

COMP2211 Exploring Artificial Intelligence

Naïve Bayes Classifier

Huiru Xiao

Department of Computer Science and Engineering
The Hong Kong University of Science and Technology

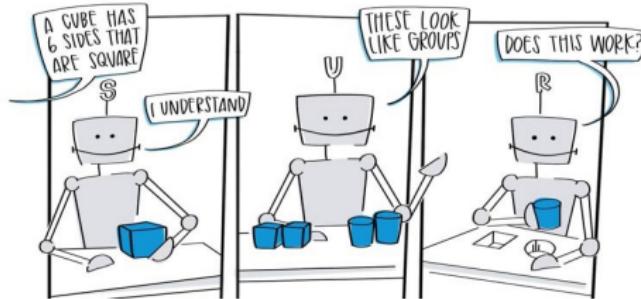
Fundamentals of Machine Learning

What is Machine Learning?

Definition

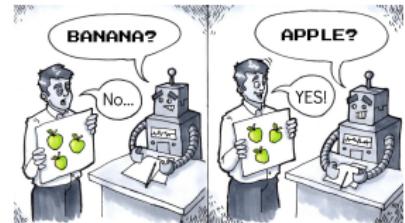
Machine Learning is the science and engineering of getting computers to act without being explicitly programmed.

- Types:
 - Supervised learning
 - Unsupervised learning
 - Reinforcement learning (Will not be covered in this course, only for self-study)



Supervised Learning

- We form a **training dataset** consisting of **examples with inputs and correct outputs** (we called these **labels**).
 - This training dataset is called a **labelled dataset**
- The labelled dataset is then used to train a computer/machine to yield the desired output, i.e., allowing the machine to **learn the relationship/mapping** between the inputs and outputs.
- Supervised learning can be further divided into two types:
 - ① **Classification:** **Classify the data into categories** (also called classes – represented using different labels).
For example, yes or no, positive or negative or neutral, etc.
 - ② **Regression:** **Give a real or continuous value** as output.
For example, salary based on work experience, weight based on height, etc.



An Example Demonstrating the Idea of Supervised Learning

- Suppose we want to train our computer/machine to do maths.
- Instead of explicitly programming it, we pass a set of input data and the corresponding answers to the computer/machine and then ask it to learn/figure out the relationship between the data and the answers.

1.	2	4	5	=	3	5.	6	2	2	=	10
2.	5	2	8	=	2	6.	3	1	1	=	2
3.	2	2	1	=	3	7.	5	3	4	=	11
4.	4	2	2	=	6	8.	1	8	1	=	7

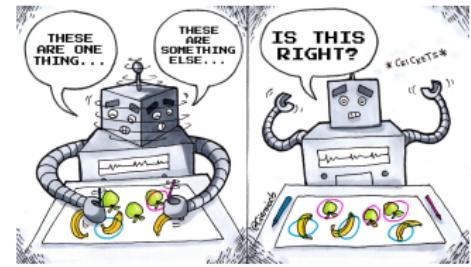
Question: For this example, what is the training dataset? What is label?
Can you figure out the relationship between the data and labels?

Real-Life Applications of Supervised Learning

- Image classification
 - E.g., Facebook can recognize your friend in a picture from an album of tagged photos.
- Visual recognition
 - E.g., The ability of a machine to identify objects, places, people, actions, and images.
- Signature recognition
 - E.g., Automatic signature recognition for automating bank cheques processing.
- Spam detection
 - E.g., Google has now advanced to the point that it can detect and filter out spam and phishing emails with about 99.9% accuracy.
- Weather forecasting
 - E.g., HK Observatory incorporates machine learning in strengthening the techniques in monitoring and forecasting.

Unsupervised Learning

- We use an **unlabelled dataset** (i.e., without giving the correct outputs) and train the machines to **learn on themselves**.
- The computer/machine tries to **find a pattern in the unlabeled data** and respond.
- Unsupervised learning can be further divided into two groups:
 - ① **Clustering**: **Group data into clusters** that are similar between them and are dissimilar to the objects belonging to another cluster. For example, finding out which customers made similar product purchases.
 - ② **Association**: A rule-based machine learning to **discover the probability of the co-occurrence of items** in a collection. For example, finding out which products were purchased together.



An Example Demonstrating the Idea of Unsupervised Learning

- Once again, suppose we want to train our computer/machine to do maths.
- Instead of explicitly programming it, we pass a set of input data, but NOT the corresponding answers, to the computer/machine and then ask it to learn/figure out the patterns.

