



DS 310 Probabilistic Models

Vasant G. Honavar

Dorothy Foehr Huck and J. Lloyd Huck Chair in Biomedical Data Sciences and Artificial Intelligence
Professor of Data Sciences, Informatics, Computer Science and Engineering, Bioinformatics & Genomics,
Public Health Sciences and Neuroscience
Director, Center for Artificial Intelligence Foundations and Scientific Applications
Associate Director, Institute for Computational and Data Sciences
Pennsylvania State University

vhonavar@psu.edu
<http://faculty.ist.psu.edu/vhonavar>
<http://ailab.ist.psu.edu>

Bayesian recipe for classification

Note that
$$P(\omega_i | \mathbf{x}) = \frac{P(\mathbf{x} | \omega_i)P(\omega_i)}{P(\mathbf{x})}$$

Model $P(\mathbf{x} | \omega_1)$, $P(\mathbf{x} | \omega_2)$, $P(\omega_1)$, and $P(\omega_2)$

Using Bayes rule, choose ω_1 if $P(\mathbf{x} | \omega_1)P(\omega_1) > P(\mathbf{x} | \omega_2)P(\omega_2)$

Otherwise choose ω_2

Multiple disjoint classes

$$\text{Estimate } P(\omega_i|X) = \frac{P(X|\omega_i)P(\omega_i)}{P(X)}$$

$$\omega = \text{argmax } P(\omega_i|X)$$

Assign sample to the most probable class!

Summary of Bayesian recipe for classification

- The Bayesian recipe is simple, optimal, and in principle, straightforward to apply
- To use this recipe in practice, we need to know
 - $P(X|\omega_i)$, the **generative model for data** for each class and $P(\omega_i)$, the **prior probabilities of classes**
 - **Because these probabilities are unknown, we need to estimate them from data – or learn them!**
- X is typically high-dimensional or may have complex structure
- Need to estimate $P(X|\omega_i)$ from data

Full joint 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	*
			...		

Problem:

- all possible combination of features
- ~10,000 binary features
- Sample space size: 2^{10000}

Full joint 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 probability 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
- The assumptions simplify the modeling task



to the rescue

Variables are **independent** if knowing the value of one
g about the value of the other

ndent variables, knowing the value of one does not
bability distribution of the other variable (or the
ny individual event)

the toss of a coin is independent of a roll of a die
tea in England is independent of the whether or not



Independent or dependent?

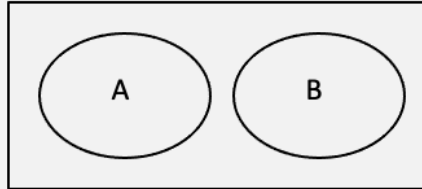
Age and being allergic to cats

Education and driving habits

Success as a mathematician

Success as a basketball player

Independent random variables



If A and B are independent (written $A \perp\!\!\!\perp B$)

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

How does independence help us?

How does independence help us?

If A and B are independent

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

Independence

- Reduces the space needed to store distributions
- Reduces the complexity of the distribution
- Reduces the number of probabilities we need to estimate



Independence

Variables can become independent given certain other

weight

weight given genetics

Conditionally independent given C

$$= P(A|C)P(B|C)$$

$$= P(A|C)$$

$$= P(B|C)$$

$$\neq P(A)P(B)$$

Conditional Independence

- X is **conditionally independent** of Y given Z if the probability distribution governing X is independent of the value of Y given the value of Z :
- $P(X | Y, Z) = P(X | Z)$ that is

