235982084151272023162702352987080237879
4460004651996019099530984538645257892546513204107022110253546458647431585227076999373340842842
722420012281878260072931082617043194484266392077841250999968601694360066600112098175792966787
819625237700655294757256678055809293844627218640216108862600816099132874749204352087401101862
690842327501724602311293955235059054544214554772509509096307889478094683592939574112569473438
6191215296848474344406741204174020887540371869421701550220735398381224299258743537536161041593
43594557666561701790904172597025336526662682021808493892812699709528570890696375575541434487608
8248369941993802415197514510125127043829087280919538476302837811854024099998995964192277601255
3654911562403499947144160905730842429313962119536793701294479560024833357073898392029910322
3465980389530690429801740098017325210691307971242016963397230218353007589784519525848553710885
8195631737000743805167411189134617501484521767984296782842287373127422122022517597535994839257
029877907706355334790244938435386605125910795672914312162977887848185522928196541766009803989
97991481404749384215743515802030811510828640678970483829220546427576550737656054750702714
4662263487685709621261074762705203049488907208978593689047063428548531668665653271746660658185
609064849508012761754614572161769555751992117507514067775104496728590822558547771447242334900
764026321760892113525261241194338702680299044001838585057671926968975936612135888838680023840
9255673807750189147030466215099498385389520715495963392370267592041517264907070077833625108
3209283964807237954887069546621688046652112493076290091990717423550391351174415329737479300
8955830518884135334798464113680004999403724560035428811232632821866113106455077289922996946
915601858083982074170466852124388152026099584696588161375826382921029547343888632163627122302
921227953846843548353710400407789177417026363662072695437517780741313455101810094688094
0781122057380335371124632958916237089580476224595091825301636909236240671411644331656159828058
3720783439888562390892028440902553829376
    
```

Any problems with this?

Full distribution tables

x ₁	x ₂	x ₃	...	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?

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 whatever 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 dice
- price of tea in England is independent of the whether or not you pass AI

independent or dependent?

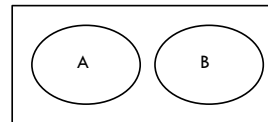
Catching a cold and having cat-allergy

Miles per gallon and driving habits

Height and longevity of life

Independent variables

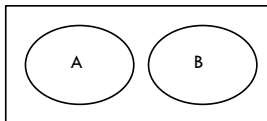
How does independence affect our probability equations/properties?



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 of 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 other features 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

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

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

$p(x_i | y)$ 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 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

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 $p(\text{feature}, \text{label})$?

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

How do we deal with overfitting?

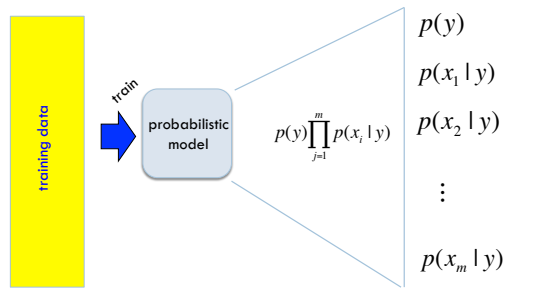
Obtaining probabilities



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* that based on data, though:

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

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.

What is the MLE for heads?

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

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.

What is the likelihood of the data under this model (each coin flip is a data point)?

MLE example

You flip a coin 100 times. 60 times you get heads.

MLE for heads: $p(\text{head}) = 0.60$

What is the likelihood of the data under this model (each coin flip is a data point)?

$$\text{likelihood}(\text{data}) = \prod_i p(x_i)$$

$$\log(0.60^{60} * 0.40^{40}) = -67.3$$

MLE example

Can we do any better?

$$\text{likelihood}(\text{data}) = \prod_i p(x_i)$$

$p(\text{heads}) = 0.5$

$$\log(0.50^{60} * 0.50^{40}) = -69.3$$

$p(\text{heads}) = 0.7$

$$\square \log(0.70^{60} * 0.30^{40}) = -69.5$$