

COMP3055

Machine Learning

Topic 7 – Bayesian Learning

Dr. Zheng LU
2018 Autumn

Probability

- The world is a very uncertain place.
- 30 years of Artificial Intelligence research danced around this fact.
- And then a few AI researchers decided to use some ideas from the eighteenth century.

Random Variables

- A variable whose possible values are outcomes of a random phenomenon.
- Denotes a quantity (event) that is uncertain.
- Maybe result of experiment (flipping a coin) or a real world observation (measuring attendance rate).
- If observe several instances of a random variable, we get different values.
- Some values occur more than others and this information is captured by a probability distribution.

Random Variables

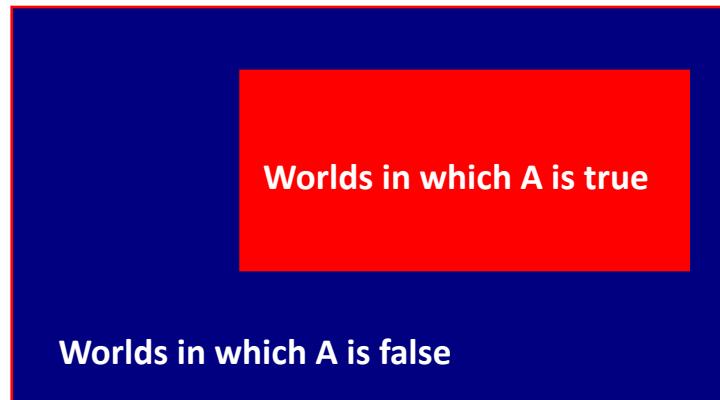
- A is a Boolean-valued random variable
 - if A denotes such event that there is some degree of uncertainty as to whether A occurs.
- Examples
 - A = The US president in 2023 will be male
 - A = You wake up tomorrow with a headache
 - A = You have Ebola

Probabilities

We write $P(A)$ as

“the fraction of possible worlds in which A is true”

Event space of all
possible worlds



The whole area is 1

→

$P(A) = \text{Area of}$
 red rectangle

All probabilities between 0 and 1

$$0 \leq P(A) \leq 1$$

True proposition has probability 1, false has probability 0.

$$P(\text{not } A) = P(\sim A) = 1 - P(A)$$

Multivalued Random Variables

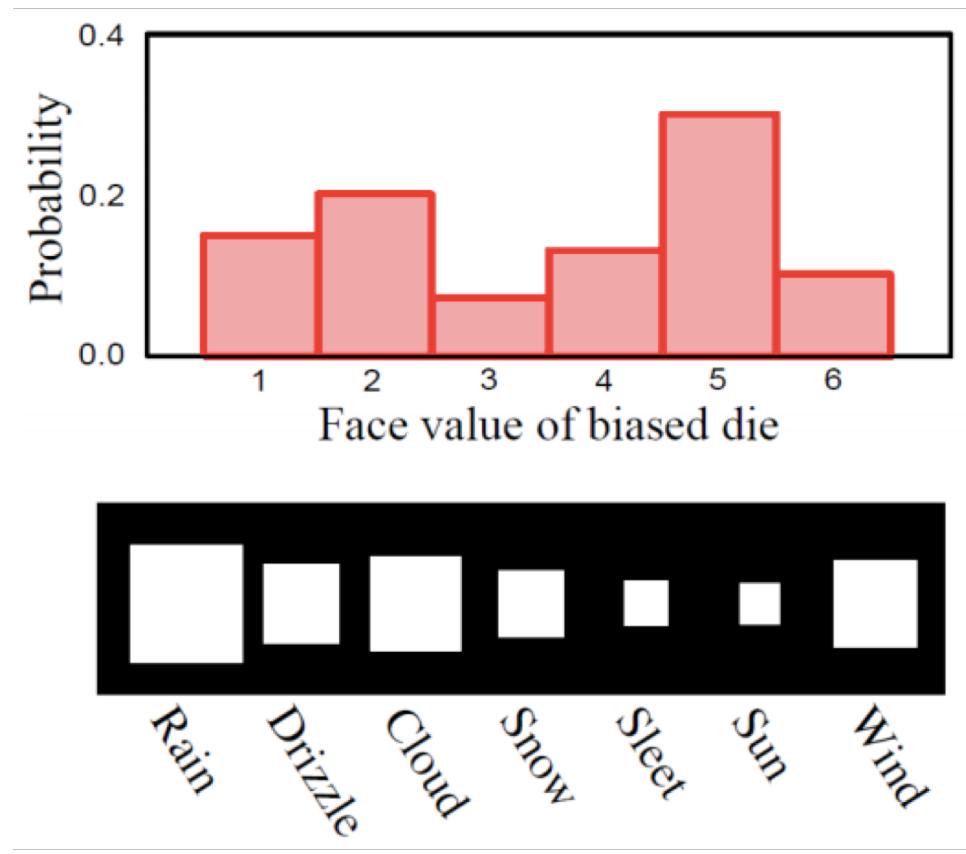
Suppose a random variable A can take on exactly one value out of $\{v_1, v_2, \dots, v_k\}$