$$(\forall x_i, y_j, z_k) P(X = x_i | Y = y_j, Z = z_k) = P(X = x_i | Z = z_k)$$

Quick proof that independence is symmetric

- Assume: $P(X|Y, Z) = P(X|Y)$
- X and Z are independent given Y

$$P(Z | X, Y) = \frac{P(X, Y | Z)P(Z)}{P(X, Y)} \quad (\text{Bayes's Rule})$$

$$= \frac{P(Y | Z)P(X | Y, Z)P(Z)}{P(X | Y)P(Y)} \quad (\text{Chain Rule})$$

$$= \frac{P(Y | Z)P(X | Y)P(Z)}{P(X | Y)P(Y)} \quad (\text{By Assumption})$$

$$= \frac{P(Y | Z)P(Z)}{P(Y)} = P(Z | Y) \quad (\text{Bayes's Rule})$$

Naïve Bayes assumption

$$P(\text{features}, \text{label}) = p(y) \prod_{i=1}^m p(x_i | y, x_1 \cdots x_{i-1})$$

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

What does this mean?

Naïve Bayes assumption

$$P(\text{features}, \text{label}) = p(y) \prod_{i=1}^m p(x_i | y, x_1 \cdots x_{i-1})$$

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

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

Naïve Bayes assumption

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

- Assumes ith feature is independent of the 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 the word “automatic” occurring is independent of whether or not “petrol” occurs given that the review is about “car”

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 “San” occurs will probably make it *more likely* that “Francisco” occurs (or take a compound phrase like “Penn State”)
- 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

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

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 count features
- for real valued features

Naïve Bayes Classifier

- How to learn $P(X|\omega_i)$?
- **Naïve Bayes solution:** Assume that the random variables in X are conditionally independent given the class.
- **Result: Naïve Bayes classifier which performs optimally under certain assumptions**
- A simple, practical learning algorithm grounded in Probability Theory

When to use

- Attributes that describe instances are likely to be conditionally independent given classification
- The data is insufficient to estimate all the probabilities reliably if we do not assume independence

Implications of Independence

- Suppose we have 5 Binary attributes and a binary class label
- Without independence, in order to specify the joint distribution, we need to specify a probability for each possible assignment of values to each variable resulting in a table of size $2^6=64$
- Suppose the features are independent given the class label – we only need $5(2 \times 2)=20$ entries
- The reduction in the number of probabilities to be estimated is even more striking when N , the number of attributes is large – from $O(2^N)$ to $O(N)$

Naive Bayes Classifier

Consider a discrete valued target function $f : \chi \rightarrow \Omega$
where an instance $X = (X_1, X_2, \dots, X_n) \in \chi$ is described
in terms of attribute values $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$
where $x_i \in \text{Domain}(X_i)$

$$\begin{aligned} \omega_{MAP} &= \arg \max_{\omega_j \in \Omega} P(\omega_j | X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) \\ &= \arg \max_{\omega_j \in \Omega} \frac{P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | \omega_j) P(\omega_j)}{P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)} \\ &= \arg \max_{\omega_j \in \Omega} P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | \omega_j) P(\omega_j) \end{aligned}$$

ω_{MAP} is called the *maximum a posteriori* classification

Naive Bayes Classifier

$$\begin{aligned}\omega_{MAP} &= \arg \max_{\omega_j \in \Omega} P(\omega_j | X_1 = x_1, X_2 = x_2 \dots X_n = x_n) \\ &= \arg \max_{\omega_j \in \Omega} P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | \omega_j) P(\omega_j)\end{aligned}$$

If the attributes are *independent* given the class, we have

$$\begin{aligned}\omega_{MAP} &= \arg \max_{\omega_j \in \Omega} \prod_{i=1}^n P(X_i = x_i | \omega_j) P(\omega_j) \\ &= \omega_{NB} \\ &= \arg \max_{\omega_j \in \Omega} P(\omega_j) \prod_{i=1}^n P(X_i = x_i | \omega_j)\end{aligned}$$

Naive Bayes Learner

For each possible value ω_j of Ω ,

$$\hat{P}(\Omega = \omega_j) \leftarrow \text{Estimate}(P(\Omega = \omega_j), D)$$

