

Naïve Bayes classifier

Lecture 5

Mathematical predictions

- We can 'predict' where the spacecraft will be at noon in 2 months from now
- We cannot predict where you will be tomorrow at noon
- But, based on numerous observations, we can estimate the probability

Bayesian beliefs

- How do we judge that something is true?
- Can mathematics help make judgments more accurate?
- Bayes: our beliefs should be updated as new evidences become available



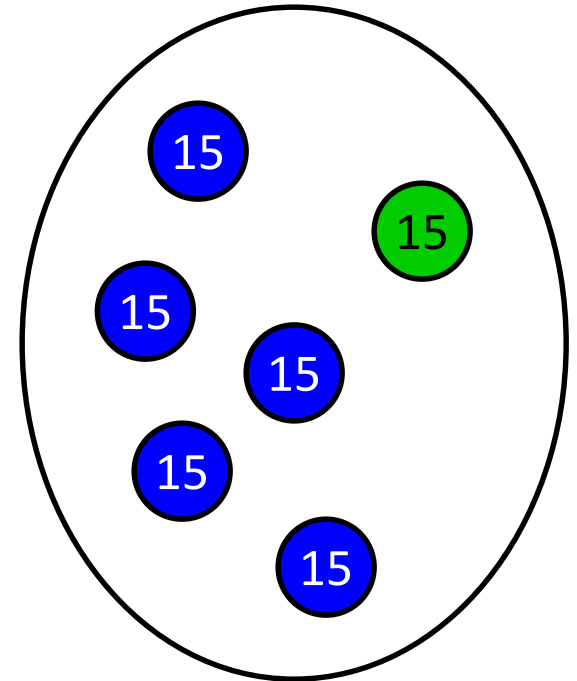
T. Bayes.

Bayes' method

- There are 2 events: **A** and not A (**B**) which you believe occur with probabilities $P(\mathbf{A})$ and $P(\mathbf{B})$. Estimation $P(\mathbf{A}):P(\mathbf{B})$ represents odds of A vs. B.
- Collect evidence data **E**.
- Re-estimate $P(\mathbf{A} | \mathbf{E}):P(\mathbf{B} | \mathbf{E})$ and update your beliefs.

Example (fictitious): hit-and-run

- 75 blue cabs (**B**) and 15 green cabs (**G**)
- $P(\mathbf{B}):P(\mathbf{G})=5:1$
- At night: hit-and-run episode
- Witness: “I saw a green cab”: X_G
- Witness is tested at night conditions:
identifies correct color 4 times out of 5
- Question: what is more probable:
B or **G**
?



Probability

- Basic element: **random variable**
e.g., *Car* is one of $\langle \text{blue}, \neg\text{blue}(\text{green}) \rangle$
Weather is one of $\langle \text{sunny}, \text{rainy}, \text{cloudy}, \text{snow} \rangle$
- Both *Car* and *Weather* are **discrete** random variables
 - Domain values must be
 - **exhaustive** (blue and green – are all the cabs)
 - **mutually exclusive** (green is always not blue, there are no cars which are half green, half blue)
- **Elementary propositions** are constructed by the assignment of a value to a random variable:
e.g., $\text{Car} = \neg\text{blue}$,
 $\text{Weather} = \text{sunny}$

Conditional probability

- $P(A | B)$ – probability of event A given that event B has happened

- In our case we want to compare:

the car was **G** given a witness testimony X_G : $P(\mathbf{G} | \mathbf{X}_G)$

vs.

the car was **B** given a witness testimony X_G : $P(\mathbf{B} | \mathbf{X}_G)$

Prior probability and distribution

- **Prior** or **unconditional probability** associated with a proposition is the degree of belief accorded to it in the absence of any other information.

e.g.,

$$P(\text{Car} = \text{blue}) = 0.83 \quad (\text{or abbrev. } P(\text{blue}) = 0.83)$$

$$P(\text{Weather} = \text{sunny}) = 0.7 \quad (\text{or abbrev. } P(\text{sunny}) = 0.7)$$

- **Probability distribution** gives probabilities of all possible value assignments:

$$P(\text{Weather} = \text{sunny}) = 0.7$$

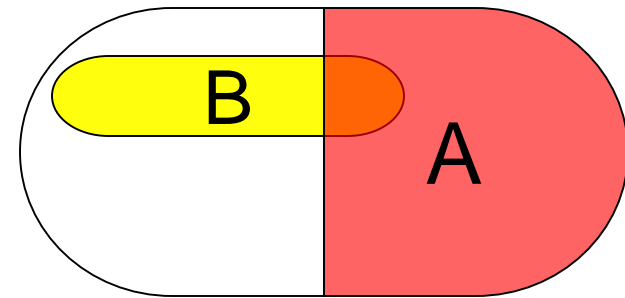
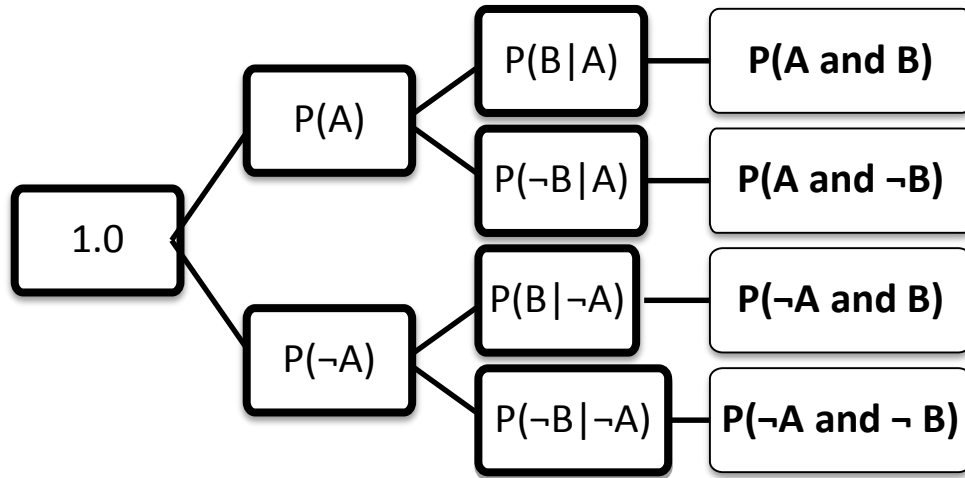
$$P(\text{Weather} = \text{rain}) = 0.2$$

