Announcements

HW3 due tonight

Tutorial feedback back tonight

Tutorial due Apr 6 (Submission)

Tutorial peer evaluation: Apr 11 (Peer evaluation)

15-388/688 - Practical Data Science: Maximum likelihood estimation, naïve Bayes

Pat Virtue Carnegie Mellon University Spring 2022

Slide credits: CMU AI, Zico Kolter

Outline

Maximum likelihood estimation

Naive Bayes

Machine learning and maximum likelihood

Outline

Maximum likelihood estimation

Naive Bayes

Machine learning and maximum likelihood

Challenge

comp likelihood en max 11 11 likelihood

Assume that exam scores are drawn independently from the same Gaussian (Normal) distribution.

Given three exam scores 75, 80, 90, which pair of parameters is a better fit?

- A) Mean 80, standard deviation 3
- B) Mean 85, standard deviation 7

Use a calculator/computer.

Gaussian PDF: $p(x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

Estimating the parameters of distributions

We're moving now from probability to statistics

Prob:
$$O \longrightarrow P$$

Stat: Data $\xrightarrow{P} \hat{O}$
 $X \longrightarrow O$

Estimating the parameters of distributions

We're moving now from probability to statistics

The basic question: given some data $x^{(1)}, \dots, x^{(m)}$, how do I find a distribution that captures this data "well"?

In general (if we can pick from the space of all distributions), this is a hard question, but if we pick from a particular *parameterized family* of distributions $p(X;\theta)$, the question is (at least a little bit) easier

Question becomes: how do I find parameters θ of this distribution that fit the data?

log ab = log a + log b Maximum likelihood estimation i.i.d.

Given a distribution $p(X; \theta)$, and a collection of observed (independent) data points $x^{(1)}, \dots, x^{(m)}$, the probability of observing this data is simply

$$\xrightarrow{p(x^{(1)}, \dots, x^{(m)}; \theta)} = \mathcal{T}_p(x^{(j)}; \theta) = \mathcal{L}(\theta)$$

Basic idea of maximum likelihood estimation (MLE): find the parameters that $l(\theta) = \log \pi p(x^{(j)}; \theta)$ $= \sum \log p(x^{(j)}; \theta)$ maximize the probability of the observed data

$$\underset{\theta}{\text{maximize}} \prod_{i=1}^{m} p(x^{(i)}; \theta) \equiv \underset{\theta}{\text{maximize}} \underset{\theta}{\overset{\alpha} \land \overset{\alpha}{ }} p(x^{(i)}; \theta) = \underset{\theta}{\overset{\alpha} \land \overset{\alpha}{ }$$

where $\ell(\theta)$ is called the **log likelihood** of the data

Seems "obvious", but there are many other ways of fitting parameters

Parameter estimation for Bernoulli

Simple example: Bernoulli distribution

$$p(X = 1; \phi) = \phi, \qquad p(X = 0; \phi) = 1 - \phi$$

Given observed data $x^{(1)}, \dots, x^{(m)}$, the "obvious" answer is:

$$\hat{\phi} = \frac{\#1's}{\#Total} = \frac{\sum_{i=1}^{m} x^{(i)}}{m} \qquad \frac{3}{5} \qquad MLE$$

But why is this the case?

Maybe there are other estimates that are just as good, i.e.?

$$\phi = \frac{\sum_{i=1}^{m} x^{(i)} + 1}{m+2}$$

20% heads

=0,2

Likelihood for Bernoulli

The likelihood for Bernoulli is given by

$$L(\phi) = \prod_{i=1}^{m} p(x^{(i)}; \phi)$$

= $\prod \phi^{m} \xi \chi^{(i)} = i \xi$
$$= \prod \phi^{m} \xi \chi^{(i)} = i \xi \qquad (1 - \phi)$$

Let's say we have a dataset of 3 heads and 2 tails:

MLE for Bernoulli 26

Maximum likelihood solution for Bernoulli is given by $\max_{\phi} \min_{i=1}^{m} p(x^{(i)}; \phi) = \max_{\phi} \min_{\phi} \min_{i=1}^{m} p(x^{(i)}; \phi) = \max_{\phi} \min_{\phi} \max_{\phi} \max_{\phi} \min_{\phi} \max_{\phi} \max_{\phi}$

Taking the negative log of the optimization objective (just to be consistent with our usual notation of optimization as minimization) $\operatorname{Min} \log L(\phi) = \operatorname{Min} \bigotimes_{i=1}^{\infty} x^{(i)} \log \phi + (1-x^{(i)}) \log (1-\phi)$

Derivative with respect to ϕ is given by

$$\frac{d}{d\phi}\ell(\phi) = \sum_{i=1}^{m} \left(\frac{x^{(i)}}{\phi} - \frac{1 - x^{(i)}}{1 - \phi}\right) = \frac{\sum_{i=1}^{m} x^{(i)}}{\phi} - \frac{\sum_{i=1}^{m} (1 - x^{(i)})}{1 - \phi}$$

MLE for Bernoulli, continued

Setting derivative to zero gives:

to zero gives:

$$\frac{\sum_{i=1}^{m} x^{(i)}}{\phi} - \frac{\sum_{i=1}^{m} (1 - x^{(i)})}{1 - \phi} \equiv \frac{a}{\phi} - \frac{b}{1 - \phi} = 0$$

$$\Rightarrow (1 - \phi)a = \phi b$$

$$\Rightarrow \phi = \frac{a}{a + b} = \frac{\sum_{i=1}^{m} x^{(i)}}{m}$$

So, we have shown that the "natural" estimate of ϕ actually corresponds to the maximum likelihood estimate

MLE for Gaussian, briefly

For Gaussian distribution

$$p(x;\mu,\sigma^2) = (2\pi\sigma^2)^{-1/2} \exp(-(1/2)(x-\mu)^2/\sigma^2)$$

Log likelihood given by:

$$\ell(\mu, \sigma^2) = -m\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2}\sum_{i=1}^m \frac{(x^{(i)} - \mu)^2}{\sigma^2}$$

Derivatives (see if you can derive these fully):

$$\frac{d}{d\mu}\ell(\mu,\sigma^2) = -\frac{1}{2}\sum_{i=1}^m \frac{x^{(i)} - \mu}{\sigma^2} = 0 \implies \mu = \frac{1}{m}\sum_{i=1}^m x^{(i)}$$
$$\frac{d}{d\sigma^2}\ell(\mu,\sigma^2) = -\frac{m}{2\sigma^2} + \frac{1}{2}\sum_{i=1}^m \frac{(x^{(i)} - \mu)^2}{(\sigma^2)^2} = 0 \implies \sigma^2 = \frac{1}{m}\sum_{i=1}^m (x^{(i)} - \mu)^2$$

Outline

Maximum likelihood estimation

Naive Bayes

Machine learning and maximum likelihood

SPAM Classification

Example

Training Data

- Spam? E-mail body
 - 1 Money is free now
 - 0 Pat teach 388
 - 0 Pat free to teach
 - 1 Sir money to teach
 - 1 Pat free money now
 - 0 Teach 388 now
 - 0 Pat to teach 301

Vocabulary	
388	X,
free	X
is	X,
money	X4
now	•
Pat	•
Sir	e
teach	
to	
tomorrow	X

10

<u>Test Data</u>	
Spam?	E-mail body
$p(Y=1 X_{i})$	Pat teach now x_{10}

Poll 1

Assume:

Y is a binary random variable representing whether or not the email is spam, and X_i is a binary random variable representing whether or not the *i*-th word is in the email.

With a vocabulary of size 10, how may probability values are in the following probability table?