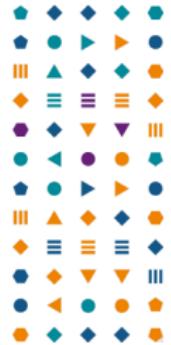
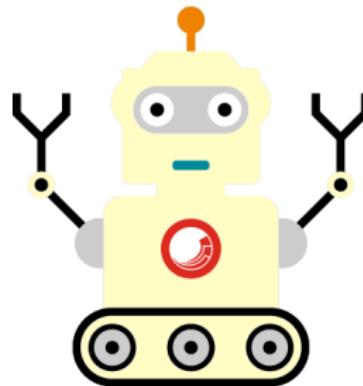
$$\begin{array}{l} 1. \quad 2 \quad 4 \quad 5 \quad = \\ 2. \quad 5 \quad 2 \quad 8 \quad = \\ 3. \quad 2 \quad 2 \quad 1 \quad = \\ 4. \quad 4 \quad 2 \quad 2 \quad = \end{array}$$

$$\begin{array}{l} 5. \quad 6 \quad 2 \quad 2 \quad = \\ 6. \quad 3 \quad 1 \quad 1 \quad = \\ 7. \quad 5 \quad 3 \quad 4 \quad = \\ 8. \quad 1 \quad 8 \quad 1 \quad = \end{array}$$

Question: For this example, is it a suitable task for the machine to learn the patterns?

Real-Life Applications of Unsupervised Learning

- Semantic Clustering
 - Semantically similar words share a similar context
 - Grouping images/videos based on similarities
 - Grouping audios based on acoustic similarity
- Identifying Accident Prone Areas
 - Used to identify accident-prone areas and introduce safety measures based on the intensity of those accidents.



Comparison between Supervised and Unsupervised Learning

	Supervised Learning	Unsupervised Learning
Training data	Labeled data Need domain expert to label data	Unlabeled data
Categories	Classification and regression	Clustering and association
Preference	Calculate outcomes	Discover underlying patterns
Training time	Shorter	Longer
Accuracy	High	Low
Optimal strategy	Depend on the data and learning algorithm	Depend on the data
Algorithms	Naïve Bayes, K-Nearest Neighbour, Decision Tree, Linear Regression, Logistics Regression, Support Vector Machine	K-Means Clustering, Hierarchical Clustering, Apriori Algorithm

Problem

Suppose we have a dataset consisting of symptoms and the results showing the patients suffering from disease 'Z' or not.

Data	Blood Pressure	Fever	Diabetes	Vomit	Suffering from disease Z
1	High	High	Yes	No	No
2	High	High	Yes	Yes	No
3	Low	High	Yes	No	Yes
4	Normal	Mild	Yes	No	Yes
5	Normal	No fever	No	No	Yes
6	Normal	No fever	No	Yes	No
7	Low	No fever	No	Yes	Yes
8	High	Mild	Yes	No	No
9	High	No fever	No	No	Yes
10	Normal	Mild	No	No	Yes
11	High	Mild	No	Yes	Yes
12	Low	Mild	Yes	Yes	Yes
13	Low	High	No	No	Yes
14	Normal	Mild	Yes	Yes	No



Can we estimate whether a person suffering from disease 'Z' based on the symptoms (Blood Pressure = high, Fever = no, Diabetes = yes, and Vomit = yes) from the above table? Yes, use Bayes' rule.

Bayes' Rule

Conditional Probability

- Conditional probability denoted by $P(B|A)$ means what the probability that an event B occurs, given that event A has already occurred.
- $P(B|A) = \frac{P(A \& B)}{P(A)}$
- Example:

Transport	Car	Foot	Bicycle	Total
Undergraduate Students	200	330	120	650
Postgraduate Students	250	320	80	650
Total	450	650	200	1300

- What is the probability $P(\text{Car}|\text{Undergraduate Students})$?
Answer: $200/650 = 0.31$
- What is the probability $P(\text{Bicycle}|\text{Postgraduate Students})$?
Answer: $80/650 = 0.12$



Bayes' Rule

Bayes' rule is named after 18th-century British mathematician Thomas Bayes. It is the **most important rule in data science**. It is the mathematical rule that describes how to update a belief, given some evidence.

$$P(B|E) = P(B) \times \frac{P(E|B)}{P(E)}$$

- E: Evidence (also called Feature)
- B: Belief (also called Class)
- $P(B|E)$: Posterior probability (How often belief happens given evidence happens)
- $P(B)$: Prior probability (How likely belief is on its own)
- $P(E|B)$: Likelihood (How often evidence happens given that belief happens)
- $P(E)$: Marginal probability (How likely evidence is on its own)