For each possible value a_{i_k} of X_i

$$\hat{P}(X_i = a_{i_k} | \omega_j) \leftarrow \text{Estimate}(P(X_i = a_{i_k} | \Omega = \omega_j), D)$$

Classify a new instance $X = (x_1, x_2, \dots, x_N)$

$$c(X) = \arg \max_{\omega_j \in \Omega} P(\omega_j) \prod_{i=1}^n P(X_i = x_i | \omega_j)$$

Estimate is a procedure for estimating the relevant probabilities from set of training examples

Learning Dating Preferences

Data samples – ordered 3-tuples of attribute values corresponding to

Height (tall, short)

Hair (dark, blonde, red)

Eye (blue, brown)

Classes: +, –

	Training Data	
	Instance	Class label
	I ₁	(t, d, l) +
	I ₂	(s, d, l) +
	I ₃	(t, b, l) –
	I ₄	(t, r, l) –
	I ₅	(s, b, l) –
	I ₆	(t, b, w) +
	I ₇	(t, d, w) +
	I ₈	(s, b, w) +

Probabilities to estimate

$$P(+)=5/8$$

$$P(-)=3/8$$

$P(\text{Height} c)$	t	s
+	3/5	2/5
-	2/3	1/3

$P(\text{Hair} c)$	d	b	r
+	3/5	2/5	0
-	0	2/3	1/3

$P(\text{Eye} c)$	l	w
+	2/5	3/5
-	1	0

Classify ($\text{Height}=t, \text{Hair}=b, \text{eye}=l$)

$$P(X|+) = (3/5)(2/5)(2/5) = (12/125)$$

$$P(X|-) = (2/3)(2/3)(1) = (4/9)$$

$$P(+|X) \propto P(+P(X|+)) = (5/8)(12/125) = 0.06$$

$$P(-|X) \propto P(-)P(X|-) = (3/8)(4/9) = 0.1667$$

Classify ($\text{Height}=t, \text{Hair}=r, \text{eye}=w$)

- Note the problem when a probability is 0
- Solution – Use Dirichlet Priors

$$P(x_i=a_{i_k}) = \frac{\alpha_{i_k} + N_{i_k}}{\sum_k \alpha_{i_k} + \sum_k N_{i_k}}$$

PennState Institute for Computational and Data Sciences
Center for Artificial Intelligence Foundations & Scientific Applications
Artificial Intelligence Research Laboratory
PennState Clinical and Translational Science Institute

Modeling $p(x|y)$

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:

Model using an appropriate distribution:

- Continuous values - Gaussian (i.e. normal) distribution
- Counts - Multinomial distribution (more on this later)
- ...

PennState Office of Information Science and Technology
Fall 2022
Vasant G Honavar

- for discrete, we could simply do a much larger table, but often that doesn't capture everything we want

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?

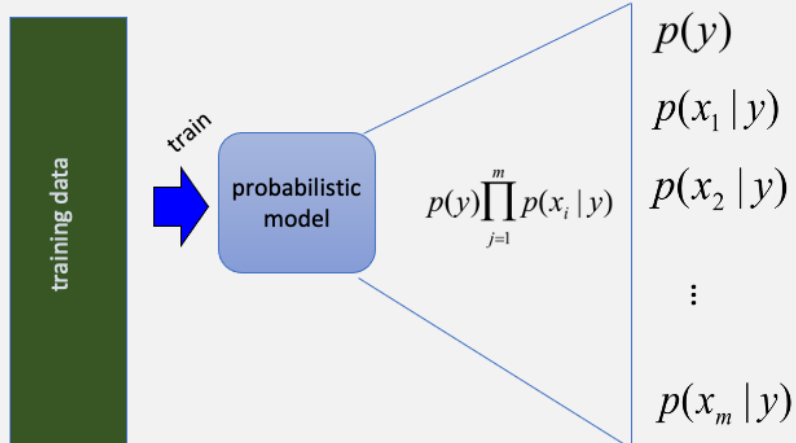
Obtaining probabilities



We've talked a lot about probabilities, but not where they come from

- How do we estimate $p(x_i|y)$ from training data?
- What is the probability of surviving the titanic?
- What is the probability that a review is about Toyota or BMW?