I	5	$D(V \mid V \mid V)$		<u>Vocabulary</u>
		$F(I \mid \Lambda_1, \dots, \Lambda_{10})$	1	388
<i>A.</i> 10		P(y-n x=0 x=0 X=0)	2	free
<i>B.</i> 11			3	is
C 110			4	money
<i>C.</i> 110			5	now
<i>D.</i> 22			6	Pat
$\sim E_{*} 2^{10}$	48%		7	Sir
(E) 211	2407		8	teach
I . <i>Z</i>	C 1 10		9	to
			10	tomorrow

Naive Bayes modeling

Naive Bayes is a machine learning algorithm that rests relies heavily on probabilistic modeling

But, it is also interpretable according to the three ingredients of a machine learning algorithm (hypothesis function, loss, optimization), more on this later

Basic idea is that we model input and output as random variables $X = (X_1, X_2, ..., X_n)$ (several Bernoulli, categorical, or Gaussian random variables), and Y (one Bernoulli or categorical random variable), goal is to find p(Y|X)

Naive Bayes assumptions $P(X, X_2, X_3...X_b | Y)$ We're going to find p(Y|X) via Bayes' rule $p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)} = \frac{p(X|Y)p(Y)}{\sum_{y} p(X|y)p(y)} \quad x = \frac{p(X|Y)P(Y)}{\sum_{y} p(X|y)p(y)}$

The denominator is just the sum over all values of Y of the distribution specified by the numeration, so we're just going to focus on the p(X|Y)p(Y) term

Modeling full distribution p(X|Y) for high-dimensional X is not practical, so we're going to make the **naive Bayes assumption**, that the elements X_i are conditionally independent given Y

$$p(X|Y) = \prod_{i=1}^{n} \frac{p(X_i|Y)}{p(X_i|Y)} + p(X_i|Y) + p(X_i|Y)$$

Poll 2

Assume:

Y is a binary random variable representing whether or not the email is spam, and X_i is a binary random variable representing whether or not the *i*-th word is in the email.

True or False: $P(X_1 = 1 | Y = 0) = P(X_1 = 1 | Y = 1)$

Vocabulary 388 1 2 free 3 is 4 money 5 now 6 Pat Sir 7 8 teach 9 to 10 tomorrow

Modeling individual distributions

We're going to explicitly model the distribution of each $p(X_i|Y)$ as well as p(Y)

We do this by specifying a distribution for p(Y) and a *separate* distribution and for each $p(X_i|Y = y)$

So assuming, for instance, that Y_i and X_i are binary (Bernoulli random variables), then we would represent the distributions

$$p(Y; \phi_{Y=1}), \quad p(X_i | Y = 0; \phi_{Y=0,i}), \quad p(X_i | Y = 1; \phi_{Y=1,i})$$

We then estimate the parameters of these distributions using MLE, i.e.

$$\phi_{Y=1} = \frac{\sum_{j=1}^{m} y^{(j)}}{m}, \quad \phi_{y,i} = \frac{\sum_{j=1}^{m} x_i^{(j)} \cdot \mathbb{1}\{y^{(j)} = y\}}{\sum_{j=1}^{m} \mathbb{1}\{y^{(j)} = y\}}$$

Making predictions Given some new data point x, we can now compute the probability of each class

$$p(Y = y \mid x) \propto \underline{p(Y = y)} \prod_{i=1}^{n} p(x_i \mid Y = y) = \phi_y \prod_{i=1}^{n} (\phi_{y,i})^{x_i} (1 - \phi_1^y)^{1-x_i}$$

After you have computed the right-hand side, just normalize (divide by the sum over all y) to get the desired probability

Alternatively, if you just want to know the most likely Y, just compute each righhand side and take the maximum

Example

Potential issues

Problem #1: when computing probability, the product $p(y) \prod_{i=1}^{n} p(x_i | y)$ quickly goes to zero to numerical precision

Solution: compute log of the probabilities instead

$$\log p(y) + \sum_{i=1}^{n} \log p(x_i|y)$$

Problem #2: If we have never seen either $X_i = 1$ or $X_i = 0$ for a given y, then the corresponding probabilities computed by MLE will be zero

Solution: Laplace smoothing, "hallucinate" one $X_i = 0/1$ for each class $\phi_{y,i} = \frac{\sum_{j=1}^m x_i^{(j)} \cdot \mathbb{1}\{y^{(j)} = y\} + 1}{\sum_{j=1}^m \mathbb{1}\{y^{(j)} = y\} + 2} \xrightarrow{\qquad + \infty}$