Example

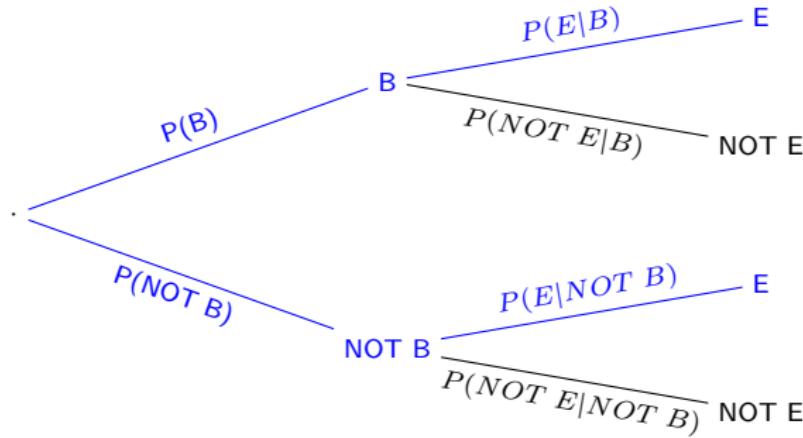
- E: Smoke, B: Fire
- $P(B)$: how often there is fire
- $P(E)$: how often there is smoke
- $P(B|E)$: how often there is fire when we can see smoke
- $P(E|B)$: how often there is smoke when there is fire
- So, Bayes' rule is a kind of formula telling us $P(\text{Fire}|\text{Smoke})$ when we know $P(\text{Smoke}|\text{Fire})$
- Dangerous fires are rare (1%): $P(B) = P(\text{Fire}) = 0.01$
- But smoke is fairly common (10%) due to BBQ: $P(E) = P(\text{Smoke}) = 0.1$
- And 90% of dangerous fires make smoke: $P(E|B) = P(\text{Smoke}|\text{Fire}) = 0.9$

Probability of dangerous Fire when there is smoke

$$P(B|E) = P(\text{Fire}|\text{Smoke}) = \frac{P(\text{Fire})P(\text{Smoke}|\text{Fire})}{P(\text{Smoke})} = \frac{(0.01)(0.9)}{0.1} = 0.09$$

How to Understand Bayes' Rule?

- We can represent the relationship between E and B as a tree.



$$P(E) = P(E|B)P(B) + P(E|Not B)P(Not B)$$

- It is okay if we know $P(E|B)$, $P(B)$, $P(E|Not B)$, $P(Not B)$. So, we have

$$P(B|E) = \frac{P(B)P(E|B)}{P(E|B)P(B) + P(E|Not B)P(Not B)} = \frac{P(B)P(E|B)}{P(E)}$$

Another Example: Virus Detection

- A doctor gives a patient a test for a particular virus.
- It is believed that 0.1% of the population has this virus.
- Based on experience, the doctor knows that, in 99% of the cases in which the virus is present, the test is positive; and in 95% of the cases in which the virus is not present, the test is negative.
- Question: If the test turns out to be positive, what is the chance that the virus is present?

Can we apply Bayes' rule to solve this problem? Let's try.

- Let V be the part of the population carrying the virus
- Let H be the remaining part

$$P(V) = 0.001, P(H) = 1 - 0.001 = 0.999$$

$$P(+|V) = 0.99, P(-|H) = 0.95$$

$$P(V|+) = ?$$

Another Example: Virus Detection

- According to Bayes' rule

$$P(V|+) = \frac{P(V)P(+|V)}{P(+)}$$

What is $P(+)$?

