Applied Machine Learning

Naive Bayes

Reihaneh Rabbany



COMP 551 (winter 2023) 1

Learning objectives

- the assumption of Naive Bayes classifier
- what does learning and prediction steps involve?
- different likelihood functions
- Bayesian parameter learning in Naive Bayes
- practical considerations

Bayes rule for classification

given

- the prior probability of each class
- likelihood of observations given the class

use Bayes rule for classification



Bayes rule for classification example

 $x \in \{-,+\}$ input: test results, a single binary feature $y \in \{ ext{yes, no}\}$ label: patient has cancer

Generative classification

traininglearn the following distributions from the data $\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$ prior probability of each class $p(y = c) \forall c \in \{1, \dots, C\}$ likelihood of data for each classp(x|y = c)predictionuse the Bayes rule to get the posterior class probability $p(y = c \mid x) \propto p(c)p(x|c)$

generative classifier because we are learning the joint data distribution p(x,y) = p(y)p(x|y)we can generate new data from this joint distribution

in a **discriminative classifier** we directly learn p(y|x)

Generative classification

prior class probability: frequency of observing this label

 $p(y = c \mid x) = \frac{p(c)p(x|c)}{p(x)}$ posterior probability of a given class $p(x) = \frac{p(c)p(x|c)}{p(x)}$ $p(x) = \frac{p(c)p(x|c)}{p(x)}$ $p(x) = \frac{p(c)p(x|c)}{p(x)}$ $p(x) = \frac{p(c)p(x|c)}{p(x)}$ $p(x) = \frac{p(c)p(x|c)}{p(x)}$

Some generative classifiers:

- Gaussian Discriminant Analysis: the likelihood is multivariate Gaussian
- Naive Bayes: decomposed likelihood

Naive Bayes model

assumption about the likelihood $\ p(x|y) = \prod_{d=1}^{D} p(x_d|y)$.

when is this assumption correct?

when features are **conditionally independent** given the label $x_i \perp x_j \mid y$

knowing the label, the value of one input feature gives us no information about the other input features

number of input features

How is the likelihood derived from this independence assumption?

chain rule of probability (true for any distribution)

 $p(x|y) = p(x_1|y)p(x_2|y,x_1)p(x_3|y,x_1,x_2) \dots p(x_D|y,x_1,\dots,x_{D-1})$ conditional independence assumption x_1, x_2 give no extra information, so $p(x_3|y,x_1,x_2) = p(x_3|y)$

Naive Bayes: objective

given the training dataset $\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$

a generative classifier maximizes the joint likelihood (or log-likelihood)

$$\begin{split} L(\pi,\theta;\mathcal{D}) &= \prod_{n\in\mathcal{D}} p(x^{(n)},y^{(n)};\pi,\theta) & \pi, \theta \text{ are the model parameters} \\ \ell(\pi,\theta) &= \sum_{n} \log p(x^{(n)},y^{(n)};\pi,\theta) & p(x^{(n)}|y^{(n)};\theta) \\ &= \sum_{n} \left[\log p(y^{(n)};\pi) + \log p(x^{(n)}|y^{(n)};\theta) \right] \\ &= \sum_{n} \log p(y^{(n)};\pi) + \sum_{n} \log p(x^{(n)}|y^{(n)};\theta) & \longleftarrow \text{ using Naive Bayes assumption here} \\ &= \sum_{n} \log p(y^{(n)};\pi) + \sum_{d} \sum_{n} \log p(x^{(n)}_{d}|y^{(n)};\theta_{d}) & p(x|y) = \prod_{d=1}^{D} p(x_{d}|y) \\ &= \sum_{n} \log p(y^{(n)};\pi) + \sum_{d} \sum_{n} \log p(x^{(n)}_{d}|y^{(n)};\theta_{d}) & p(x|y) = \prod_{d=1}^{D} p(x_{d}|y) \\ &= \sum_{n} \log p(y^{(n)};\pi) + \sum_{d} \sum_{n} \log p(x^{(n)}_{d}|y^{(n)};\theta_{d}) & p(x|y) = \prod_{d=1}^{D} p(x_{d}|y) \\ &= \sum_{n} \log p(x^{(n)}_{d};\pi) + \sum_{d} \sum_{n} \log p(x^{(n)}_{d}|y^{(n)};\theta_{d}) & p(x|y) = \prod_{d=1}^{D} \log p(x_{d}|y) \\ &= \sum_{d=1}^{D} \log$$