Obtaining probabilities



Estimating probabilities

What is the probability of a Macbook Air review?

We don't know!

We can *estimate* it based on data, though:

number of reviews labeled Macbook Air

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. one characterized by parameters, say θ) assigns to the data

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

for each example

the model parameters (e.g. probability of heads)

how probable is it under the model

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

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

for each example

the model parameters (e.g. probability of heads)

how probable is it under the model

Likelihood

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

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

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

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

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

60 heads with $p(\text{head}) = 0.6$

40 tails with $p(\text{tail}) = 0.4$

MLE example

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

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

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?

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

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

60 heads with $p(\text{head}) = 0.6$

40 tails with $p(\text{tail}) = 0.4$

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

60 heads with $p(\text{head}) = 0.5$

40 tails with $p(\text{tail}) = 0.5$

MLE example

Can we do any better?

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

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

60 heads with $p(\text{head}) = 0.6$

40 tails with $p(\text{tail}) = 0.4$

What about $p(\text{head}) = 0.7$?

MLE example

Can we do any better?

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

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

60 heads with $p(\text{head}) = 0.6$

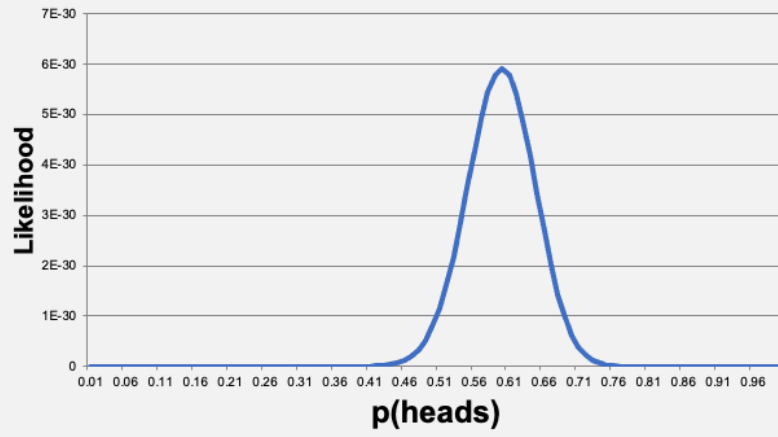
40 tails with $p(\text{tail}) = 0.4$

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

60 heads with $p(\text{head}) = 0.7$

40 tails with $p(\text{tail}) = 0.3$

MLE Example



PennState Institute for Computational and Data Sciences
Center for Artificial Intelligence Foundations & Scientific Applications
Artificial Intelligence Research Laboratory
PennState Clinical and Translational Science Institute

Maximum Likelihood Estimation (MLE)

The **maximum likelihood** estimate for a model parameter is the one that maximize 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:

$$\begin{aligned} MLE &= \arg \max_{\theta} \log\left(\prod_{i=1}^n p_{\theta}(x_i)\right) \\ &= \arg \max_{\theta} \sum_{i=1}^n \log(p(x_i)) \end{aligned}$$

Why is this ok?

PennState Office of Information Science and Technology
Fall 2022
Vasant G Honavar

- log is a strictly increasing function
- it just squishes values but does not change their order, so the max of likelihood is still the max of log-likelihood

Calculating MLE

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

$$MLE = \operatorname{argmax}_{\theta} \sum_{i=1}^n \log(p(x_i))$$

Given some training data, how do we calculate the MLE?

Training data: 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.

$$\begin{aligned}\log\text{-likelihood} &= \sum_{i=1}^n \log(p(x_i)) \\ &= 60 \log(p(\text{heads})) + 40 \log(p(\text{tails})) \\ &= 60 \log(\theta) + 40 \log(1 - \theta)\end{aligned}$$

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

How do we find the max?

