

A crash course in probability and Naïve Bayes classification

Chapter 9

1

Probability theory

Random variable: a variable whose possible values are numerical outcomes of a random phenomenon.

Examples: A person's height, the outcome of a coin toss

Distinguish between discrete and continuous variables.

The distribution of a discrete random variable:

The probabilities of each value it can take.

Notation: $P(X = x_i)$.

These numbers satisfy:

$$\sum_i P(X = x_i) = 1$$

2

Probability theory

		n_{ij}	

Marginal Probability

$$p(X = x_i) = \frac{c_i}{N}$$

Joint Probability

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

Conditional Probability

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

3

Probability theory

A joint probability distribution for two variables is a table.

		n_{ij}	

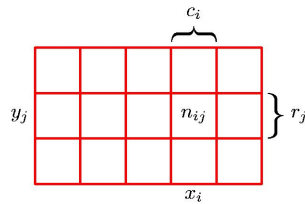
If the two variables are binary, how many parameters does it have?

What about joint probability of d variables $P(X_1, \dots, X_d)$?

How many parameters does it have if each variable is binary?

4

Probability theory



Marginalization:

$$p(X = x_i) = \frac{c_i}{N} = \frac{1}{N} \sum_{j=1}^L n_{ij}$$

$$= \sum_{j=1}^L p(X = x_i, Y = y_j)$$

Product Rule

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N}$$

$$= p(Y = y_j | X = x_i) p(X = x_i)$$

The Rules of Probability

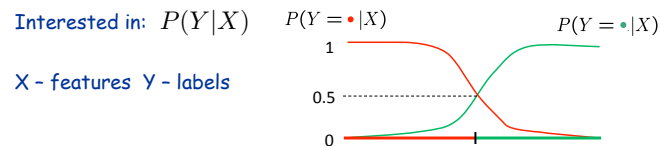
Marginalization $p(X) = \sum_Y p(X, Y)$

Product Rule $p(X, Y) = p(Y|X)p(X)$

Independence: X and Y are independent if $P(Y|X) = P(Y)$

This implies $P(X, Y) = P(X) P(Y)$

Using probability in learning



For example, when classifying spam, we could estimate $P(Y | \text{Viagra, lottery})$

We would then classify an example if $P(Y|X) > 0.5$.

However, it's usually easier to model $P(X | Y)$

Maximum likelihood

Fit a probabilistic model $P(x | \theta)$ to data

- Estimate θ

Given independent identically distributed (i.i.d.) data

$\mathbf{X} = (x_1, x_2, \dots, x_n)$

- Likelihood

$$P(\mathbf{X}|\theta) = P(x_1|\theta)P(x_2|\theta), \dots, P(x_n|\theta)$$

- Log likelihood

$$\ln P(\mathbf{X}|\theta) = \sum_{i=1}^n \ln P(x_i|\theta)$$

Maximum likelihood solution: parameters θ that maximize $\ln P(\mathbf{X} | \theta)$

Example

Example: coin toss

Estimate the probability p that a coin lands "Heads" using the result of n coin tosses, h of which resulted in heads.

The likelihood of the data: $P(\mathbf{X}|\theta) = p^h(1-p)^{n-h}$

Log likelihood: $\ln P(\mathbf{X}|\theta) = h \ln p + (n-h) \ln(1-p)$

Taking a derivative and setting to 0:

$$\frac{\partial \ln P(\mathbf{X}|\theta)}{\partial p} = \frac{h}{p} - \frac{(n-h)}{(1-p)} = 0$$

$$\Rightarrow p = \frac{h}{n}$$

9

Bayes' rule

From the product rule:

$$P(Y, X) = P(Y | X) P(X)$$

and:

$$P(Y, X) = P(X | Y) P(Y)$$

Therefore:

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

This is known as Bayes' rule

10

Bayes' rule

$$P(Y|X) = \frac{\overset{\text{likelihood}}{P(X|Y)} \overset{\text{prior}}{P(Y)}}{\underset{\text{posterior}}{P(X)}}$$