$$P(A = v_i \wedge A = v_j) = 0 \quad \text{if } i \neq j$$

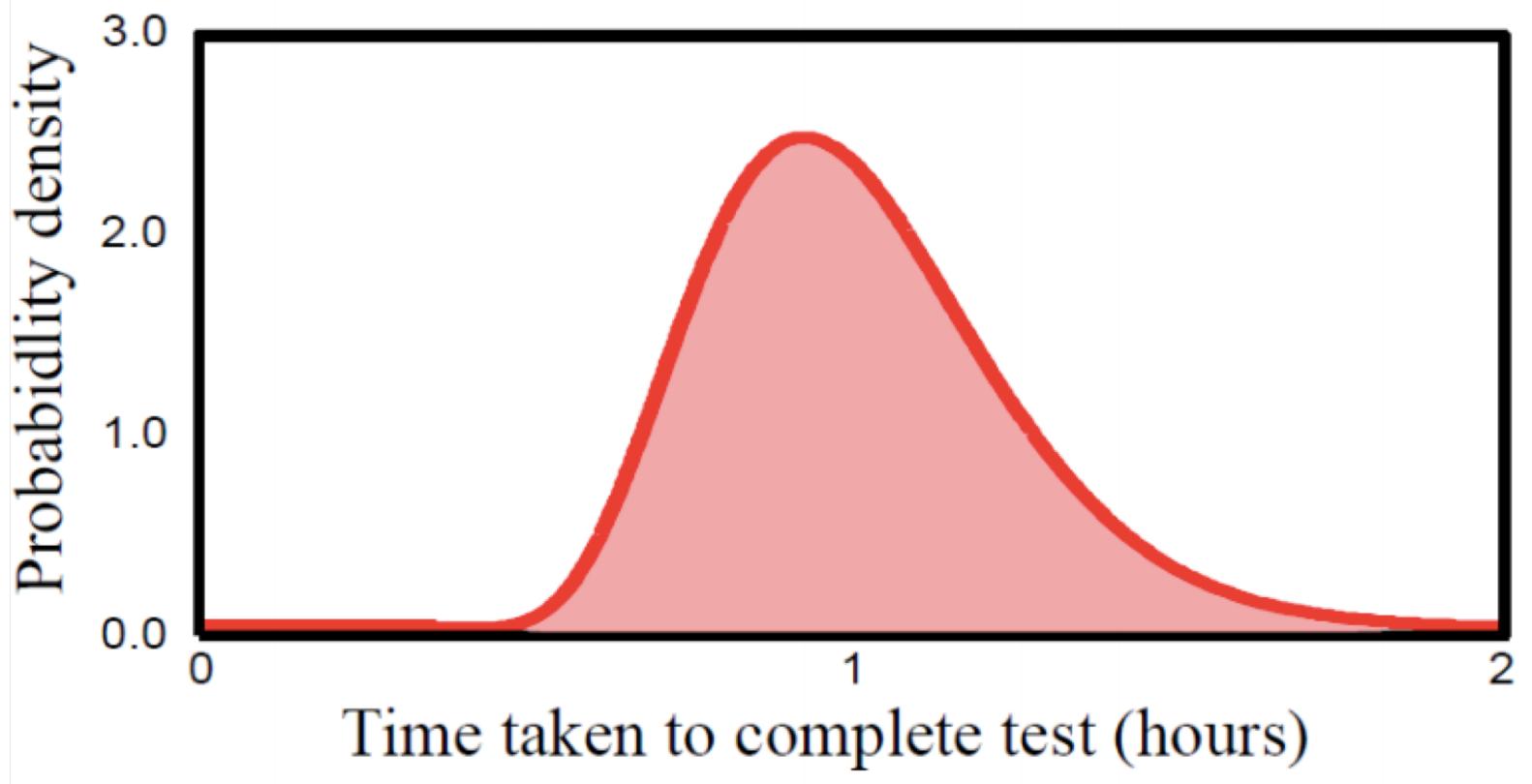
$$P(A = v_1 \vee A = v_2 \vee \dots \vee A = v_k) = 1$$