separate max-likelihood problems for prior and each feature x_d given the label

Prior class probabilities

class probabilities prior to looking at the features

for binary classification, class probability is given by Bernoulli $p(y; \pi) = \pi^y (1 - \pi)^{1-y}$

recall the max-likelihood estimate for Bernoulli

$$rg\max_{\pi}\sum_n \log p(y^{(n)};\pi) = rac{1}{N}\sum_n y^{(n)}$$

for multi-class classification, class probability is given by categorical distribution $p(y;\pi) = \prod_{c=1}^C \pi_c {}^{\mathbb{I}(y=c)} = \pi_y$ note that in this case π is a vector

max-likelihood estimate is again given by empirical frequencies

$$rg\max_{\pi_c}\sum_{n \in T_c} \log p(y^{(n)};\pi) = rac{N(y=c)}{N}$$
 frequency of class c in our dataset $\pi^* = [rac{N_1}{N}, \dots, rac{N_C}{N}]$

In both cases we learn the prior simply as the class frequencies in the training data

Naive Bayes: objective

given the training dataset $\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$

a generative classifier maximizes the joint likelihood (or log-likelihood)



separate max-likelihood problems for prior and each feature x_d given the label

Likelihood terms

likelihood terms $\, p(x_d | y; heta_d) \,$

- encode our assumption about the *generative process*
- different types of features require different forms of likelihood
 - Bernoulli for binary features
 - Categorical for categorical features
 - Multinomial for "count" features
 - Gaussian is one option for continuous feature
- each feature x_d may use a different likelihood form
- separate maximum conditional likelihood estimate for each feature

$$rg\max_{ heta_d}\sum_{n=1}^N \log p(x_d^{(n)} \mid y^{(n)}; heta_d)$$

Bernoulli Naive Bayes

for a binary **feature** likelihood is Bernoulli

$$\left\{ egin{array}{l} p(x_d \mid y=0; heta_d) = ext{Bernoulli}(x_d; heta_{d,0}) \ p(x_d \mid y=1; heta_d) = ext{Bernoulli}(x_d; heta_{d,1}) \end{array}
ight.$$

one parameter per label

short form: $p(x_d \mid y; \theta_d) = \text{Bernoulli}(x_d; \theta_{d,y})$

max-likelihood estimation is similar to what we saw for the prior



0

Covid-19 classification example

each patient has seven binary features $\,x\in\{0,1\}^7$

we have a dataset of N=1000 patients, where 200 had covid-19

learning:

learn the prior: $\pi = \frac{N(y=1)}{N} = .2$ Bernoulli $(y;\pi)$

for each

prob

$$\theta_{1,1}$$
 $\theta_{2,1}$
Fever 88% Fatigue 38% Sore throat 15% Dry cough 68% Muscle pain 5% J9% $\theta_{7,1}$ second to the second secon

symptom d:
ability of symptom
$$x_d = 1$$
 given $y = \begin{cases} 1 \quad \theta_{d,1} = \frac{N(y=1,x_d=1)}{N(y=1)} & \text{Bernoulli}(x_d|y=1;\theta_{d,1}) \\ 0 \quad \theta_{d,0} = \frac{N(y=0,x_d=1)}{N(y=0)} & \text{Bernoulli}(x_d|y=0;\theta_{d,0}) \end{cases}$

prediction:

for a new patient x calculate unnormalized posterior

$$\begin{cases} \tilde{p}(y=0|x) = \text{Bernoulli}(0;\pi) \prod_{d} \text{Bernoulli}(x_{d};\theta_{d,0}) & \text{normalize it } p(y=1|x) = \frac{\tilde{p}(y=1|x)}{\tilde{p}(y=0|x) + \tilde{p}(y=1|x)} \\ \tilde{p}(y=1|x) = \text{Bernoulli}(1;\pi) \prod_{d} \text{Bernoulli}(x_{d};\theta_{d,1}) & \text{normalize it } p(y=1|x) = \frac{\tilde{p}(y=1|x)}{\tilde{p}(y=0|x) + \tilde{p}(y=1|x)} \end{cases}$$

Disease diagnosis example

what changes in **multi-class** setting?



earn the prior:
$$\ \pi_{m{c}} = rac{N(y=c)}{N}$$

for each symptom d:

 $p(x_d|y; heta_d)$

learn the conditional likelihood: probability of symptom $x_d = 1$ given $y = \begin{cases} 0 \\ \vdots \\ C \end{cases}$

how many parameters in our model?