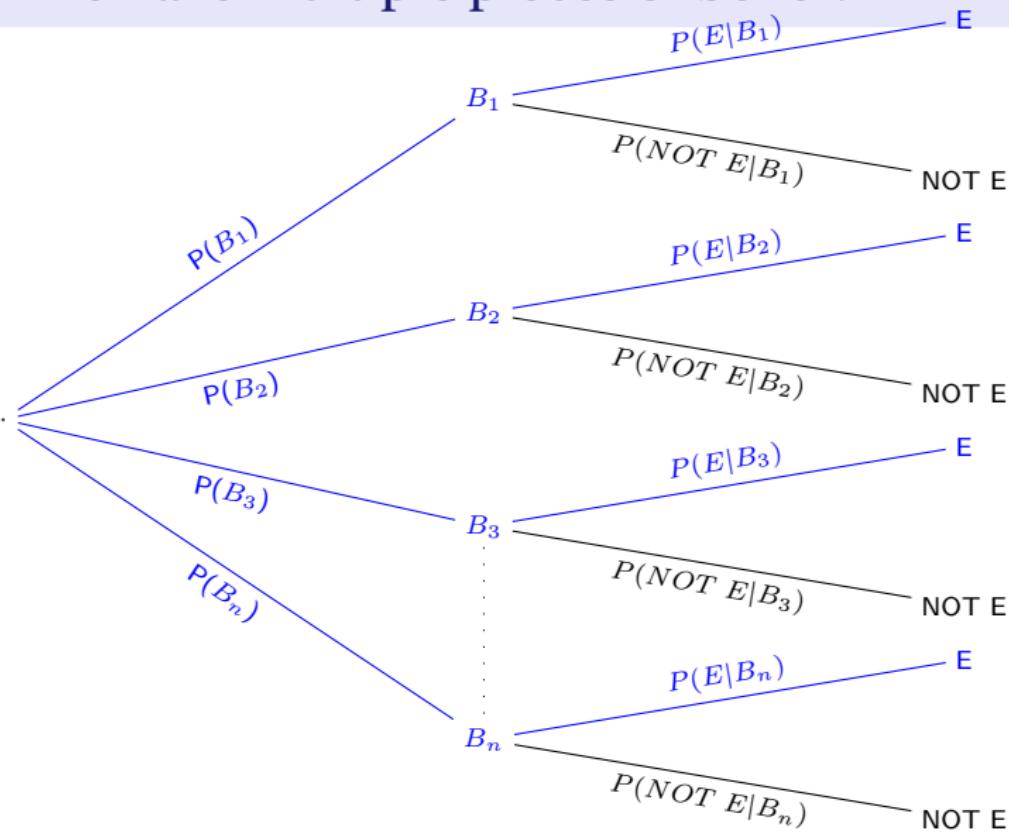
- In fact,

$$\begin{aligned} P(+) &= P(+|V)P(V) + P(+|H)P(H) \\ &= (0.99)(0.001) + (1 - 0.95)(0.999) \\ &= (0.99)(0.001) + (0.05)(0.999) \\ &= 0.05094 \end{aligned}$$

- So,

$$\begin{aligned} P(V|+) &= \frac{(0.001)(0.99)}{0.05094} \\ &\approx 0.02 \end{aligned}$$

What if we have multiple pieces of belief?



General Form of Bayes' Rule

$$\begin{aligned} P(B_i|E) &= \frac{P(B_i)P(E|B_i)}{P(E|B_1)P(B_1) + P(E|B_2)P(B_2) + \cdots + P(E|B_n)P(B_n)} \\ &= \frac{P(B_i)P(E|B_i)}{\sum_{j=1}^n P(E|B_j)P(B_j)} \end{aligned}$$

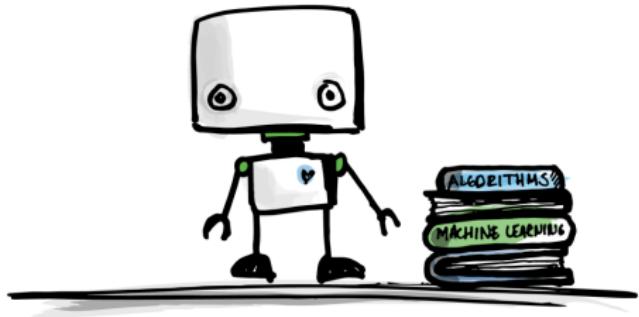
- Each belief, B_i , is also called a class. For instance, V (population carrying the virus) and H (the population that does not carry the virus) are two different classes in the virus detection example.



What if we have multiple pieces of evidence also?

- Say, we have d pieces of evidence $E = \{e_1, e_2, e_3, \dots, e_d\}$

$$\begin{aligned} P(B_i|E) &= \frac{P(B_i)P(E|B_i)}{P(E|B_1)P(B_1) + P(E|B_2)P(B_2) + \dots + P(E|B_n)P(B_n)} \\ &= \frac{P(B_i)P((e_1, e_2, e_3, \dots, e_d)|B_i)}{\sum_{j=1}^n P((e_1, e_2, e_3, \dots, e_d)|B_j)P(B_j)} \end{aligned}$$



Naïve Bayes Model

Naïve Bayes

- The Bayes' rule is referred to as "Naive Bayes" when we assume that each piece of evidence conditional independently and equally contributes to the probability of a sample belonging to a certain belief. In other words, the presence of one particular evidence in a belief does not affect the presence of another one.
 - Conditional independent:
For example, given a certain outcome (e.g., going for a walk), the temperature being "Hot" has nothing to do with the humidity being "high" or outlook being "Rainy" does not affect the winds. Hence the features are assumed to be conditional independent.
 - Equal contribution:
For example, knowing only temperature and humidity alone cannot predict the outcome accurately. Each evidence contributes equally to the outcome.

The assumption made by Naïve Bayes are NOT generally correct in real-world situations. In fact, the independence assumption is never correct but often works well in practice.

Naïve Bayes

- With each evidence makes an independent and equal contribution given the belief, Bayes' rule can be transformed to the following.

$$\begin{aligned} P(B_i|E) &= \frac{P(B_i)P((e_1, e_2, e_3, \dots, e_d)|B_i)}{\sum_{j=1}^n P((e_1, e_2, e_3, \dots, e_d)|B_j)P(B_j)} \\ &= \frac{P(B_i)P(e_1|B_i)P(e_2|B_i)P(e_3|B_i)\dots P(e_d|B_i)}{\sum_{j=1}^n P(e_1|B_j)P(e_2|B_j)P(e_3|B_j)\dots P(e_d|B_j)P(B_j)} \end{aligned}$$

Question: We are now given a test data sample with evidence E , we want to know which belief (class) it belongs to, so we calculate $P(B_1|E), P(B_2|E), \dots, P(B_n|E)$.