Discrete vs Continuous Random Variable

Depends on whether outcome values are discrete or continuous



Continuous Random Variables

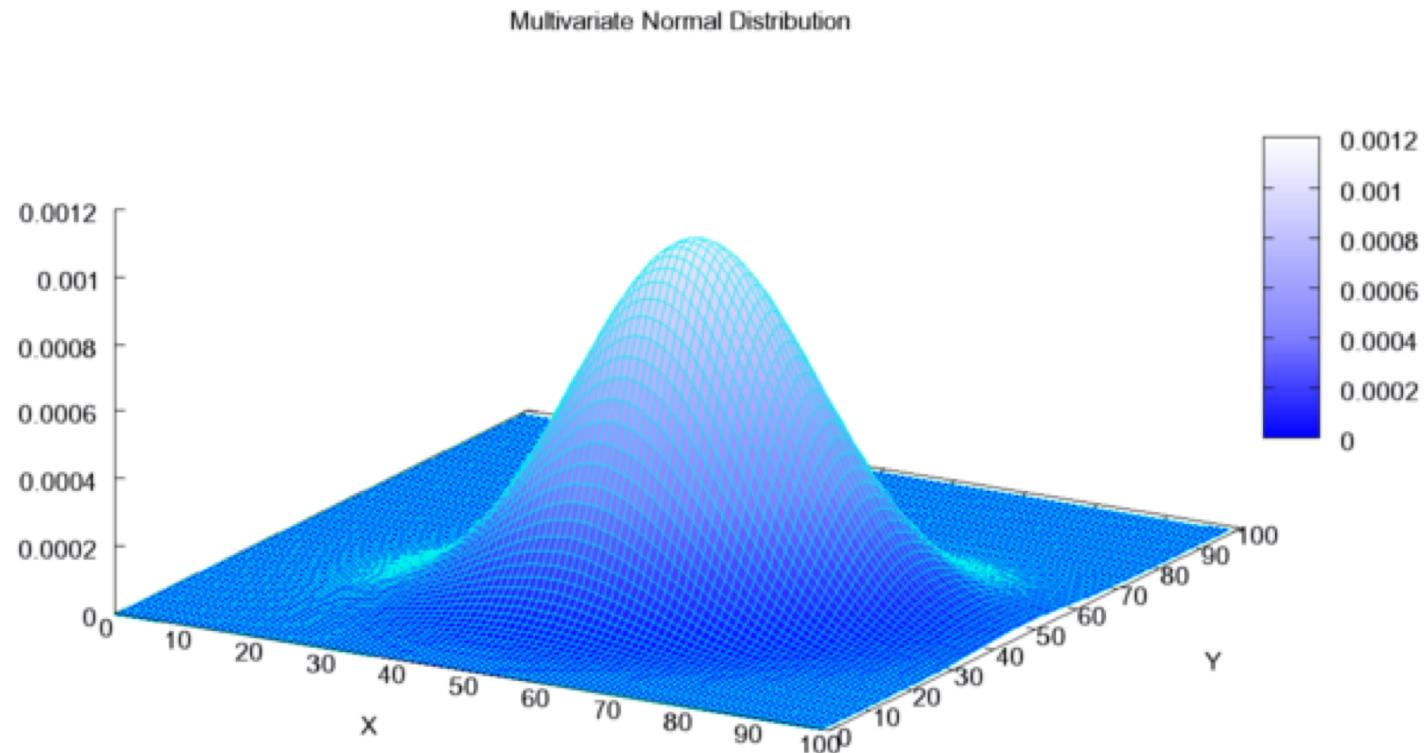


Note: Instead of summing over all possible outcome in discrete random variable, integrating the probability over all outcome should be 1 in continuous random variable.

Joint Probability

- Consider two random variables A and B.
- If we observe multiple paired instances, some combinations of outcomes are more likely than others.
- This is captured by joint probability $P(A \wedge B)$.