binary classification, binary features 1+2D

multi-class classification, binary features C+CD

Symptoms are common Symptoms occur sometimes										
OSym	ptoms are uncommon	Symptoms are	rare 🚫 (Doesn't have	e these symptoms					
	Symptom	COVID-19	Flu	Cold	Seasonal allergies					
	Body aches	\bigcirc		\bigcirc	and the second s					
	Cough				Θ					
Å	Diarrhea	0	\bigcirc	and the second s	\bigcirc					
Z^{z^2}	Fatigue			\bigcirc	Θ					
	Fever			And the second s	and the second s					
	Headaches	\bigcirc		and the second	Θ					
0,	Itchy or watery eyes	S Salara	\bigcirc	\bigcirc						
	Loss of smell or tas	te 😜	And the second s	and the second s						
	Nausea or vomiting	\circ	\bigcirc	And the second s	\bigcirc					
	Runny / stuffy nose		\bigcirc							
A	Shortness of breath	n	\bigcirc	and the second	Θ					
	Sneezing	and the second s	\bigcirc							
	Sore throat	Θ	\bigcirc		Θ					
Sourc	e: WHO, CDC			(NEWS					

Document classification

example

e.g., spam filtering

words in our vocabulary

N = 5

 $x^{(n)} \in \{0,1\}^D$ each document (email) is one instance

 $x^{(n)}_{\scriptscriptstyle A}=1$ if the word d appears in document n

classify the documents based on this bag of words representation

learning:

MLE for the prior $\text{Bernoulli}(y; \pi)$ (spam frequency in our dataset)

MLE for the likelihood terms $\text{Bernoulli}(x; \theta_{d,y})$

(frequency of word (d) in spam/non-spam documents)

prediction:

calculate the posterior $p(y|x) \propto \text{Bernoulli}(y;\pi) \prod_d \text{Bernoulli}(x_d;\theta_{d,y})$

D=7



Document classification example



label

Why Naive Bayes assumption?

Naive Bayes assumption $\ p(x|y) = \prod_d p(x_d|y)$

what if we did not make this assumption?

consider the **spam filtering example:**

- D can be very large
- with the Naive Bayes assumption: learn the frequency of each word (d) in spam/non-spam documents
- without it: learn the frequency of each possible subset of words in spam/non-spam documents

e.g., for x = [1, 1, 0, 0, 0, 1, 0] we need to estimate p(x|y)

problems

- many combinations of words may not appear in even one document
- we need exponentially more parameters
- even for large datasets, this could lead to overfitting



Bayesian Naive Bayes

using MLE in Naive Bayes can lead to overfitting

example let's classify this new sentence: that dog was my puppy

$$egin{aligned} ilde{p}(y=1|x) &= rac{4}{5} imes rac{1}{4} imes rac{0}{4} imes \ldots = 0 \ ilde{p}(y=0|x) &= rac{1}{5} imes rac{0}{1} imes rac{0}{1} imes \ldots = 0 \end{aligned}$$

the problem is that the word "is" appears in all instances max-likelihood estimate $\ heta_{1,1}= heta_{1,0}=1$

we can solve this by being Bayesian in parameter learning:

instead of maintaining a point estimates $\pi, \theta_{d,y}$ we maintain distributions $p(\pi), p(\theta_{d,y})$ for $y \in \{0, 1\}, d$ start from separate prior for each parameter $p(\pi), p(\theta_{d,y})$

calculate the likelihood $\prod_n p(y^{(n)}|\pi)$

update with observed frequencies in the dataset



Bayesian Naive Bayes

start from separate prior for each parameter $p(\pi) = \text{Beta}(\pi; \alpha^{\pi}, \beta^{\pi}) \quad p(\theta_{d,y}) = \text{Beta}(\theta; \alpha^{\theta}, \beta^{\theta})$

calculate the posterior $p(\pi | \mathcal{D}) = ext{Beta}(\pi; lpha^\pi + N(y=1), eta^\pi + N(y=0))$