When calculating these posterior probabilities, we follow the above Naïve Bayes formula.

However, do we need to strictly follow the formula and the actual probability values?

Naïve Bayes

- Suppose we further assume that we are only interested in knowing which belief has the highest probability value rather than the actual probability value. In that case, we then we can further simplify the equation by removing the denominator as follows:

$$B_{NB} = \arg \max_{B_i} P(B_i)(P(e_1|B_i)P(e_2|B_i)P(e_3|B_i)\cdots P(e_d|B_i))$$

- It is because the denominator remains constant for a given input.

Now consider a feature e_j , suppose e_j has m possible values:
 $e_j \in \{v_1, v_2, \dots, v_m\}$. Then the likelihood

$$P(e_j = v_k | B_i) = \frac{P(e_j = v_k, B_i)}{P(B_i)} = \frac{\text{Count of } (e_j = v_k, B_i)}{\text{Count of } B_i}$$

Question: What if there is no observation of $(e_j = v_k, B_i)$, i.e.,
 $\text{Count of } (e_j = v_k, B_i) = 0$?

Zero Frequency Problem and α -Laplace Smoothing

Problem

If the categorical feature has a category in the test data set, which was not observed in the training data set, then the model will assign a 0 probability and will be unable to make a prediction. This is often known as **Zero Frequency**.

Solution: α -Laplacian Smoothing

Idea: Pretend you have seen every value of the feature given every belief α more time than you actually did, i.e.,

$$P(e_j = v_k | B_i) = \frac{\text{Count of } (e_j=v_k, B_i) + \alpha}{\text{Count of } B_i + m \cdot \alpha}$$

where m is the number of possible values of e_j .

$\alpha = 1$ or 1-Laplace Smoothing: Pretend that you have seen 1 more time you actually did.

Predicting Suffering from Disease 'Z' using Naïve Bayes Classifier

Data	Blood Pressure	Fever	Diabetes	Vomit	Suffering from disease Z
1	High	High	Yes	No	No
2	High	High	Yes	Yes	No
3	Low	High	Yes	No	Yes
4	Normal	Mild	Yes	No	Yes
5	Normal	No fever	No	No	Yes
6	Normal	No fever	No	Yes	No
7	Low	No fever	No	Yes	Yes
8	High	Mild	Yes	No	No
9	High	No fever	No	No	Yes
10	Normal	Mild	No	No	Yes
11	High	Mild	No	Yes	Yes
12	Low	Mild	Yes	Yes	Yes
13	Low	High	No	No	Yes
14	Normal	Mild	Yes	Yes	No

For a patient with symptoms (High, No, Yes, Yes), what is the probability of suffering from disease 'Z'?

Predicting Suffering from Disease 'Z' using Naïve Bayes Classifier

Counts:

Blood Pressure			Fever			Diabetes			Vomit			Disease 'Z'	
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No
High	2	3	High	2	2	Yes	3	4	Yes	3	3	9	5
Normal	3	2	Mild	4	2	No	6	1	No	6	2		
Low	4	0	No Fever	3	1								

Relative Frequency:

Blood Pressure			Fever			Diabetes			Vomit			Disease 'Z'	
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No
High	$\frac{2}{9}$	$\frac{3}{5}$	High	$\frac{2}{9}$	$\frac{2}{5}$	Yes	$\frac{3}{9}$	$\frac{4}{5}$	Yes	$\frac{3}{9}$	$\frac{3}{5}$	$\frac{9}{14}$	$\frac{5}{14}$
Normal	$\frac{3}{9}$	$\frac{2}{5}$	Mild	$\frac{4}{9}$	$\frac{2}{5}$	No	$\frac{6}{9}$	$\frac{1}{5}$	No	$\frac{6}{9}$	$\frac{2}{5}$		
Low	$\frac{4}{9}$	$\frac{0}{5}$	No Fever	$\frac{3}{9}$	$\frac{1}{5}$								

Predicting Suffering from Disease ‘Z’ using Naïve Bayes Classifier

- Given the naïve bayes' formula

$$B_{NB} = \operatorname{argmax}_{B_i} P(B_i)(P(e_1|B_i)P(e_2|B_i)P(e_3|B_i)\cdots P(e_d|B_i))$$

we have

$$\begin{aligned} &P(\text{Disease}Z = \text{Yes})(P(\text{BloodPressure}|\text{Yes})P(\text{Fever}|\text{Yes})P(\text{Diabetes}|\text{Yes})P(\text{Vomit}|\text{Yes})) \\ &P(\text{Disease}Z = \text{No})(P(\text{BloodPressure}|\text{No})P(\text{Fever}|\text{No})P(\text{Diabetes}|\text{No})P(\text{Vomit}|\text{No})) \end{aligned}$$