Joint Probability



Marginalization

$$P(A = v_1 \vee A = v_2 \vee \cdots \vee A = v_i) = \sum_{j=1}^i P(A = v_j)$$

$$P(B \wedge [A = v_1 \vee A = v_2 \vee \cdots \vee A = v_i]) = \sum_{j=1}^i P(B \wedge A = v_j)$$

$$P(B) = \sum_{j=1}^k P(B \wedge A = v_j)$$

Marginalization: We can recover probability distribution of any variable in a joint distribution by summing over the other variables

Conditional Probability

$P(A|B)$ = Fraction of worlds in which B is true that also have A true.

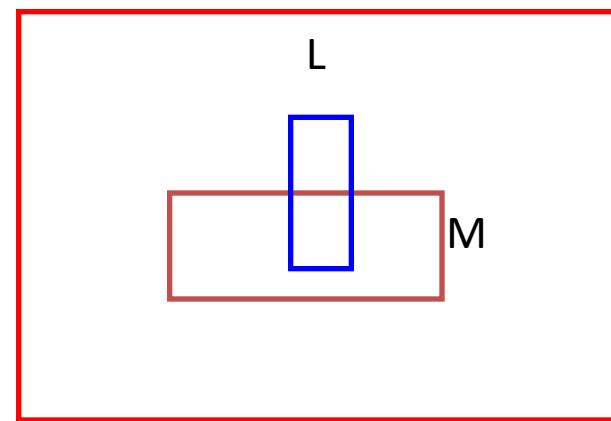
M = “Miss 9am lecture”

L = “Go to bed late”

$$P(M) = 1/4$$

$$P(L) = 1/10$$

$$P(M|L) = 1/2$$



“Missing 9am class is rare and going to bed late is rarer. But if you go to bed late, there is a 50-50 chance you will miss 9am lecture.”

Conditional Probability

$P(A|B)$ = Fraction of worlds in which B is true that also have A true.

$P(M|L)$ = Fraction of late sleeper worlds in which also miss 9am lecture

$$= \frac{\text{#worlds of late sleeper and miss 9am class}}{\text{-----}}$$

#worlds late sleeper

$$= \frac{\text{Area of "M and L" region}}{\text{-----}}$$

Area of "L" region

$$= \frac{P(M \wedge L)}{\text{-----}}$$

$P(L)$

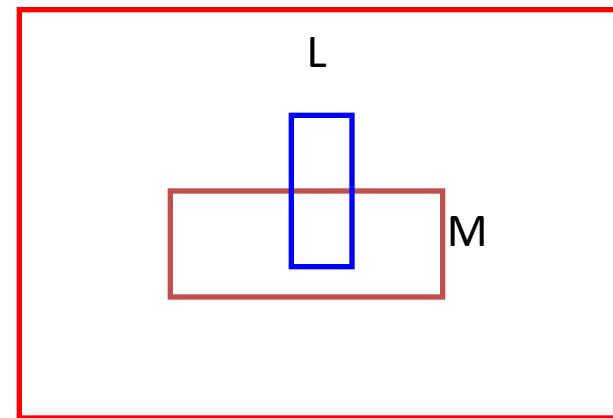
M = “Miss 9am lecture”

L = “Go to bed late”

$$P(M) = 1/4$$

$$P(L) = 1/10$$

$$P(M|L) = 1/2$$



Definition of Conditional Probability

Conditional probability definition

$$P(A|B) = \frac{P(A \wedge B)}{P(B)}$$

Hence, the Chain Rule

$$P(A \wedge B) = P(A|B)P(B)$$

Probabilistic Inference

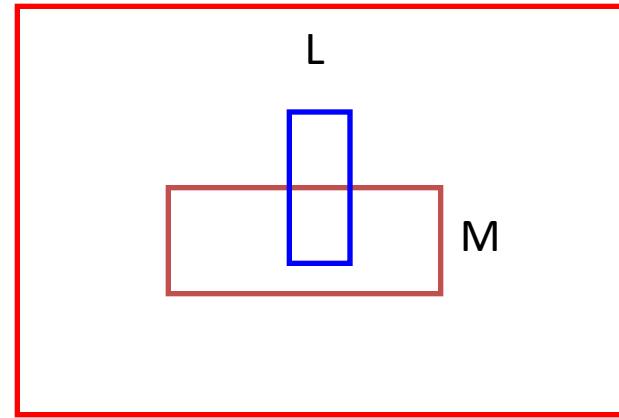
M = “Miss 9am lecture”

L = “Go to bed late”

$$P(M) = 1/4$$

$$P(L) = 1/10$$

$$P(M | L) = 1/2$$



One day you miss my 9am class.