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

$P(X)$ can be computed as:

$$P(X) = \sum_Y P(X|Y)P(Y)$$

But is not important for inferring a label

Maximum a-posteriori and maximum likelihood

The **maximum a posteriori (MAP)** rule:

$$y_{MAP} = \arg \max_Y P(Y|X) = \arg \max_Y \frac{P(X|Y)P(Y)}{P(X)} = \arg \max_Y P(X|Y)P(Y)$$

If we ignore the prior distribution or assume it is uniform we obtain the **maximum likelihood** rule:

$$y_{ML} = \arg \max_Y P(X|Y)$$

A classifier that has access to $P(Y|X)$ is a **Bayes optimal** classifier.

12

Naïve Bayes classifier

We would like to model $P(X | Y)$, where X is a feature vector, and Y is its associated label.

Task: Predict whether or not a picnic spot is enjoyable

Training Data: $X = (X_1 \quad X_2 \quad X_3 \quad \dots \quad X_d) \quad Y$

Sky	Temp	Humid	Wind	Water	Forecast	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

n rows

How many parameters?

Prior: $P(Y)$ $k-1$ if k classes

Likelihood: $P(X | Y)$ $(2^d - 1)k$ for binary features

13

Naïve Bayes classifier

We would like to model $P(X | Y)$, where X is a feature vector, and Y is its associated label.

Simplifying assumption: conditional independence: given the class label the features are independent, i.e.

$$P(\mathbf{X}|Y) = P(x_1|Y)P(x_2|Y), \dots, P(x_d|Y)$$

How many parameters now?

14

Naïve Bayes classifier

We would like to model $P(X | Y)$, where X is a feature vector, and Y is its associated label.

Simplifying assumption: conditional independence: given the class label the features are independent, i.e.

$$P(\mathbf{X}|Y) = P(x_1|Y)P(x_2|Y), \dots, P(x_d|Y)$$

How many parameters now? $dk + k - 1$

15

Naïve Bayes classifier

Naïve Bayes decision rule:

$$y_{NB} = \arg \max_Y P(\mathbf{X}|Y)P(Y) = \arg \max_Y \prod_{i=1}^d P(x_i|Y)P(Y)$$

If conditional independence holds, NB is an optimal classifier!

16

Training a Naïve Bayes classifier

Training data: Feature matrix X ($n \times d$) and labels y_1, \dots, y_n

Maximum likelihood estimates:

Class prior: $\hat{P}(y) = \frac{|\{i : y_i = y\}|}{n}$

Likelihood: $\hat{P}(x_i|y) = \frac{\hat{P}(x_i, y)}{\hat{P}(y)} = \frac{|\{i : X_{ij} = x_i, y_i = y\}|/n}{|\{i : y_i = y\}|/n}$

17

Example

Email classification

Suppose our vocabulary contains three words a , b and c , and we use a multivariate Bernoulli model for our e-mails, with parameters

$$\theta^{\oplus} = (0.5, 0.67, 0.33) \quad \theta^{\ominus} = (0.67, 0.33, 0.33)$$

This means, for example, that the presence of b is twice as likely in spam (+), compared with ham.

The e-mail to be classified contains words a and b but not c , and hence is described by the bit vector $\mathbf{x} = (1, 1, 0)$. We obtain likelihoods

$$P(\mathbf{x}|\oplus) = 0.5 \cdot 0.67 \cdot (1 - 0.33) = 0.222 \quad P(\mathbf{x}|\ominus) = 0.67 \cdot 0.33 \cdot (1 - 0.33) = 0.148$$

The ML classification of \mathbf{x} is thus spam.

18

Example

Email classification: training data

E-mail	$a?$	$b?$	$c?$	Class
e_1	0	1	0	+
e_2	0	1	1	+
e_3	1	0	0	+
e_4	1	1	0	+
e_5	1	1	0	-
e_6	1	0	1	-
e_7	1	0	0	-
e_8	0	0	0	-

What are the parameters of the model?