- Now, (Blood Pressure, Fever, Diabetes, Vomit) = (High, No, Yes, Yes)

$$\begin{aligned} &P(\text{Disease}Z = \text{Yes}) \cdot (P(\text{BP} = \text{High}|\text{Yes}) \cdot P(\text{Fever} = \text{No}|\text{Yes}) \cdot P(\text{Diabetes} = \text{Yes}|\text{Yes}) \\ &\quad \cdot P(\text{Vomit} = \text{Yes}|\text{Yes})) = \left(\frac{9}{14}\right) \left(\frac{2}{9}\right) \left(\frac{3}{9}\right) \left(\frac{3}{9}\right) = 0.005291 \end{aligned}$$

$$\begin{aligned} &P(\text{Disease}Z = \text{No}) \cdot (P(\text{BP} = \text{High}|\text{No}) \cdot P(\text{Fever} = \text{No}|\text{No}) \cdot P(\text{Diabetes} = \text{Yes}|\text{No}) \\ &\quad \cdot P(\text{Vomit} = \text{Yes}|\text{No})) = \left(\frac{5}{14}\right) \left(\frac{3}{5}\right) \left(\frac{1}{5}\right) \left(\frac{4}{5}\right) \left(\frac{3}{5}\right) = 0.020571 \end{aligned}$$

- Result:** The patient is not suffering from Disease ‘Z’. Since $0.020571 > 0.005291$.
- After normalization, we have $0.020571 / (0.020571 + 0.005291) = 0.795414121$, $0.005291 / (0.020571 + 0.005291) = 0.204585879$.

Naïve Bayes Implementation using Scikit-Learn

```
import numpy as np
# CategoricalNB is for doing Naive Bayes classifier for categorical features
from sklearn.naive_bayes import CategoricalNB

# Forming training data
# Features: Blood Pressure (0=High, 1=Normal, 2=Low), Fever (0=High, 1=Mild, 2=No Fever),
#             Diabetes (0=Yes, 1=No), Vomit (0=Yes, 1=No)
training = np.array([[0, 0, 0, 1], [0, 0, 0, 0], [2, 0, 0, 1], [1, 1, 0, 1], [1, 2, 1, 1],
                    [1, 2, 1, 0], [2, 2, 1, 0], [0, 1, 0, 1], [0, 2, 1, 1], [1, 1, 1, 1],
                    [0, 1, 1, 0], [2, 1, 0, 0], [2, 0, 1, 1], [1, 1, 0, 0]])

# Forming the label set
# Labels: 0=Yes, 1=No
outcome = np.array([1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1])

# This is the new data to be evaluated using the trained AI (model)
new_sample = np.array([[0, 2, 0, 0]])

# Train Naive Bayes classifier according to training, outcome
clf = CategoricalNB(alpha=1.0e-10).fit(training, outcome)
# Perform classification on new_sample, and get the probability estimates
pred_class = clf.predict(new_sample); prob = clf.predict_proba(new_sample)
# Print the results
print(pred_class, prob)
```

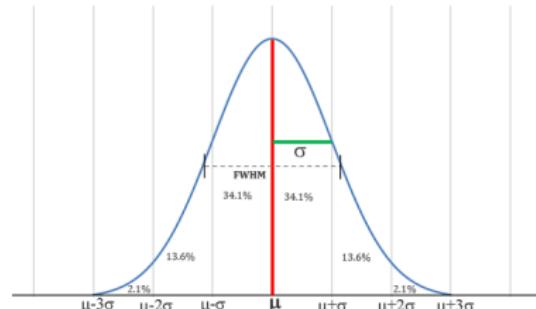
Calculation for Continuous Variables

- In the last example, the symptoms – Blood Pressure and Fever are discrete values. What if they are **continuous**?
- If the Blood Pressure and Fever were continuous, the calculation would have been different.
- Assume the **probability of likelihood** follows **Gaussian distribution**, i.e.,

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

where

- μ : Mean
- σ : Standard deviation



Calculation for Continuous Variables

- We need to calculate the mean and standard deviation of Blood Pressure and Fever for each belief (i.e., suffering from disease 'Z' and not suffering from disease 'Z').
- Suppose,
 - For Blood Pressure
 - Suffering from disease 'Z' is Yes: Mean (μ) = 73, Standard deviation (σ) = 6.2
 - Suffering from disease 'Z' is No: Mean (μ) = 75, Standard deviation (σ) = 7.9
 - For Fever
 - Suffering from disease 'Z' is Yes: Mean (μ) = 79, Standard deviation (σ) = 10.2
 - Suffering from disease 'Z' is No: Mean (μ) = 86, Standard deviation (σ) = 9.7

Question: Any potential problem of this approach?