I think, “50% of late sleepers miss my 9am lecture,
so this student (you) must have a 50-50 chance of
going to bed late last night.”

Is this thought reasonable?

Probabilistic Inference

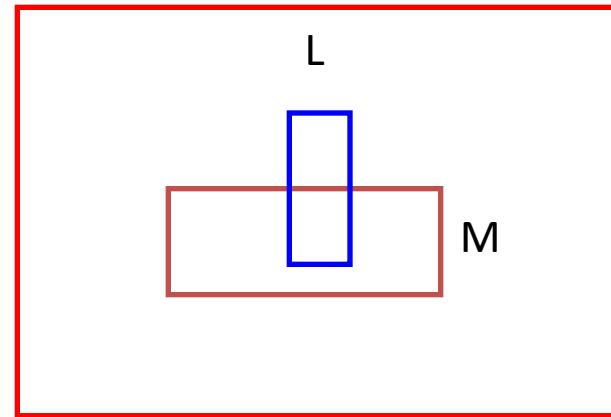
M = “Miss 9am lecture”

L = “Go to bed late”

$$P(M) = 1/4$$

$$P(L) = 1/10$$

$$P(M|L) = 1/2$$



One day you miss my 9am class.

I think, “50% of late sleepers miss my 9am lecture, so this student (you) must have a 50-50 chance of going to bed late last night.” $\leftarrow P(L|M)$

Is this thought reasonable?

$$\begin{aligned} P(M \wedge L) &= P(M|L)P(L) \\ &= \frac{1}{2} * \frac{1}{10} = \frac{1}{20} \end{aligned}$$

$$\begin{aligned} P(L|M) &= \frac{P(L \wedge M)}{P(M)} \\ &= \frac{1/20}{1/4} = \frac{1}{5} \end{aligned}$$

Probabilistic Inference

M = “Miss 9am lecture”

L = “Go to bed late”

$$P(M) = 1/4$$

$$P(L) = 1/10$$

$$P(M | L) = 1/2$$

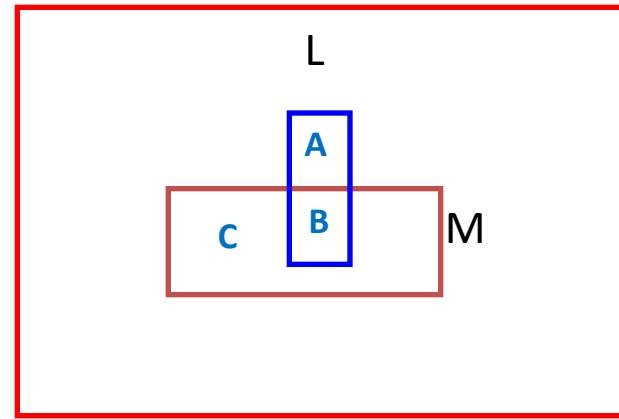
Area wise we have

$$P(L) = A + B$$

$$P(M) = B + C$$

$$P(M | L) = B / (A + B)$$

$$P(L | M) = B / (B + C) = P(M | L) P(L) / P(M)$$



Bayes Rule

- **Bayes, Thomas (1763)** An essay towards solving a problem in the doctrine of chances. Philosophical Transactions of the Royal Society of London, **53:370-418**
- **Bayes Rule**

$$p(A | B) = \frac{P(A \wedge B)}{P(B)} = \frac{P(B | A)P(A)}{P(B)}$$



Bayesian Learning

$$p(h | x) = \frac{P(x | h)P(h)}{P(x)}$$

Understanding Bayes' rule

x = data

h = hypothesis (model)

- rearranging

$$p(h | x)P(x) = P(x | h)P(h)$$

$$P(x, h) = P(x, h)$$

the same joint probability
on both sides

$P(h)$: prior belief (probability of hypothesis h before seeing any data)

$P(x | h)$: likelihood (probability of the data if the hypothesis h is true)

$P(x) = \sum_h P(x | h)P(h)$: data evidence (marginal probability of the data)

$P(h | x)$: posterior (probability of hypothesis h after having seen the data d)

Choosing Hypotheses

- Generally, we want the most probable hypothesis(class label) given the observed data
 - Maximum a posteriori (**MAP**) hypothesis
 - Maximum likelihood (**ML**) hypothesis

Maximum A Posteriori (MAP)

- Maximum a posteriori (MAP) hypothesis

$$p(h | x) = \frac{P(x | h)P(h)}{P(x)}$$

$$h_{MAP} = \arg \max_{h \in H} p(h | x) = \arg \max_{h \in H} \frac{P(x | h)P(h)}{P(x)} = \arg \max_{h \in H} P(x | h)P(h)$$

Note $P(x)$ is independent of h , hence can be ignored.