19

Example

Email classification: training data

E-mail	$a?$	$b?$	$c?$	Class
e_1	0	1	0	+
e_2	0	1	1	+
e_3	1	0	0	+
e_4	1	1	0	+
e_5	1	1	0	-
e_6	1	0	1	-
e_7	1	0	0	-
e_8	0	0	0	-

What are the parameters of the model?

$$\hat{P}(y) = \frac{|\{i : y_i = y\}|}{n}$$

$$\hat{P}(x_i|y) = \frac{\hat{P}(x_i, y)}{\hat{P}(y)} = \frac{|\{i : X_{ij} = x_i, y_i = y\}|/n}{|\{i : y_i = y\}|/n}$$

20

Example

Email classification: training data

E-mail	a?	b?	c?	Class
e ₁	0	1	0	+
e ₂	0	1	1	+
e ₃	1	0	0	+
e ₄	1	1	0	+
e ₅	1	1	0	-
e ₆	1	0	1	-
e ₇	1	0	0	-
e ₈	0	0	0	-

What are the parameters of the model?

$$P(+) = 0.5, P(-) = 0.5$$

$$\hat{P}(y) = \frac{|\{i : y_i = y\}|}{n}$$

$$P(a|+) = 0.5, P(a|-) = 0.75$$

$$P(b|+) = 0.75, P(b|-) = 0.25$$

$$P(c|+) = 0.25, P(c|-) = 0.25$$

$$\hat{P}(x_i|y) = \frac{\hat{P}(x_i, y)}{\hat{P}(y)} = \frac{|\{i : X_{ij} = x_i, y_i = y\}|/n}{|\{i : y_i = y\}|/n}$$

21

Comments on Naïve Bayes

Usually features are not conditionally independent, i.e.

$$P(\mathbf{X}|Y) \neq P(x_1|Y)P(x_2|Y), \dots, P(x_d|Y)$$

And yet, one of the most widely used classifiers. Easy to train!

It often performs well even when the assumption is violated.

Domingos, P., & Elkan, C. (1997). Bayes' Theorem as a Linear Classifier. *Machine Learning*, 29, 103-130.

22

When there are few training examples

What if you never see a training example where $x_1=a$ when $y=\text{spam}$?

$$P(x | \text{spam}) = P(a | \text{spam}) P(b | \text{spam}) P(c | \text{spam}) = 0$$

What to do?

23

When there are few training examples

What if you never see a training example where $x_1=a$ when $y=\text{spam}$?

$$P(x | \text{spam}) = P(a | \text{spam}) P(b | \text{spam}) P(c | \text{spam}) = 0$$

What to do?

Add "virtual" examples for which $x_1=a$ when $y=\text{spam}$.

24

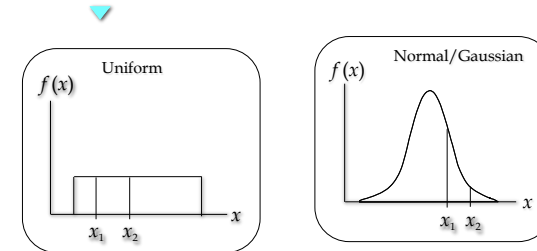
Naïve Bayes for continuous variables

Need to talk about continuous distributions!

25

Continuous Probability Distributions

The probability of the random variable assuming a value within some given interval from x_1 to x_2 is defined to be the area under the graph of the probability density function between x_1 and x_2 .



Expectations

Discrete variables

$$\mathbb{E}[f] = \sum_x p(x) f(x)$$

Continuous variables

$$\mathbb{E}[f] = \int p(x) f(x) dx$$

$$\mathbb{E}_x[f|y] = \sum_x p(x|y) f(x)$$

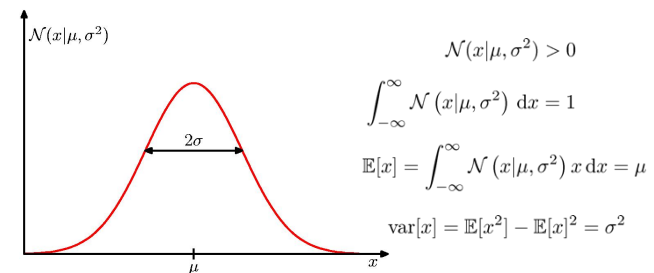
Conditional expectation
(discrete)