Calculating MLE

- How do you find the maximum of a function of θ ?

$$\frac{d}{d\theta} 60 \log(\theta) + 40 \log(1-\theta) = 0$$

$$\frac{60}{\theta} - \frac{40}{1-\theta} = 0$$

$$\frac{40}{1-\theta} = \frac{60}{\theta}$$

$$40\theta = 60 - 60\theta$$

$$100\theta = 60$$

$$\theta = \frac{60}{100}$$

Calculating MLE

You flip a coin n times.

a times you get heads and b times you get tails.

$$\frac{d}{d\theta} a \log(\theta) + b \log(1-\theta) = 0$$

• • •

$$\theta = \frac{a}{a+b}$$

MLE estimation for Naïve Bayes

training data

train

probabilistic
model

$$p(y) \prod_i^n p(x_i|y)$$

$p(y)$ $p(x_i | y)$

What are the MLE estimates
for these?

PennState Institute for Computational and Data Sciences
Center for Artificial Intelligence Foundations & Scientific Applications
Artificial Intelligence Research Laboratory
PennState Clinical and Translational Science Institute

Maximum likelihood estimates

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

number of examples with label
total number of examples

$$p(x_i | y) = \frac{\text{count}(x_i, y)}{\text{count}(y)}$$

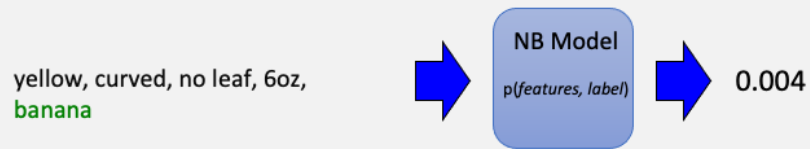
number of examples with the label with feature
number of examples with label

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

PennState Office of Information Science and Technology
Fall 2022
Vasant G Honavar

- just involves iterating over the data and aggregating these counts!

Naïve Bayes classification



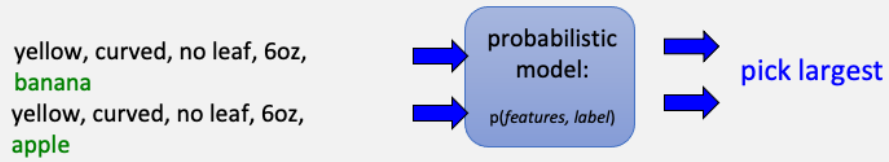
$$p(y) \prod_i^n p(x_i|y)$$

Given an unlabeled example:

predict the label

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

Probabilistic models




$$p(y) \prod_i^n p(x_i|y)$$

$$\text{label} = \arg \max p(y) \prod_i^n p(x_i|y)$$

PennState Institute for Computational and Data Sciences
Center for Artificial Intelligence Foundations & Scientific Applications
Artificial Intelligence Research Laboratory
PennState Clinical and Translational Science Institute

Generative Story



- Given model and a sample, e.g., a document
 - We can calculate the probability of the sample being generated under the model
 - We can also ask how a given model would **generate** a document - “**generative story**” for a model
- The generative story can
 - help you understand the model
 - help you develop and refine a model

PennState Office of Research and Innovation
Fall 2022
Vasant G Honavar

- although we don't generally "generate" a document from a model, it's often useful to look at the generative story of a model (i.e. how the model says a document was generate) to help us understand why the model assigns certain probabilities

Naïve Bayes generative story



$$p(\omega) \prod_i^n P(x_i | \omega)$$

What is the generative story for the NB model?

NB generative story

$$p(\omega) \prod_i^n P(x_i | \omega)$$



1. Pick a label according to $p(\omega)$
 - Roll a biased, m sided die
 - Set the class label that shows up
2. Given the class, for each of the n features feature:
 - Roll a biased die with as many sides as possible values of the feature.
 - Set the feature to the value that shows up