Calculation for Continuous Variables

- For Blood Pressure = 66, what will be $P(\text{Blood Pressure} = 66 | \text{Disease } Z = \text{Yes})$?

$$P(\text{Blood Pressure} = 66 | \text{Disease } Z = \text{Yes}) = \frac{1}{\sqrt{2\pi}(6.2)} \exp\left(-\frac{(66 - 73)^2}{2 * (6.2)^2}\right) = 0.0340187$$

- For Fever = 97, what will be $P(\text{Fever} = 97 | \text{Disease } Z = \text{Yes})$?

$$P(\text{Fever} = 97 | \text{Disease } Z = \text{Yes}) = \frac{1}{\sqrt{2\pi}(10.2)} \exp\left(-\frac{(97 - 79)^2}{2 * (10.2)^2}\right) = 0.00824276$$

- Posterior

$$\begin{aligned} & P(\text{Disease } Z = \text{Yes})(P(\text{BP} = 66 | \text{Yes})P(\text{Fever} = 97 | \text{Yes})P(\text{Diabetes} = \text{Yes} | \text{Yes})P(\text{Vomit} = \text{Yes})) \\ &= \left(\frac{9}{14}\right)(0.0340187)(0.00824276)\left(\frac{3}{9}\right)\left(\frac{3}{9}\right) = 0.00002 \end{aligned}$$

* Please note that the probability of a continuous random variable taking an exact value is 0. For the sake of simplicity, we have used the notations loosely here. For details, please refer to the supplementary notes: Why can $P(e_i | B_i)$ be found by substituting $x = e_i$ into $f(x)$?

Calculation for Continuous Variables

- For Blood Pressure = 66, what will be $P(\text{Blood Pressure} = 66 | \text{Disease } Z = \text{No})$?

$$P(\text{Blood Pressure} = 66 | \text{Disease } Z = \text{No}) = \frac{1}{\sqrt{2\pi}(7.9)} \exp\left(-\frac{(66 - 75)^2}{2 * (7.9)^2}\right) = 0.0263909$$

- For Fever = 97, what will be $P(\text{Fever} = 97 | \text{Disease } Z = \text{No})$?

$$P(\text{Fever} = 97 | \text{Disease } Z = \text{No}) = \frac{1}{\sqrt{2\pi}(9.7)} \exp\left(-\frac{(97 - 86)^2}{2 * (9.7)^2}\right) = 0.0216215$$

- Posterior

$$\begin{aligned} P(\text{Disease } Z = \text{No})(P(\text{BP} = 66 | \text{No})P(\text{Fever} = 97 | \text{No})P(\text{Diabetes} = \text{Yes} | \text{No})P(\text{Vomit} = \text{Yes} | \text{No})) \\ = \left(\frac{5}{14}\right)(0.0263909)(0.0216215)\left(\frac{4}{5}\right)\left(\frac{3}{5}\right) = 0.0000978 \end{aligned}$$

- Result:** The Patient is not suffering from Disease 'Z'. Since $0.0000978 > 0.00002$.
- After normalization, we have $0.0000978 / (0.0000978 + 0.00002) = 0.8302207$, $0.00002 / (0.0000978 + 0.00002) = 0.1697793$.

* Please note that the probability of a continuous random variable taking an exact value is 0. For the sake of simplicity, we have used the notations loosely here. For details, please refer to the supplementary notes: Why can $P(e_i | B_i)$ be found by substituting $x = e_i$ into $f(x)$?

Underflow Prevention

- Naïve Bayes classifier:

$$B_{NB} = \arg \max_{B_i} P(B_i)(P(e_1|B_i)P(e_2|B_i)P(e_3|B_i)\cdots P(e_d|B_i))$$

- Multiplying lots of probabilities between 0 and 1 by definition, can result in **floating-point underflow**.
- $\log(xy) = \log(x) + \log(y)$, it is better to perform all computations by **summing logs of probabilities** rather than multiplying probabilities.
- Belief with the highest final log probability score is still the most probable.

$$B_{NB} = \operatorname{argmax}_{B_i} \left(\log P(B_i) + \sum_{j=1}^d \log P(e_j|B_i) \right)$$

Logarithmic Formulas

$$\log[ab] = \log a + \log b$$

$$\log\left[\frac{a}{b}\right] = \log a - \log b$$

$$\log(a^b) = b \log a \quad a^{\log a^b} = b$$

Applications of Naïve Bayes Algorithm

- **Real-time prediction:** Naïve Bayes Algorithm is **fast** and always ready to learn hence best suited for real-time predictions.
- **Multi-class prediction:** The probability of multi-classes of any target variable can be predicted using a Naïve Bayes algorithm.
- **Text classification** where Naïve Bayes is mostly used is **Spam Filtering** in emails (Naïve Bayes is widely used for text classification). Due to its better performance with multi-class problems and its independence rule, the Naïve Bayes algorithm performs better or have a higher success rate in text classification. Therefore, it is used in **sentiment analysis and spam filtering**.
- **Recommendation system:** Naïve Bayes classifier and collaborative filtering together build a recommendation system that uses machine learning and data mining techniques to filter unseen information and predict whether a user would like a given resource or not.