$$P(\text{Weather} = \text{cloudy}) = 0.08$$

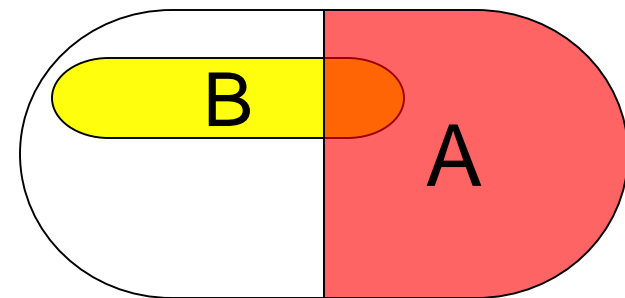
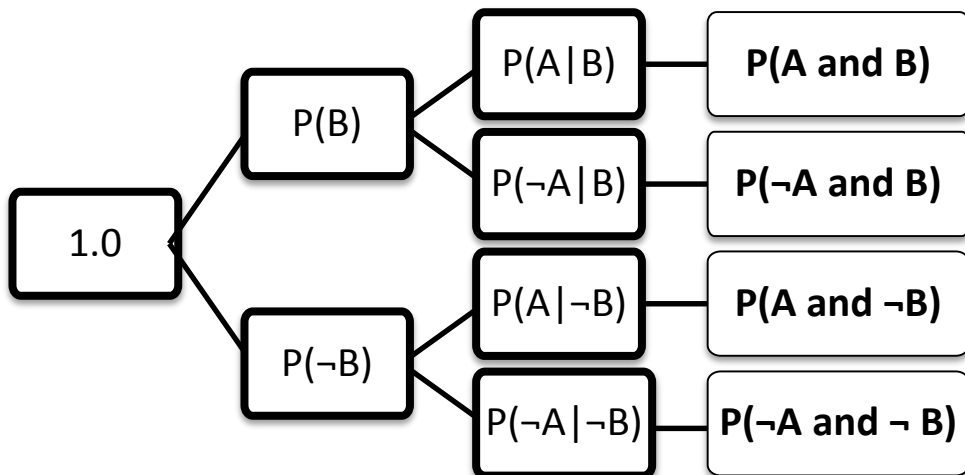
$$P(\text{Weather} = \text{snow}) = 0.02$$

- Sums up to 1.0

Two random events (not independent) happen at the same time – $P(A \text{ and } B)$



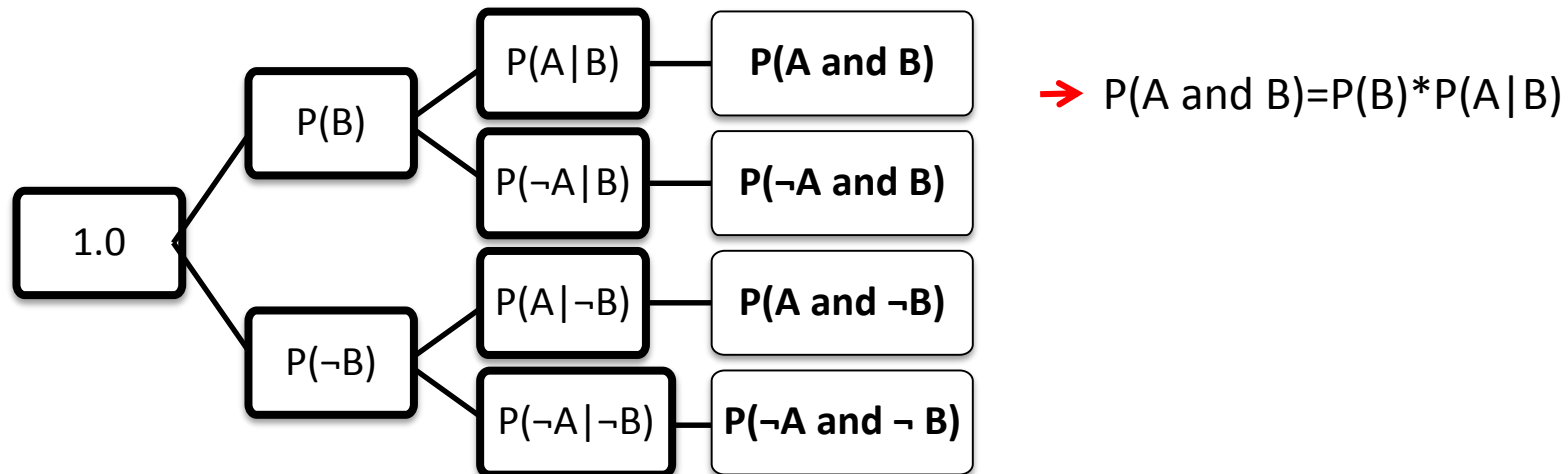
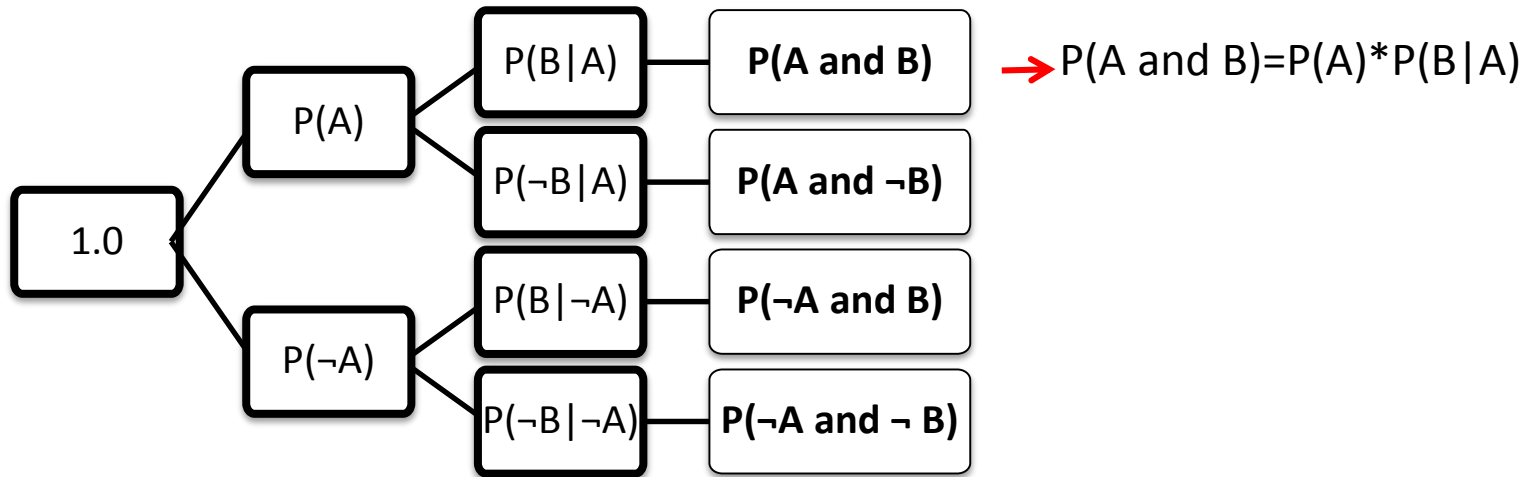
Possible event combinations when we know the outcome of event A:
 $P(B|A)=1/12$ and $P(A)=1/2$



Possible event combinations when we know the outcome of event B:
 $P(A|B)=1/4$ and $P(B)=1/6$