Categorical class

Let Y be the random variable for a class that takes on one of K possible categories $\{1, ..., K\}$ (rather than binary as we were doing before)

$$P(Y = y) = \phi_y = \frac{\sum_{j=1}^m \mathbb{1}\{y^{(j)} = y\}}{m}$$

Y	<i>X</i> ₁	<i>X</i> ₂	Y	<i>X</i> ₁	<i>X</i> ₂
cat			1		
dog			2		
rat			3		
rat			3		
cat			1		
cat			1		

$$\phi_1 = \frac{3}{6}$$

 $\phi_2 = \frac{3}{6}$
 $\phi_2 = \frac{3}{6}$
 $\phi_2 = \frac{3}{6}$

Categorical feature conditioned on class

Assume the *i*-th feature takes on one of *K* possible categories $\{1, ..., K\}$ (rather than binary as we were doing before)

$$P(X_i = k \mid Y = y) = \phi_{y,i,k} = \frac{\sum_{j=1}^m \mathbb{1}\left\{x_i^{(j)} = k\right\} \cdot \mathbb{1}\{y^{(j)} = y\}}{\sum_{j=1}^m \mathbb{1}\{y^{(j)} = y\}}$$

Y	<i>X</i> ₁	<i>X</i> ₂		Y	<i>X</i> ₁	<i>X</i> ₂
cat	blue	wood		→ 1		3
dog	blue	metal		2	1	1
rat	green	metal		3	2	1
rat	red	paper		3	3	2
cat	red	wood	~	→ 1	3	3
cat	blue	wood	-	→ 1	(1)	3

 $\phi_{y=1,i=1} = \frac{2}{3}$

Though naive Bayes is often presented as "just" counting, the value of the maximum likelihood interpretation is that it's clear how to model $p(X_i|Y)$ for non-categorical random variables

Example: if x_i is real-valued, we can model $p(X_i|Y = y)$ as a Gaussian $p(x_i|y; \mu^y, \sigma_y^2) = \mathcal{N}(x_i; \mu^y, \sigma_y^2)$

with maximum likelihood estimates

$$\mu_{y} = \frac{\sum_{j=1}^{m} x_{i}^{(j)} \cdot \mathbb{1}\{y^{(j)} = y\}}{\sum_{j=1}^{m} \mathbb{1}\{y^{(j)} = y\}}, \ \sigma_{y}^{2} = \frac{\sum_{j=1}^{m} (x_{i}^{(j)} - \mu^{y})^{2} \cdot \mathbb{1}\{y^{(j)} = y\}}{\sum_{j=1}^{m} \mathbb{1}\{y^{(j)} = y\}}$$

All probability computations are exactly the same as before (it doesn't matter that some of the terms are probability densities)

Gaussian features conditioned on class

$$\mu_{y} = \frac{\sum_{j=1}^{m} x_{i}^{(j)} \cdot \mathbb{1}\{y^{(j)} = y\}}{\sum_{j=1}^{m} \mathbb{1}\{y^{(j)} = y\}}, \ \sigma_{y}^{2} = \frac{\sum_{j=1}^{m} (x_{i}^{(j)} - \mu^{y})^{2} \cdot \mathbb{1}\{y^{(j)} = y\}}{\sum_{j=1}^{m} \mathbb{1}\{y^{(j)} = y\}}$$



 $\mu_{3,2000}$ M3, time = 22.5 $0^{2}_{3, score}$ = 25 $=((70-25)^{2}+(60-65)^{2})/2$

Outline

Maximum likelihood estimation

Naive Bayes

Machine learning and maximum likelihood

Machine learning via maximum likelihood

Many machine learning algorithms (specifically the loss function component) can be interpreted probabilistically, as maximum likelihood estimation



Logistic probability model





i.e., the least-squares loss function can be viewed as MLE under Gaussian errors

Other approaches possible too: absolute loss function can be viewed as MLE under Laplace errors