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^N f(x_n)$$

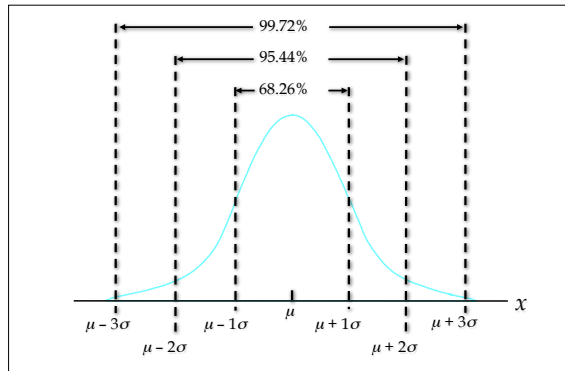
Approximate expectation
(discrete and continuous)

The Gaussian (normal) distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$



Properties of the Gaussian distribution



Standard Normal Distribution

▶ A random variable having a normal distribution with a mean of 0 and a standard deviation of 1 is said to have a standard normal probability distribution.

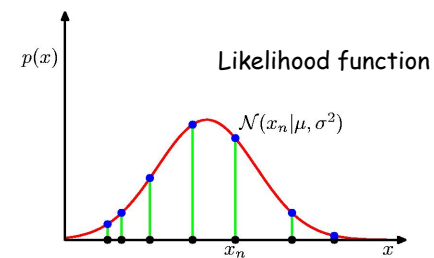
Standard Normal Probability Distribution

- ▶ Converting to the Standard Normal Distribution

$$z = \frac{x - \mu}{\sigma}$$

We can think of z as a measure of the number of standard deviations x is from μ .

Gaussian Parameter Estimation



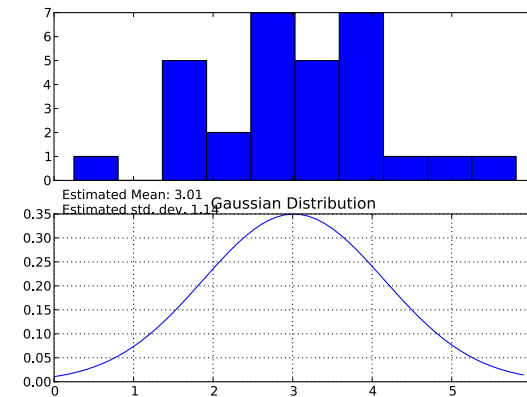
$$p(\mathbf{x}|\mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \sigma^2)$$

Maximum (Log) Likelihood

$$\ln p(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n \quad \sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2$$

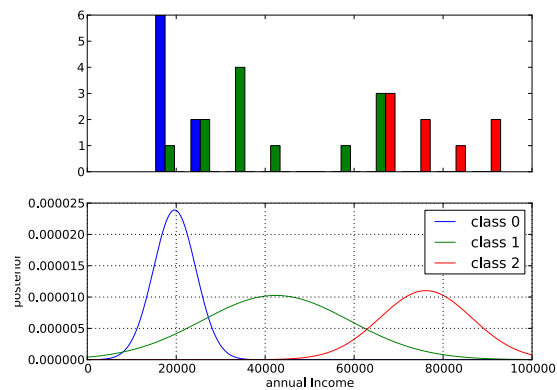
Example



34

Gaussian models

Assume we have data that belongs to three classes, and assume a likelihood that follows a Gaussian distribution



35

Gaussian Naïve Bayes

Likelihood function:

$$P(X_i = x|Y = y_k) = \frac{1}{\sqrt{2\pi}\sigma_{ik}} \exp\left(-\frac{(x - \mu_{ik})^2}{2\sigma_{ik}^2}\right)$$

Need to estimate mean and variance for each feature in each class.

36

Summary

Naïve Bayes classifier:

- ✧ What's the assumption
- ✧ Why we make it
- ✧ How we learn it

Naïve Bayes for discrete data

Gaussian naïve Bayes