But in both cases $P(A \text{ and } B)$ is the same: orange area in the diagram

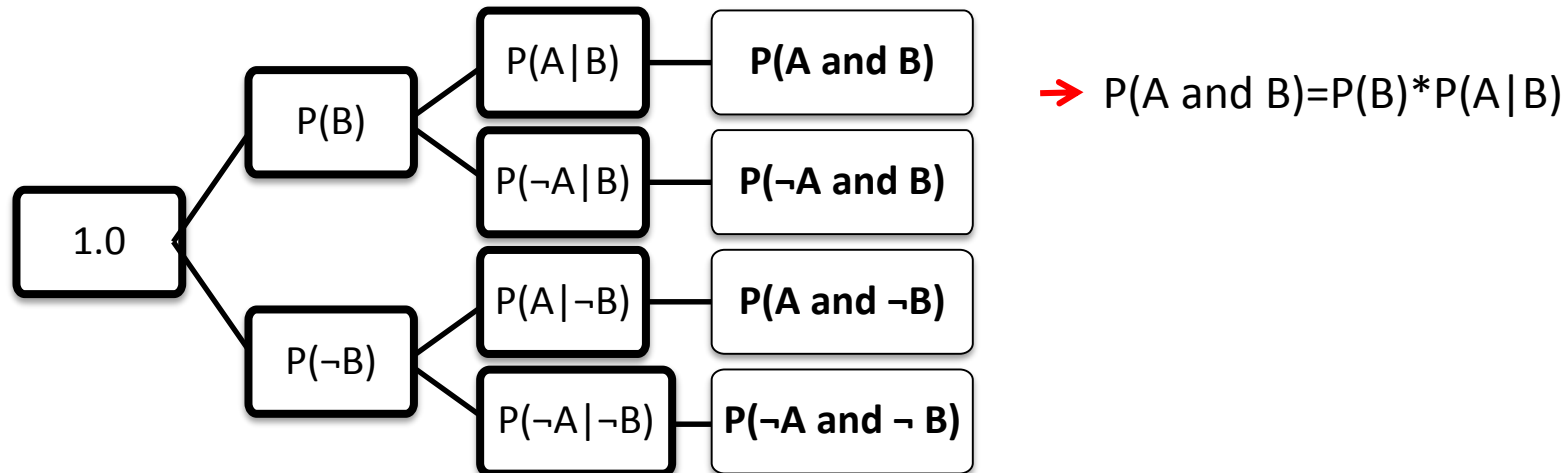
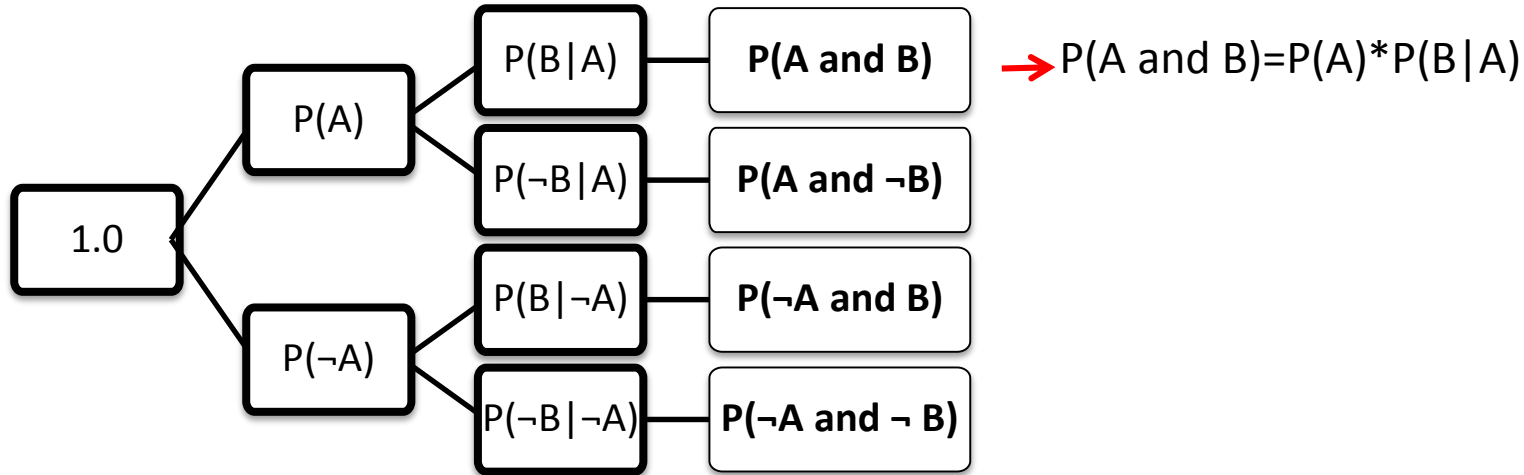
Intuition for Bayes's theorem



$$P(A \text{ and } B) = P(A) * P(B|A) = P(B) * P(A|B)$$

$$P(\neg A \text{ and } B) = P(\neg A) * P(B|\neg A) = P(B) * P(\neg A|B)$$

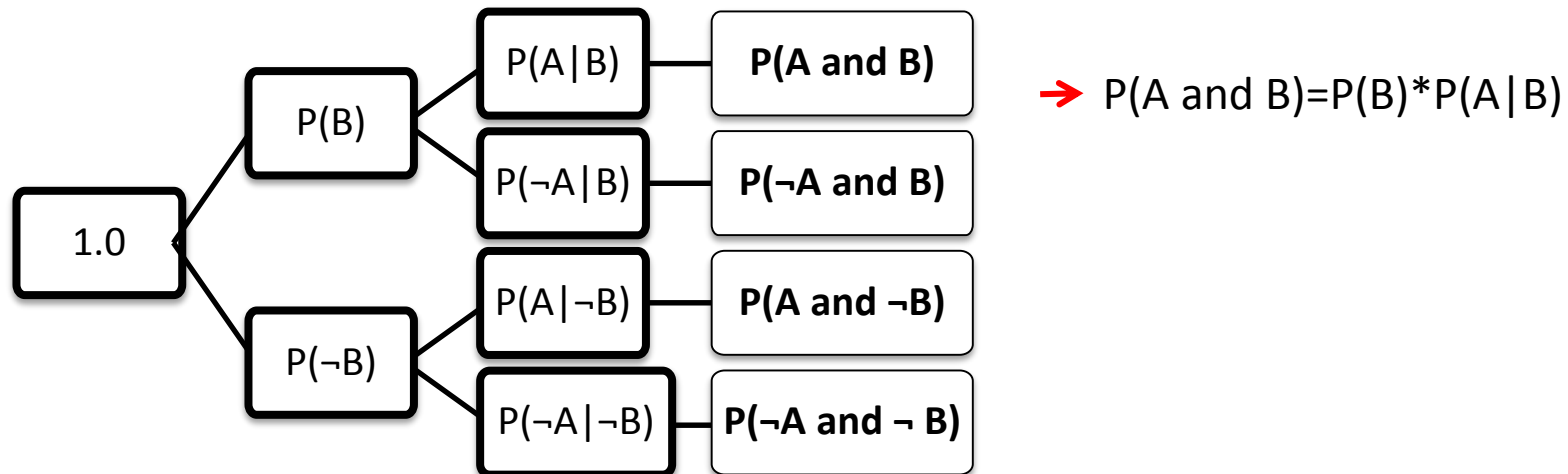
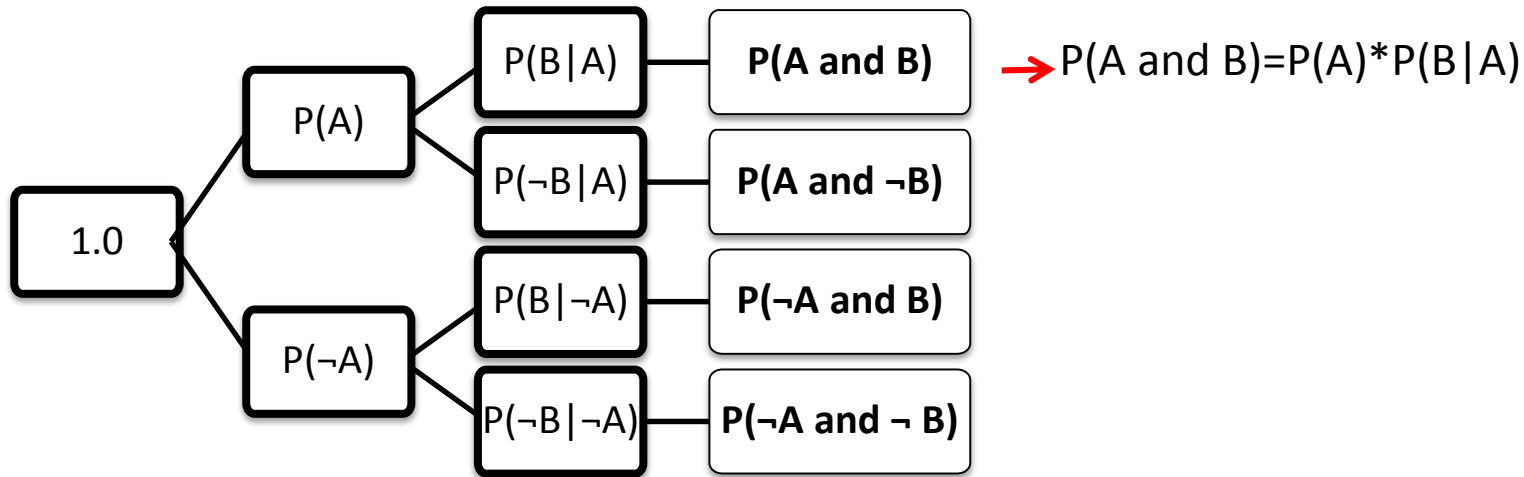
Bayes' theorem