 $p(heta_{d,ar{y}}|\mathcal{D}) = ext{Beta}(heta_{d,ar{y}};lpha^ heta+N(y=ar{y},x_d=1),eta^ heta+N(y=ar{y},x_d=0))$

use posterior predictive for a new instance (x) $p(y = 1|x, D) = \int_{\theta, \pi} p(y = 1|\pi) p(\pi|D) \prod_d p(x_d|\theta_{d,1}) p(\theta_{d,1}|D) d\theta d\pi$

individual posterior predictives
$$= \left(\int_{\pi} p(y=1|\pi)p(\pi|\mathcal{D})d\pi\right) \prod_{d} \left(\int_{\theta_{d,1}} p(x_{d}|\theta_{d,1})p(\theta_{d,1}|\mathcal{D})d\theta\right)$$

for Beta distribution, we simply used the posterior mean (and dropped the integral)

$$ilde{p}(y=1|x,\mathcal{D}) = rac{lpha^{\pi}+N(y=1)}{lpha^{\pi}+eta^{\pi}+N} \prod_{d} \left(rac{lpha^{ heta}+N(y=1,x_d=1)}{lpha^{ heta}+eta^{ heta}+N(y=1)}
ight)^{x_d} \left(rac{eta^{ heta}+N(y=1,x_d=0)}{lpha^{ heta}+eta^{ heta}+N(y=1)}
ight)^{(1-x_d)}$$
 recall: Laplace smoothing

compare with our previous prediction (using MLE)

$$ilde{p}(y=1|x,\mathcal{D}) = rac{N(y=1)}{N} \prod_d \left(rac{N(y=1,x_d=1)}{N(y=1)}
ight)^{x_d} \left(rac{N(y=1,x_d=0)}{N(y=1)}
ight)^{(1-x_d)}$$
 we are simply adding a constant to various frequencies

Bayesian Naive Bayes

let's classify this new sentence using Laplace smoothing:

example

puppy it is a puppy it is a kitten 0 0 0 it is a cat 0 0 that is a dog and this is a pen 0 0 0 it is a matrix 0 0 y=0 $rac{1}{1}$ $rac{1}{1}$ $rac{0}{1}$ $rac{0}{1}$ $rac{0}{1}$

is

cat

 $\frac{3}{4}$ $\frac{4}{4}$ $\frac{1}{4}$ $\frac{1}{4}$ $\frac{1}{4}$ $\frac{1}{4}$ $\frac{4}{4}$

pen

it

 $ilde{p}(y=0|x) = rac{1+1}{5+2} imes rac{0+1}{1+2} imes rac{0+1}{1+2} imes rac{0+1}{1+2} imes rac{1+1}{1+2} imes rac{1+1}{1+2} imes rac{0+1}{1+2} imes rac{0+1}{1+2} pprox 0.0052$

 $ilde{p}(y=1|x) = rac{4+1}{5+2} imes rac{1+1}{4+2} imes rac{0+1}{4+2} imes rac{1+1}{4+2} imes rac{3+1}{4+2} imes rac{3+1}{4+2} imes rac{0+1}{4+2} imes rac{1+1}{4+2} pprox .00032$

 $lpha^{\pi}=eta^{\pi}=lpha^{ heta}=eta^{ heta}=1$

$$p(y=0|x)=rac{.00052}{.00032+.00052}pprox .62$$

this dog was my puppy

note that if D is large we have to calculate the product of many terms

d = 1

y = 1

this

0

0

0

0

 $\frac{0}{1}$

 $\frac{1}{4}$

d = 7

 $\frac{1}{1}$

Log-Sum-Exp trick

In estimating unnormalized posteriors we could get numerical problems (underflow) when calculating the posterior for new instances, we work with in the **log-domain**:

$$\log ilde{p}(y|x;\pi, heta) = \log p(y;\pi) + \sum_d \log p(x_d|y; heta_d)$$

to get the final probabilities we need to normalize \tilde{p} we can do this **normalization in the log domain** as well:

$$\log p(y|x;\pi, heta) = \log ilde p(y|x;\pi, heta) - \log \sum_{c=1}^C \exp(\log ilde p(c|x;\pi, heta))$$

we could run into very large or small numbers inside the exponential

we can do this **normalization in the log domain** as well:

$$\log p(y|x;\pi, heta) = \log ilde p(y|x;\pi, heta) - \log \sum_{c=1}^C \exp(\log ilde p(c|x;\pi, heta))$$

observation $\log \sum_c \exp a_c = \log \left(\exp(a_0) \left(\sum_c \exp(a_c - a_0) \right) = a_0 + \log \sum_c \exp(a_c - a_0) \right)$

to make log-sum-exp numerically stable, bring the numbers $~a_c~$ close to zero

for example choose $a_0 \leftarrow \max_c a_c$

Multinomial likelihood

what if we wanted to use word frequencies in document classification?