Maximum Likelihood (ML)

$$h_{MAP} = \arg \max_{h \in H} P(x | h)P(h)$$

- Assuming that each hypothesis in H is equally probable, i.e., $P(h_i) = P(h_j)$, for all i and j, then we can drop $P(h)$ in MAP. $P(x|h)$ is often called the likelihood of data x given h . Any hypothesis that maximizes $P(x|h)$ is called the maximum likelihood hypothesis

$$h_{ML} = \arg \max_{h \in H} P(x | h)$$

An Illustrating Example

Classifying days according to whether someone will play tennis.

Each day is described by the attributes, Outlook, Temperature, Humidity and Wind.

Based on the training data in the table, classify the following instance

Outlook = sunny

Temperature = cool

Humidity = high

Wind = strong

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

An Illustrating Example

Training sample pairs (X, D)

$X = (x_1, x_2, \dots, x_n)$ is the feature vector representing the instance.

$D = (d_1, d_2, \dots, d_m)$ is the desired (target) output of the classifier

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

An Illustrating Example

Training sample pairs (X, D)

$X = (x_1, x_2, \dots, x_n)$ is the feature vector representing the instance.

$n = 4$

$x_1 = \text{outlook} = \{\text{sunny}, \text{overcast}, \text{rain}\}$

$x_2 = \text{temperature} = \{\text{hot}, \text{mild}, \text{cool}\}$

$x_3 = \text{humidity} = \{\text{high}, \text{normal}\}$

$x_4 = \text{wind} = \{\text{weak}, \text{strong}\}$

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

An Illustrating Example

Training sample pairs (X, D)

$D = (d_1, d_2, \dots d_m)$ is the desired
(target) output of the classifier

$m = 1$

$d = \text{Play Tennis} = \{\text{yes}, \text{no}\}$

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

Bayesian Classifier

- The Bayesian approach to classifying a new instance X is to assign it to the most probable target value Y (MAP classifier)

$$\begin{aligned}Y &= \arg \max_{d_i \in d} p(d_i | X) \\&= \arg \max_{d_i \in d} p(d_i | x_1, x_2, x_3, x_4) \\&= \arg \max_{d_i \in d} \frac{p(x_1, x_2, x_3, x_4 | d_i) P(d_i)}{p(x_1, x_2, x_3, x_4)} \\&= \arg \max_{d_i \in d} p(x_1, x_2, x_3, x_4 | d_i) P(d_i)\end{aligned}$$

Bayesian Classifier

$$Y = \arg \max_{d_i \in d} p(x_1, x_2, x_3, x_4 | d_i) P(d_i)$$

$P(d_i)$ is easy to calculate: simply counting how many times each target value d_i occurs in the training set.

$$P(d = yes) = 9/14$$

$$P(d = no) = 5/14$$

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

Bayesian Classifier

$$Y = \arg \max_{d_i \in d} p(x_1, x_2, x_3, x_4 | d_i) P(d_i)$$

$P(x_1, x_2, x_3, x_4 | d_i)$ is much more difficult to estimate.

In this simple example, there are $3 \times 3 \times 2 \times 2 \times 2 = 72$ possible terms.

To obtain a reliable estimate, we need to see each terms many times.

Hence, we need a very, very large training set! (which in most cases is impossible to get).

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

Naïve Bayes Classifier

Naïve Bayes classifier is based on the simplifying assumption that the attribute values are conditionally independent given the target value.

This means, we have

$$P(x_1, x_2, \dots, x_n | d_i) = \prod_i P(x_i | d_i)$$

Naïve Bayes Classifier

$$Y = \arg \max_{d_i \in d} P(d_i) \prod_{k=1}^4 P(x_k | d_i)$$

Back to the Example

Naïve Bayes Classifier

$$Y = \arg \max_{d_i \in d} \prod_{k=1}^4 P(x_k | d_i) P(d_i)$$

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

$$Y = \arg \max_{d_i \in \{yes, no\}} P(d_i) P(x_1 = suny | d_i) P(x_2 = cool | d_i) P(x_3 = high | d_i) P(x_4 = strong | d_i)$$