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

$$P(\neg A) * P(B|\neg A) = P(B) * P(\neg A|B)$$

In other words:



$$P(A|B) = P(A) * P(B|A) / P(B)$$

$$P(\neg A|B) = P(\neg A) * P(B|\neg A) / P(B)$$

Bayes' Rule for updating beliefs

$$P(A | B) = P(A) * P(B | A) / P(B)$$

$$P(\neg A | B) = P(\neg A) * P(B | \neg A) / P(B)$$

- We want to compare $P(A | B)$ and $P(\neg A | B)$, i.e. given evidence B what probability is higher: that A occurred or that $\neg A$ occurred?
- We know $P(A)$ and $P(\neg A)$ – prior probabilities
- We know $P(B | A)$ and $P(B | \neg A)$
- From Bayes' theorem:

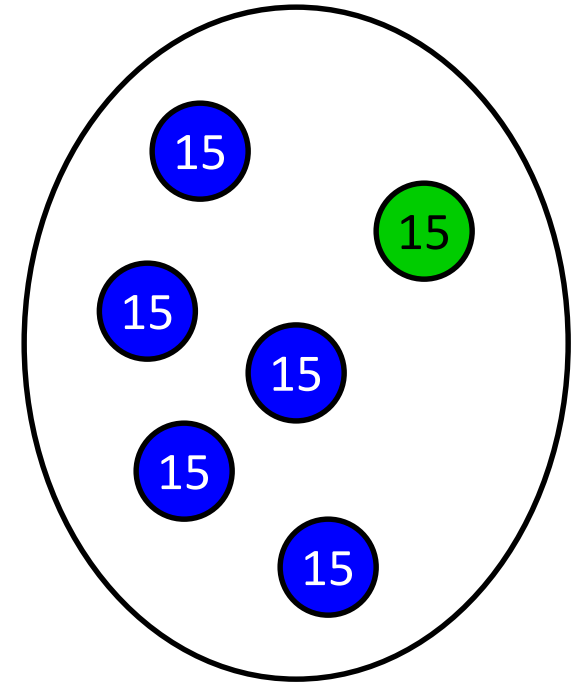
$$P(A | B) = P(A) * P(B | A) / P(B)$$

$$P(\neg A | B) = P(\neg A) * P(B | \neg A) / P(B)$$

Back to hit-and-run

What is more probable: **B** or **G**?

- All cabs were on the streets:
Prior probabilities: $P(\mathbf{B}) = 5/6$, $P(\mathbf{G}) = 1/6$
- The eyewitness test has shown:
 $P(\mathbf{X}_G | \mathbf{G}) = 4/5$ (correctly identified)
 $P(\mathbf{X}_G | \mathbf{B}) = 1/5$ (incorrectly identified)



$$P(\mathbf{G} | \mathbf{X}_G) = P(\mathbf{G}) * P(\mathbf{X}_G | \mathbf{G}) / P(\mathbf{X}_G)$$

$$P(\neg \mathbf{G} | \mathbf{X}_G) = P(\neg \mathbf{G}) * P(\mathbf{X}_G | \neg \mathbf{G}) / P(\mathbf{X}_G)$$

Bayes rule

Hit-and-run: solution

$$P(\mathbf{B}) = 5/6, \quad P(\mathbf{G}) = 1/6$$

$$P(\mathbf{X}_G | \mathbf{G}) = 4/5 \quad P(\mathbf{X}_G | \mathbf{B}) = 1/5$$

- Probability that car was **green** given the evidence X_G :

$$P(\mathbf{G} | \mathbf{X}_G) = P(\mathbf{G}) * P(\mathbf{X}_G | \mathbf{G}) / P(\mathbf{X}_G) = [1/6 * 4/5] / P(\mathbf{X}_G) = 4/30P(\mathbf{X}_G)$$

//- 4 parts of 30P(X_G)

- Probability that car was **blue** given the evidence X_G :

$$P(\mathbf{B} | \mathbf{X}_G) = P(\mathbf{B}) * P(\mathbf{X}_G | \mathbf{B}) / P(\mathbf{X}_G) = [5/6 * 1/5] / P(\mathbf{X}_G) = 6/30P(\mathbf{X}_G)$$

//- 6 parts of 30P(X_G)

6:4 odds that the car was **B**!

Probabilistic classifier

- Given the evidence (data), can we certainly derive the **diagnostic rule**:
if Toothache=true then Cavity=true ?

Name	Toothache	...	Cavity
Smith	true	...	true
Mike	true	...	true
Mary	false	...	true
Quincy	true	...	false
...

- This rule isn't right always.
 - Not all patients with toothache have cavities; some of them have gum disease, an abscess, etc.
- We could try an inverted rule:
if Cavity=true then Toothache=true
- But this rule isn't necessarily right either; not all cavities cause pain.

Certainty and Probability

- The connection between toothaches and cavities is not a certain logical consequence in either direction.
- However, we can provide a **probability** that given an evidence (toothache) the patient has cavity.
- For this we need to know:
 - Prior probability of having cavity: how many times dentist patients had cavities: $P(\text{cavity})$
 - The number of times that the evidence (toothache) was observed among all cavity patients: $P(\text{toothache} | \text{cavity})$

Bayes' Rule

for diagnostic probability

Bayes' rule:

$$P(A | B) = P(A) * P(B | A) / P(B)$$

- Useful for assessing **diagnostic** probability from **symptomatic** probability as:

$$P(\text{Cause} | \text{Symptom}) = P(\text{Symptom} | \text{Cause}) P(\text{Cause}) / P(\text{Symptom})$$

- Bayes's rule is useful in practice because there are many cases where we do have good probability estimates for these three numbers and need to compute the fourth.