 $x_d^{(n)}$ is the number of times word d appears in document n



$$p(x|y) = ext{Mult}(x; heta_y) = rac{(\sum_d x_d)!}{\prod_{d=1}^D x_d!} \prod_{d=1}^D heta_{d,y}^{x_d}$$
 probability of word d appearing x_d time

the max-likelihood estimate is again given by the relative frequency

$$heta_{d,c}^{MLE} = rac{\sum_n x_d^{(n)} \mathbb{I}(y^{(n)} = c)}{\sum_n \sum_{d'} x_{d'}^{(n)} \mathbb{I}(y^{(n)} = c)}$$

counts of word d in all documents labelled ${\boldsymbol{\mathsf{c}}}$

total word count in all documents labelled c

	it	is	рирру	cat	pen	а	this
it is a puppy	1	1	1	0	0	1	0
it is a kitten	1	1	0	0	0	1	0
it is a cat	1	1	0	1	0	1	0
that is a dog and this is a pen	0	(2)	0	0	1	(2)	1
it is a matrix	1	1	0	0	0	1	0

Univariate Gaussian density

Gaussian probability density function (pdf)

$$\mathcal{N}(x;\mu,\sigma^2) = rac{1}{\sqrt{2\pi\sigma^2}} e^{-rac{(x-\mu)^2}{2\sigma^2}}$$

two parameters are μ, σ^2

turn out to be the mean and variance

 $\mathbb{E}[x] = \mu$ $\mathbb{E}[(x-\mu)^2] = \sigma^2$

this is a random variable; we are using the same notation for a random variable and a particular value of that variable



Univariate Gaussian density

Gaussian probability density function (pdf)

$$\mathcal{N}(x;\mu,\sigma) = rac{1}{\sqrt{2\pi\sigma^2}}e^{-rac{(x-\mu)^2}{2\sigma^2}}$$

given a dataset $\mathcal{D} = \{x^{(1)}, \dots, x^{(N)}\}$

maximum likelihood estimate of μ, σ^2 are empirical mean and variance

$$egin{aligned} \mu^{MLE} &= rac{1}{N}\sum_n x^{(n)} \ & \sigma^{2^{MLE}} &= rac{1}{N}\sum_n (x^{(n)}-\mu^{MLE})^2 \end{aligned}$$





Univariate Gaussian density

two reasons why Gaussian is an important dist.

- maximum entropy dist. with a fixed variance
- central limit theorem



the average (and sum) of IID random variables has a Gaussian distribution

justifies use of Gaussian for observations that are mean or sum of some random values



Gaussian Naive Bayes

for continuous features one option is the Gaussian conditional likelihood

$$p(x_d \mid y) = \mathcal{N}(x_d; \mu_{d,y}, \sigma_{d,y}^2) = rac{1}{\sqrt{2\pi\sigma_{d,y}^2}} e^{-rac{(x_d - \mu_{d,y})^2}{2\sigma_{d,y}^2}}$$

Maximum likelihood estimates:

empirical mean & variance of feature $\,x_d\,$ across instances with label $\,y\,$

$$egin{aligned} \mu_{d,c} &= rac{1}{N(y=c)} \sum_{n=1}^N x_d^{(n)} \mathbb{I}(y^{(n)}=c) \ \sigma_{d,c}^2 &= rac{1}{N(y=c)} \sum_{n=1}^N \mathbb{I}(y^{(n)}=c) (x_d^{(n)}-\mu_{d,y})^2 \end{aligned}$$

Gaussian Naive Bayes



classification on Iris flowers dataset:

- we use categorical class prior (3 classes)
- Gaussian likelihood since the features are continuous (we use D=2 features)

the **decision boundary** found by Gaussian NB three means are identified using X

- note that we have a mean $\mu_{d,c}$ and variance $\sigma_{d,c}^2$ for each class-feature combination
- in the plot each X is showing the *combined* mean of two features, sepal length and sepal width.



Summary

- generative classification:
 - Iearn the class prior and likelihood
 - Bayes rule for conditional class probability
- Naive Bayes
 - assumes conditional independence

• e.g., word appearances indep. of each other given document type

- class prior: Bernoulli or Categorical
- Iikelihood: Bernoulli, Gaussian, Multinomial...
- MLE has closed-form, estimated separately for each feature and each label
- Bayesian Naive Bayes helps with overfitting
 - with frequent or rare feature values