Advantages and Disadvantages of Naïve Bayes

• Advantages

- It is **very easy to implement** and obtain very good results on large datasets.
- It has **computational efficiency**. The assumption that all features are independent makes a Naïve Bayes algorithm very fast compared to complicated algorithms. In some cases, speed is preferred over higher accuracy.
- It **works well on a large size dataset**.
- It can also be used to **predict multiple classes**.
- It can also be used in the **Natural Language Processing text classification**.

• Disadvantages

- The strong assumption about the features to be **independent** which is **hardly true in real-life applications**.
- If the dataset size is **small**, it will **decrease the precision value**.
- When the goal is to predicting **probability** instead of classification, then the method provides very **biased results**.
- Having trouble with the **zero frequency problem**.

Practice Problem

	Animals	Size	Color	Can we pet them?
1	Dog	Medium	Black	Yes
2	Dog	Big	White	No
3	Rat	Small	White	Yes
4	Cow	Big	White	Yes
5	Cow	Small	Brown	No
6	Cow	Big	Black	Yes
7	Rat	Big	Brown	No
8	Dog	Small	Brown	Yes
9	Dog	Medium	Brown	Yes
10	Cow	Medium	White	No
11	Dog	Small	Black	Yes
12	Rat	Medium	Black	No
13	Rat	Small	Brown	No
14	Cow	Big	White	Yes

For the animal (Cow, Medium, Black), can we pet it?

Practice Problem

- Given the naïve bayes' formula

$$P(B_i|E) = P(B_i)(P(e_1|B_i)P(e_2|B_i)P(e_3|B_i)\cdots P(e_n|B_i))$$

we have

$$\begin{aligned} &P(Pet = Yes | (Animals, Size, Color)) \\ &= P(Pet = Yes)(P(Animals|Yes)P(Size|Yes)P(Color|Yes)) \end{aligned}$$

- Now, (Animals, Size, Color) = (Cow, Medium, Black)

$$\begin{aligned} &P(Pet = Yes)(P(Cow|Yes)P(Medium|Yes)P(Black|Yes)) \\ &= \left(\frac{8}{14}\right) \left(\frac{3}{8}\right) \left(\frac{2}{8}\right) \left(\frac{3}{8}\right) = 0.02 \end{aligned}$$

$$\begin{aligned} &P(Pet = No)(P(Cow|No)P(Medium|No)P(Black|No)) \\ &= \left(\frac{6}{14}\right) \left(\frac{2}{6}\right) \left(\frac{2}{6}\right) \left(\frac{1}{6}\right) = 0.0079 \end{aligned}$$

- Result: We can pet it. Since $0.02 > 0.0079$.
- After normalization, we have $0.02/(0.02 + 0.0079) = 0.7168459$, $0.0079/(0.02 + 0.0079) = 0.2831541$.

Naïve Bayes Implementation using Scikit-Learn

```
import numpy as np
# CategoricalNB is for doing Naive Bayes classifier for categorical features
from sklearn.naive_bayes import CategoricalNB

# Forming training data
# Features: Animals (0=Dog, 1=Rat, 2=Cow), Size (0=Big, 1=Medium, 2=Small),
#             Color (0=Black, 1=White, 2=Brown)
training = np.array([[0, 1, 0], [0, 0, 1], [1, 2, 1], [2, 0, 1], [2, 2, 2],
                     [2, 0, 0], [1, 0, 2], [0, 2, 2], [0, 1, 2], [2, 1, 1],
                     [0, 2, 0], [1, 1, 0], [1, 2, 2], [2, 0, 1]])

# Forming the label set
# Labels: 0=Yes, 1=No
outcome = np.array([0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 1, 0])

# This is the new data to be evaluated using the trained AI (model)
new_sample = np.array([[2, 1, 0]])

# Train Naive Bayes classifier according to training, outcome
clf = CategoricalNB(alpha=1.0e-10).fit(training, outcome)
# Perform classification on new_sample and get the probability estimates
pred_class = clf.predict(new_sample); prob = clf.predict_proba(new_sample)
# Print the results
print(pred_class, prob)
```

That's all!

Any question?



Welcome
Back!

Acknowledgments

- The lecture notes are developed based on Dr. Desmond Tsoi's lecture slides.