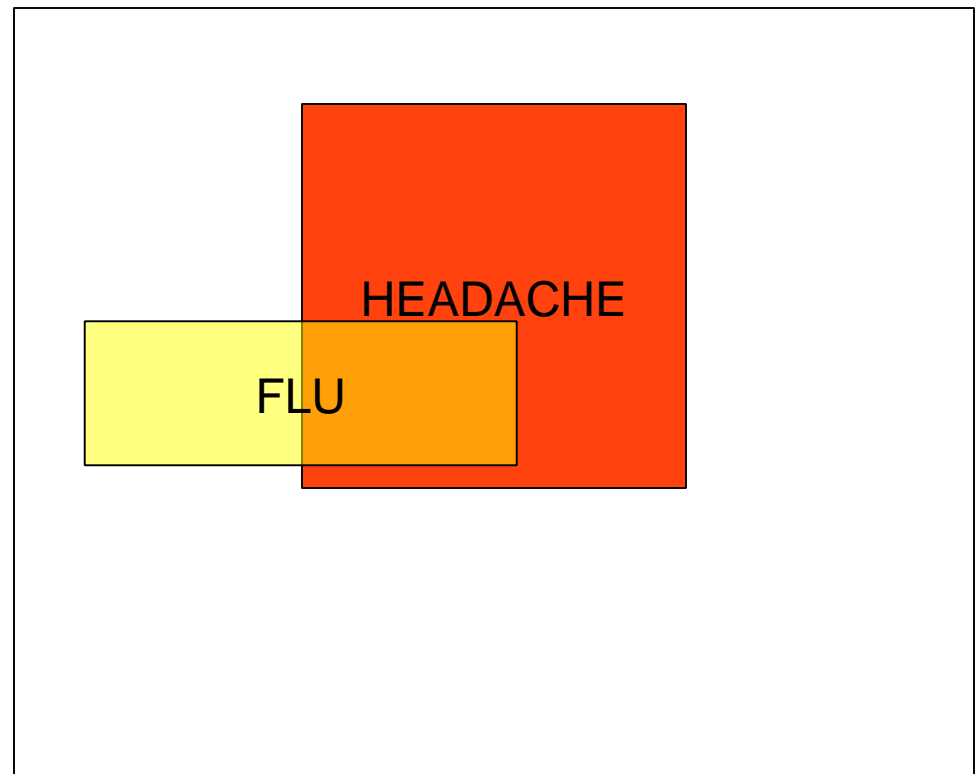
Bayes rule application. Example 1

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

$$P(F) = 1/40$$

$$P(H|F) = 1/2$$

$$P(F|H) = ?$$



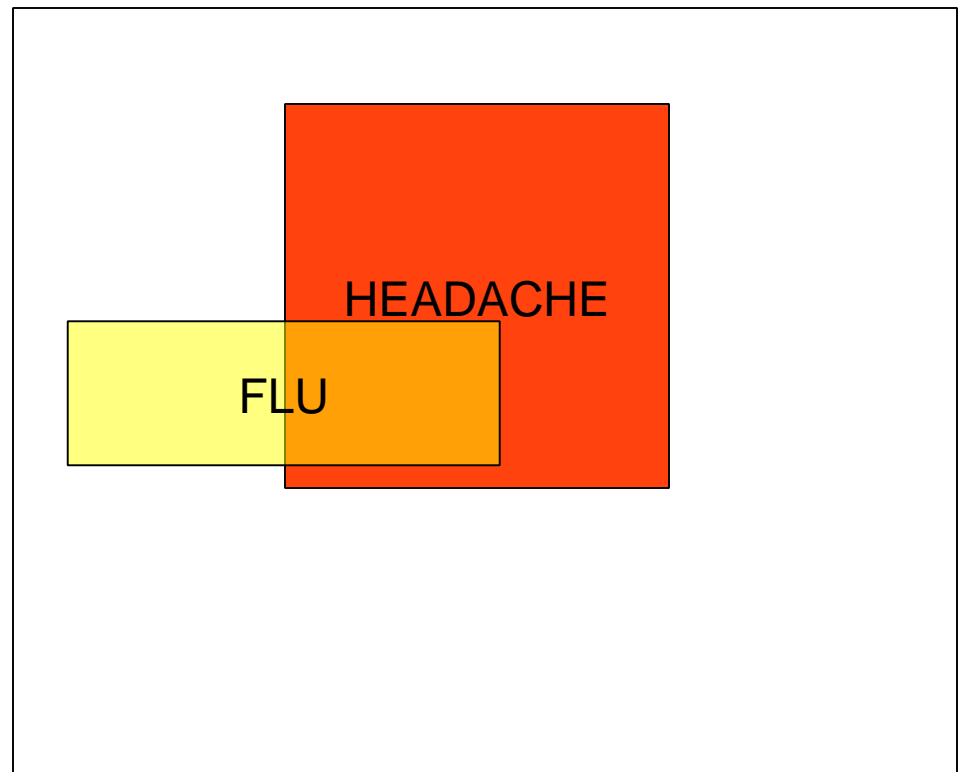
Bayes rule application. Example 1

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

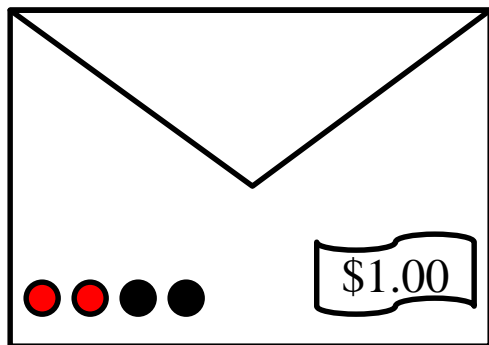
$$P(F)=1/40$$

$$P(H|F)=1/2$$

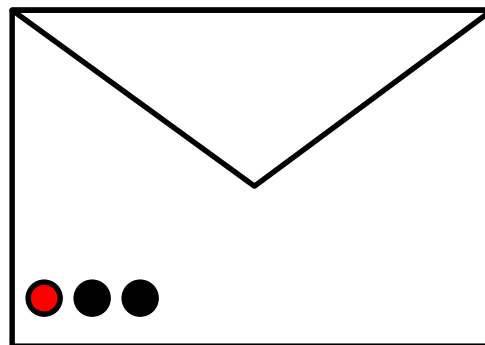
$$P(F|H) = P(H|F)P(F)/P(H)$$
$$= 1/2 * 1/40 * 10 = 1/8$$



Bayes rule application. Example 2



WIN envelope

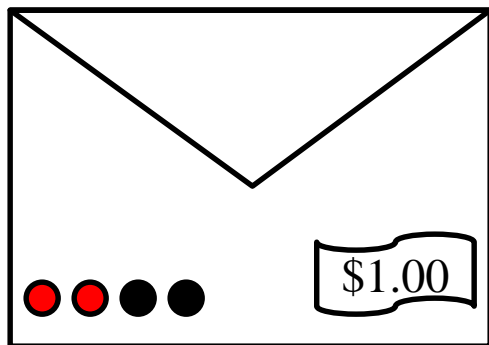


LOSE envelope

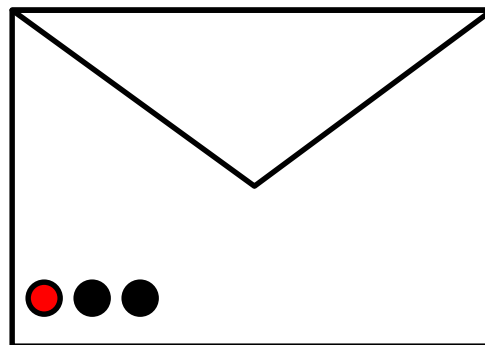
Someone draws an envelope at random and offers to sell it to you.
How much should you pay?

The probability to win is 1:1. Pay no more than 50c.