Back to the Example

$$Y = \arg \max_{d_i \in \{yes, no\}} P(d_i)P(x_1 = \text{sunny} | d_i)P(x_2 = \text{cool} | d_i)P(x_3 = \text{high} | d_i)P(x_4 = \text{strong} | d_i)$$

$$P(d=yes) = 9/14 = 0.64$$

$$P(d=no) = 5/14 = 0.36$$

$$P(x_1 = \text{sunny} | yes) = 2/9$$

$$P(x_1 = \text{sunny} | no) = 3/5$$

$$P(x_2 = \text{cool} | yes) =$$

$$P(x_2 = \text{cool} | no) =$$

$$P(x_3 = \text{high} | yes) =$$

$$P(x_3 = \text{high} | no) =$$

$$P(x_4 = \text{strong} | yes) = 3/9$$

$$P(x_4 = \text{strong} | no) = 3/5$$

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

Back to the Example

$$Y = \arg \max_{d_i \in \{yes, no\}} P(d_i)P(x_1 = \text{sunny} | d_i)P(x_2 = \text{cool} | d_i)P(x_3 = \text{high} | d_i)P(x_4 = \text{strong} | d_i)$$

$$P(\text{yes})P(x_1 = \text{sunny} | \text{yes})P(x_2 = \text{cool} | \text{yes})P(x_3 = \text{high} | \text{yes})P(x_4 = \text{strong} | \text{yes}) = 0.0053$$

$$P(\text{no})P(x_1 = \text{sunny} | \text{no})P(x_2 = \text{cool} | \text{no})P(x_3 = \text{high} | \text{no})P(x_4 = \text{strong} | \text{no}) = 0.0206$$

$Y = \text{Play Tennis} = \text{no}$

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

Estimating Probabilities

- So far, we estimate the probabilities by the fraction of times the event is observed to occur over the entire opportunities
- In the above example, we estimated

$$P(\text{wind=strong}|\text{play tennis=no}) = N_c/N,$$

where $N = 5$ is the total number of training samples for which $\text{play tennis} = \text{no}$, and N_c is the number of these for which wind=strong

- What happens if $N_c = 0$?

Estimating Probabilities

- When N_c is small, however, such approach provides poor estimation. To avoid this difficulty, we can adopt the **m-estimate** of probability

$$\frac{N_c + mP}{N + m}$$

where P is the prior estimate of the probability we wish to estimate, m is a constant called the equivalent sample size.

A typical method for choosing P in the absence of other information is to assume uniform priors: If an attribute has k possible values we set $P=1/K$.

For example, $P(\text{wind} = \text{strong} | \text{play tennis} = \text{no})$, we note wind has two possible values, so uniform priors means $P = 1/2$

Another Illustrative Example

- **Car theft Example**
 - Attributes are Color , Type , Origin, and the subject, stolen can be either yes or no.

Example No.	Color	Type	Origin	Stolen?
1	Red	Sports	Domestic	Yes
2	Red	Sports	Domestic	No
3	Red	Sports	Domestic	Yes
4	Yellow	Sports	Domestic	No
5	Yellow	Sports	Imported	Yes
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	No
10	Red	Sports	Imported	Yes

Another Illustrative Example

- **Car theft Example**

- We want to classify a Red Domestic SUV.
- Note there is no example of a Red Domestic SUV in our data set.

Example No.	Color	Type	Origin	Stolen?
1	Red	Sports	Domestic	Yes
2	Red	Sports	Domestic	No
3	Red	Sports	Domestic	Yes
4	Yellow	Sports	Domestic	No
5	Yellow	Sports	Imported	Yes
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	No
10	Red	Sports	Imported	Yes

Another Illustrative Example

- We need to estimate

$$P(x_i | d_j) = \frac{N_c + mP}{N + m}$$

N = the number of training examples for which $d = d_j$

N_c = number of examples for which $d = d_j$ and $x = x_i$

P = a priori estimate for $P(x_i|d_j)$

m = the equivalent sample size

Another Illustrative Example

- To classify a Red, Domestic, SUV, we need to estimate

$$Y = \arg \max_{d_i \in \{yes, no\}} P(d_i)P(x_1 = RED | d_i)P(x_2 = SUV | d_i)P(x_3 = Domestic | d_i)$$

Another Illustrative Example

Yes:

Red:

$N = 5$

$N_c = 3$

$P = .5$

$m = 3$

SUV:

$N = 5$

$N_c = 1$

$P = .5$

$m = 3$

No:

Red:

$N = 5$

$N_c = 2$

$P = .5$

$m = 3$

SUV:

$N = 5$

$N_c = 3$

$P = .5$

$m = 3$

Example No.	Color	Type	Origin	Stolen?
1	Red	Sports	Domestic	Yes
2	Red	Sports	Domestic	No
3	Red	Sports	Domestic	Yes
4	Yellow	Sports	Domestic	No
5	Yellow	Sports	Imported	Yes
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	No
10	Red	Sports	Imported	Yes

Domestic:

$N = 5$

$N_c = 2$

$P = .5$

$m = 3$

Domestic:

$N = 5$

$N_c = 3$

$P = .5$

$m = 3$

Another Illustrative Example

Example No.	Color	Type	Origin	Stolen?
1	Red	Sports	Domestic	Yes
2	Red	Sports	Domestic	No
3	Red	Sports	Domestic	Yes
4	Yellow	Sports	Domestic	No
5	Yellow	Sports	Imported	Yes
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	No
10	Red	Sports	Imported	Yes

$$P(\text{Red}|\text{Yes}) = \frac{3 + 3 * .5}{5 + 3} = .56$$

$$P(\text{Red}|\text{No}) = \frac{2 + 3 * .5}{5 + 3} = .43$$

$$P(\text{SUV}|\text{Yes}) = \frac{1 + 3 * .5}{5 + 3} = .31$$

$$P(\text{SUV}|\text{No}) = \frac{3 + 3 * .5}{5 + 3} = .56$$

$$P(\text{Domestic}|\text{Yes}) = \frac{2 + 3 * .5}{5 + 3} = .43$$

$$P(\text{Domestic}|\text{No}) = \frac{3 + 3 * .5}{5 + 3} = .56$$

Another Illustrative Example

- To classify a Red, Domestic, SUV, we need to estimate

$$Y = \arg \max_{d_i \in \{yes, no\}} P(d_i)P(x_1 = RED | d_i)P(x_2 = SUV | d_i)P(x_3 = Domestic | d_i)$$

$$\begin{aligned} & P(yes)P(x_1 = RED | yes)P(x_2 = SUV | yes)P(x_3 = Domestic | yes) \\ &= 0.5 * 0.56 * 0.31 * 0.43 = 0.037 \end{aligned}$$

$$\begin{aligned} & P(no)P(x_1 = RED | no)P(x_2 = SUV | no)P(x_3 = Domestic | no) \\ &= 0.5 * 0.43 * 0.56 * 0.56 = 0.069 \end{aligned}$$

$$Y = no$$

Naïve Bayesian Classifier

- What about using Naïve Bayesian Classifier for our handwritten digit recognition problem?
 - Each pixel is an x_i . There will be 784 x 's.
 - Digit label is d_k . Note there will be 10 possible d 's.
 - $P(d_k)$ can be calculated by counting number of training images for the digit, divided to total number of training images.
 - $P(x_i|d_k)$ can be calculated by counting number of images for a given digit, given pixel position, and given an intensity value, divided by number of training images with that digit.

Naïve Bayesian Classifier

$P(x_i|d_k)$ can be calculated by counting number of images for a given digit, given pixel position, and given an intensity value, divided by number of training images with that digit.

For example,

$$P(x_1 = 255|d = 0)$$

$$= \frac{\text{Number of images whose first pixel value is 255 and contain digit 0}}{\text{Total number of images contain digit 0}}$$

Naïve Bayesian Classifier

- For a given input image X and given digit label d_k , calculate $P(d_k)$ and all $P(x_i|d_k)$
- For each digit label d_k , calculate
$$P(d_k|X) = P(d_k)P(x_1 = 0|d_k)P(x_2 = 255|d_k) \dots P(x_{784} = 0)$$
- Choose the digit label k that give the max value according to above calculation.

Input image X

0	255	0	0	0	0	0	0	0	0	0
0	0	0	0	255	255	0	0	0	0	0
0	0	0	255	0	0	255	0	0	0	0
0	0	0	0	0	0	255	0	0	0	0
0	0	0	0	0	255	0	0	0	0	0
0	0	0	0	255	0	0	0	0	0	0
0	0	0	255	0	0	0	0	0	0	0
0	0	0	255	0	0	0	0	0	0	0
0	0	0	255	0	0	0	0	0	0	0
0	0	0	255	255	255	255	255	255	255	0

Further Readings

Chapter 6, T. M. Mitchell, Machine Learning,
McGraw-Hill International Edition, 1997