Bayes rule application. Example 2



WIN envelope



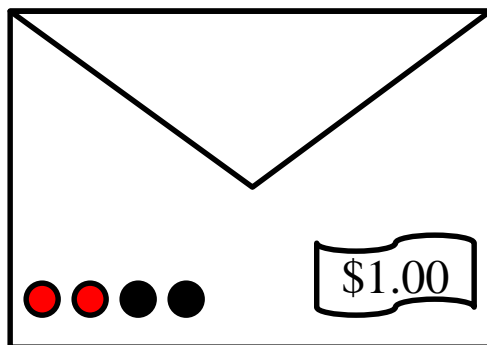
LOSE envelope

Variant: before deciding, you are allowed to see one bead drawn from the envelope.

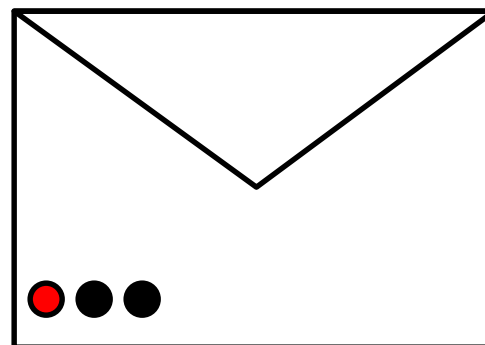
Suppose it's black: How much should you pay?

Suppose it's red: How much should you pay?

Bayes rule application. Example 2



WIN envelope



LOSE envelope

Variant: before deciding, you are allowed to see one bead drawn from the envelope.

Suppose it's black: How much should you pay?

$$P(W|b) = P(b|W)P(W)/P(b) = (1/2 * 1/2)/P(b) = 1/4 * 1/P(b)$$

$$P(L|b) = P(b|L)P(L)/P(b) = (2/3 * 1/2)/P(b) = 1/3 * 1/P(b)$$

Probability to win is now 3:4 – pay not more than $\$(3/7)$

Suppose it's red: How much should you pay? – the same logic

Classifier based on Bayes rule

- We can build a classifier which will classify a new record as class C (yes or no) by comparing probabilities
- In this case all the attributes except C are evidences E
- The data-related task is to evaluate $P(E | C)$ from historical data and based on $P(E | C)$ and prior probabilities $P(C=Yes)$ and $P(C=No)$ compare $P(C=Yes | E)$ and $P(C=No | E)$ using Bayes rule.

Single-evidence classifier: priors

event
(class)

Humidity	Play
High	No
High	No
High	Yes
High	Yes
Normal	Yes
Normal	No
Normal	Yes
High	No
Normal	Yes
Normal	Yes
Normal	Yes
High	Yes
Normal	Yes
High	No

- Prior probabilities:
 $P(\text{Play}=\text{yes})=9/14$, $P(\text{play}=\text{no})=5/14$
- From recording only 'play'/'not play' we have 5:9 odds for play to be canceled today

Single-evidence classifier: evidence

event
evidence (class)

Humidity	Play
High	No
High	No
High	Yes
High	Yes
Normal	Yes
Normal	No
Normal	Yes
High	No
Normal	Yes
Normal	Yes
Normal	Yes
High	Yes
Normal	Yes
High	No

- Priors: $P(\text{Play}=\text{yes})=9/14$, $P(\text{play}=\text{no})=5/14$

- After adding evidence about Humidity we have:

How many times Humidity=normal out of all 9 Yes's: 6

$$P(\text{normal} | \text{yes})=6/9$$

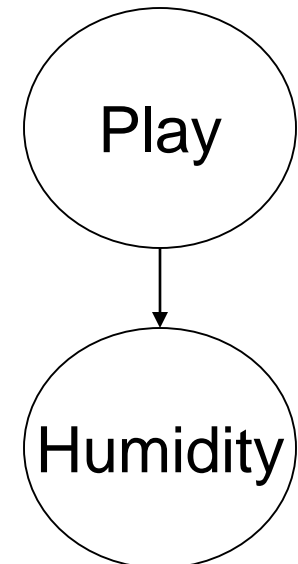
How many times Humidity=normal out of all 5 No's: 1

$$P(\text{normal} | \text{no})=1/5$$

- Similarly:

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

$$P(\text{high} | \text{no})=4/5$$



Single-evidence classifier: prediction

evidence (class)

Humidity	Play
High	No
High	No
High	Yes
High	Yes
Normal	Yes
Normal	No
Normal	Yes
High	No
Normal	Yes
Normal	Yes
Normal	Yes
High	Yes
Normal	Yes
High	No

- $P(\text{yes})=9/14$, $P(\text{no})=5/14$
- $P(\text{high} | \text{yes})=3/9$
- $P(\text{high} | \text{no})=4/5$

Today is a **high** humidity day, what is the probability to play?

- $P(\text{yes} | \text{high})=P(\text{yes}) * P(\text{high} | \text{yes}) / P(\text{high})$
- $P(\text{no} | \text{high})=P(\text{no}) * P(\text{high} | \text{no}) / P(\text{high})$

Single-evidence classifier: prediction

evidence (class)

Humidity	Play
High	No
High	No
High	Yes
High	Yes
Normal	Yes
Normal	No
Normal	Yes
High	No
Normal	Yes
Normal	Yes
Normal	Yes
High	Yes
Normal	Yes
High	No

$$P(\text{yes})=9/14, P(\text{no})=5/14$$

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

$$P(\text{high} | \text{no})=4/5$$

Today is a **high** humidity day, what is the probability to play?

$$P(\text{yes} | \text{high})=P(\text{yes}) * P(\text{high} | \text{yes}) / P(\text{high}) = [9/14 * 3/9] * 1/P(\text{high}) = 3/14 \alpha$$

$$P(\text{no} | \text{high})=P(\text{no}) * P(\text{high} | \text{no}) / P(\text{high}) = [5/14 * 4/5] * 1/P(\text{high}) = 4/14 \alpha$$

4:3 odds not to play given high humidity

(vs. 5:9 before evidence)

Bayes' rule – two evidences

Given that evidence1 is independent of evidence2:

$$P(\text{class} = A | \text{evidence1}, \text{evidence2})$$

$$= \frac{P(\text{evidence1} | \text{class} = A) * P(\text{evidence2} | \text{class} = A) * P(\text{class} = A)}{P(\text{evidence1}) * P(\text{evidence2})}$$

$$= \propto P(\text{evidence1} | \text{class} = A) * P(\text{evidence2} | \text{class} = A) * P(\text{class} = A)$$

The same – let's call it $1/\alpha$

$$P(\text{class} = B | \text{evidence1}, \text{evidence2})$$

$$= \frac{P(\text{evidence1} | \text{class} = B) * P(\text{evidence2} | \text{class} = B) * P(\text{class} = B)}{P(\text{evidence1}) * P(\text{evidence2})}$$

$$= \propto P(\text{evidence1} | \text{class} = B) * P(\text{evidence2} | \text{class} = B) * P(\text{class} = B)$$

Bayes' rule – multiple evidences

Generalized for N evidences

$$P(\text{class} = A | \text{evidence1}, \text{evidence2}, \dots, \text{evidenceN})$$

$$= \frac{P(\text{evidence1} | \text{class} = A) \cdots P(\text{evidenceN} | \text{class} = A) * P(\text{class} = A)}{P(\text{evidence1}) \cdots P(\text{evidenceN})}$$

$$= \propto P(\text{evidence1} | \text{class} = A) * \cdots * P(\text{evidenceN} | \text{class} = A) * P(\text{class} = A)$$

- Two assumptions:
 - Attributes (evidences) are:
 - equally important
 - conditionally independent (given the class value)
- This means that knowledge about the value of a particular attribute doesn't tell us anything about the value of another attribute given the class value

Naïve Bayes classifier

To predict class value for a set of attribute values (evidences)
for each class value compute and compare:

$$\begin{aligned} &P(\text{class} = A | \text{evidence}_1, \text{evidence}_2, \dots, \text{evidence}_N) \\ &= \frac{P(\text{evidence}_1 | \text{class} = A) \cdots P(\text{evidence}_N | \text{class} = A) * P(\text{class} = A)}{P(\text{evidence}_1) \cdots P(\text{evidence}_N)} \\ &= \propto P(\text{evidence}_1 | \text{class} = A) * \cdots * P(\text{evidence}_N | \text{class} = A) * P(\text{class} = A) \end{aligned}$$

- Naïve – assumes independence of variables
- Although based on assumptions that are almost never correct, this scheme works well in practice!

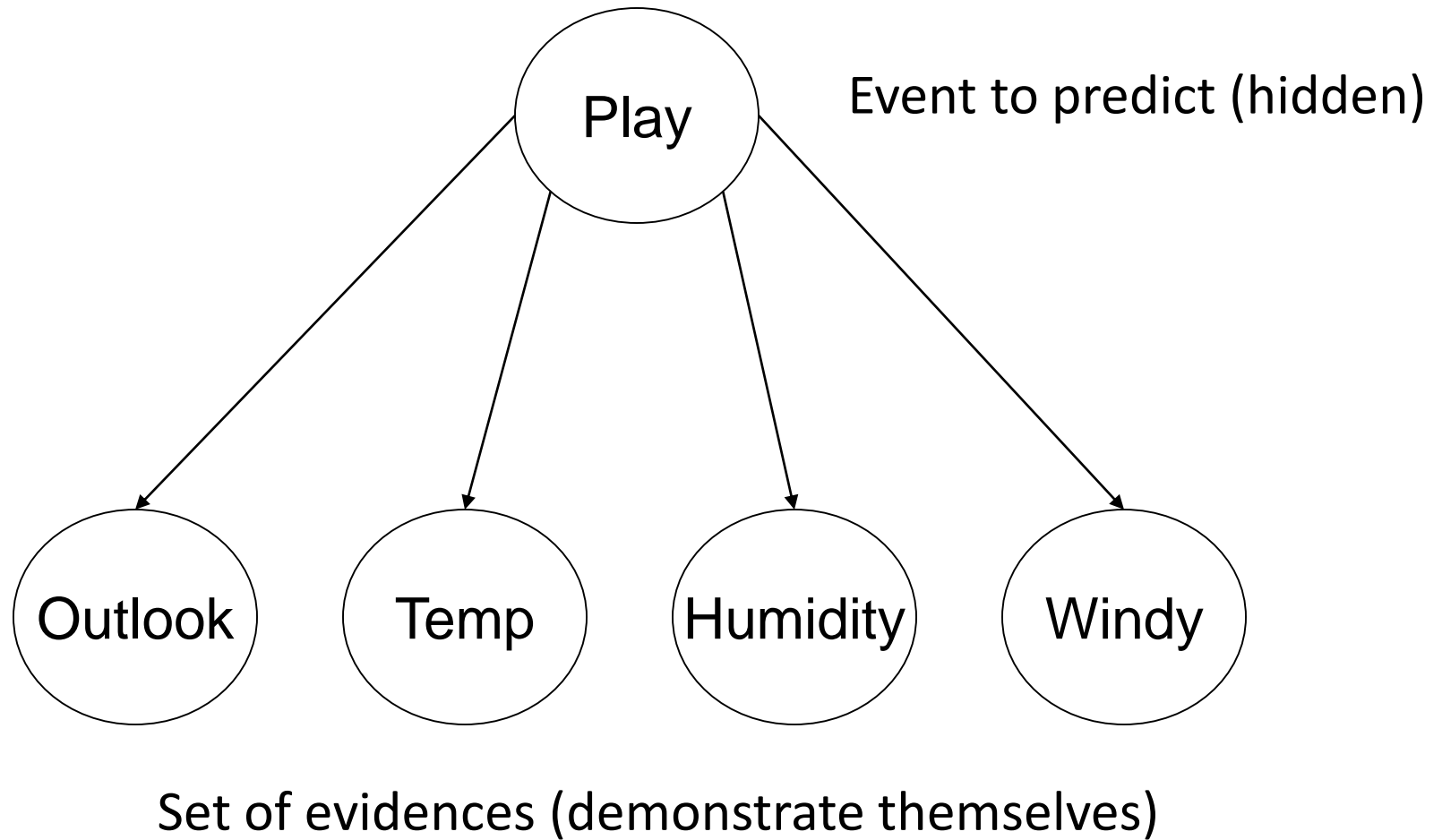
The weather data example

Outlook	Temp.	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

■ A new day:

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

Multi-evidence classifier



The weather data example: probabilities

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

← *Evidence E*

Play	Sunny	Cool	High humidity	Windy=true
Yes: 9	2/9	3/9	3/9	3/9
No: 5	3/5	1/5	4/5	3/5
Total	5	4	7	6

Outlook	Temp.	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

The weather data example: yes

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

← *Evidence E*

$$P(\text{yes} \mid E) =$$

$$P(\text{Sunny} \mid \text{yes}) *$$

$$P(\text{Cool} \mid \text{yes}) *$$

$$P(\text{Humidity}=\text{High} \mid \text{yes}) *$$

$$P(\text{Windy}=\text{True} \mid \text{yes}) *$$

$$P(\text{yes}) / P(E) =$$

$$= (2/9) *$$

$$(3/9) *$$

$$(3/9) *$$

$$(3/9) *$$

$$(9/14) / P(E) = 0.0053 / P(E)$$

Play	Sunny	Cool	High humidity	Windy=true
Yes: 9	2/9	3/9	3/9	3/9
No: 5	3/5	1/5	4/5	3/5
Total	5	4	7	6

Don't worry for the $1/P(E)$; It's alpha, the normalization constant.

The weather data example: no

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

← *Evidence E*

$$P(\text{no} \mid E) =$$

$$P(\text{Sunny} \mid \text{no}) *$$

$$P(\text{Cool} \mid \text{no}) *$$

$$P(\text{Humidity}=\text{High} \mid \text{no}) *$$

$$P(\text{Windy}=\text{True} \mid \text{no}) *$$

$$P(\text{no}) / P(E) =$$

$$= (3/5) *$$

$$(1/5) *$$

$$(4/5) *$$

$$(3/5) *$$

$$(5/14) / P(E) = 0.0206 / P(E)$$

Play	Sunny	Cool	High humidity	Windy=true
Yes: 9	2/9	3/9	3/9	3/9
No: 5	3/5	1/5	4/5	3/5
Total	5	4	7	6

The weather data example: decision

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

← *Evidence E*

$$P(\text{yes} \mid E) = 0.0053 / P(E)$$

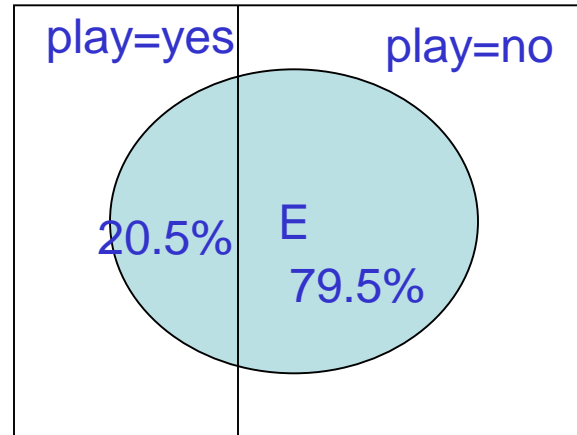
$$P(\text{no} \mid E) = 0.0206 / P(E)$$

More probable: no.

It would be nice to give the actual probability estimates

Outlook	Temp.	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Normalization constant $1/P(E)$



$$P(\text{play=yes} \mid E) + P(\text{play=no} \mid E) = 1 \quad \text{i.e.}$$

$$0.0053 / P(E) + 0.0206 / P(E) = 1 \quad \text{i.e.}$$

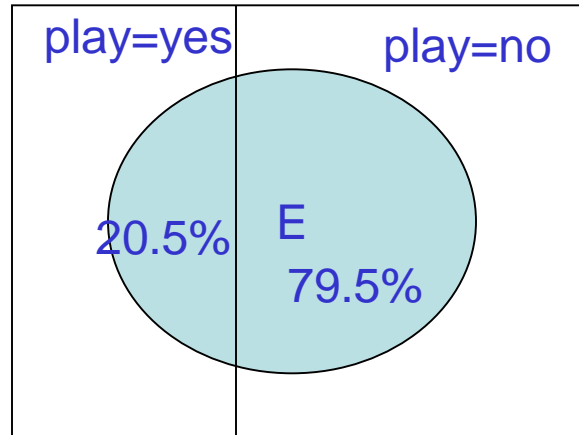
$$P(E) = 0.0053 + 0.0206$$

So,

$$P(\text{play=yes} \mid E) = 0.0053 / (0.0053 + 0.0206) = \mathbf{20.5\%}$$

$$P(\text{play=no} \mid E) = 0.0206 / (0.0053 + 0.0206) = \mathbf{79.5\%}$$

In other words:



$$P(\text{play=yes} \mid E) + P(\text{play=no} \mid E) = 1$$

$$P(\text{play=yes} \mid E) / P(\text{play=no} \mid E) = 0.0053 : 0.0206 = 0.26$$

$$0.26 * P(\text{play=no} \mid E) + P(\text{play=no} \mid E) = 1$$

$$P(\text{play=no} \mid E) = 1/1.26 = 79\%$$

The remaining goes to yes: $P(\text{play=yes} \mid E) = 21\%$

Naïve Bayes: issues

1. Zero frequency problem
2. Missing values
3. Numeric attributes

1. The “zero-frequency problem”

- What if an attribute value doesn't occur with every class value (e.g. “Humidity = High” for class “Play=Yes”)?
 - Probability $P(\text{Humidity}=\text{High} \mid \text{play}=\text{yes})$ will be zero.
- $P(\text{Play}=\text{“Yes”} \mid E)$ will also be zero!
 - No matter how likely the other values are!
- Remedy – Laplace correction:
 - Add **1** to the count for every attribute value-class combination (Laplace estimator);
 - Add k (# of possible attribute values) to the denominator.

Laplace correction

Outlook	Play	Count
Sunny	No	0
Sunny	Yes	6
Overcast	No	2
Overcast	Yes	2
Rainy	No	3
Rainy	Yes	1

+1
→

Outlook	Play	Count
Sunny	No	1
Sunny	Yes	7
Overcast	No	3
Overcast	Yes	3
Rainy	No	4
Rainy	Yes	2

It was: out of total 9 'Yes'

6 – Sunny, 2 – Overcast, 1 – Rainy

The probabilities were:

$P(\text{Sunny} | \text{yes}) = 6/9$; $P(\text{Overcast} | \text{yes}) = 2/9$; $P(\text{Rainy} | \text{yes}) = 1/9$

After correction:

7 – Sunny, 3 – Overcast, 2 – Rainy: Total 'Yes': $9+3=12$

(hence add the cardinality of the attribute to the denominator)

Laplace correction

Outlook	Play	Count
Sunny	No	0
Sunny	Yes	6
Overcast	No	2
Overcast	Yes	2
Rainy	No	3
Rainy	Yes	1

+1
→

Outlook	Play	Count
Sunny	No	1
Sunny	Yes	7
Overcast	No	3
Overcast	Yes	3
Rainy	No	4
Rainy	Yes	2

The probabilities were:

$$P(\text{Sunny} \mid \text{yes}) = 6/9; \quad P(\text{Overcast} \mid \text{yes}) = 2/9; \quad P(\text{Rainy} \mid \text{yes}) = 1/9$$

After correction the probabilities:

$$P(\text{Sunny} \mid \text{yes}) = 7/(9+3);$$

$$P(\text{Overcast} \mid \text{yes}) = 3/(9+3);$$

$$P(\text{Rainy} \mid \text{yes}) = 2/(9+3)$$

} Needs to sum up to 1.0

Laplace correction example

$$P(\text{yes} | E) =$$

$$P(\text{Outlook}=\text{Sunny} | \text{yes}) *$$

$$P(\text{Temp}=\text{Cool} | \text{yes}) *$$

$$P(\text{Humidity}=\text{High} | \text{yes}) *$$

$$P(\text{Windy}=\text{True} | \text{yes}) *$$

$$P(\text{yes}) / P(E) =$$

$$= (2/9) * (3/9) * (3/9) * (3/9) * (9/14) / P(E) = 0.0053 / P(E)$$

With Laplace correction:

Number of possible values for 'Outlook'

$$= ((2+1)/(9+3)) * ((3+1)/(9+3)) * ((3+1)/(9+2)) * ((3+1)/(9+2)) * (9/14) / P(E)$$
$$= 0.007 / P(E)$$

Number of possible values for 'Windy'

2. Missing values: in the training set

- Missing values - not a problem for Naïve Bayes
- Suppose 1 value for outlook in the training set is missing. We count only existing values. For a large dataset, the probability $P(\text{outlook}=\text{sunny}|\text{yes})$ and $P(\text{outlook}=\text{sunny}|\text{no})$ will not change much. This is because we use probabilities rather than absolute counts.

2. Missing values: in the evidence set

- The same calculation without one fraction

Outlook	Temp.	Humidity	Windy	Play
?	Cool	High	True	?

$$P(\text{yes} \mid E) =$$

$$P(\text{Temp}=\text{Cool} \mid \text{yes}) *$$

$$P(\text{Humidity}=\text{High} \mid \text{yes}) *$$

$$P(\text{Windy}=\text{True} \mid \text{yes}) *$$

$$P(\text{yes}) / P(E) =$$

$$= (3/9) * (3/9) * (3/9) * (9/14) / P(E) = 0.0238 / P(E)$$

$$P(\text{no} \mid E) =$$

$$P(\text{Temp}=\text{Cool} \mid \text{no}) *$$

$$P(\text{Humidity}=\text{High} \mid \text{no}) *$$

$$P(\text{Windy}=\text{True} \mid \text{no}) *$$

$$P(\text{play}=\text{no}) / P(E) =$$

$$= (1/5) * (4/5) * (3/5) * (5/14) / P(E) = 0.0343 / P(E)$$

2. Missing values: in the evidence set

- With missing value:

Outlook	Temp.	Humidity	Windy	Play
?	Cool	High	True	?

$$P(\text{yes} \mid E) = 0.0238 / P(E)$$

$$P(\text{no} \mid E) = 0.0343 / P(E)$$

- Without missing value:

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

$$P(\text{yes} \mid E) = 0.0053 / P(E)$$

$$P(\text{no} \mid E) = 0.0206 / P(E)$$

The numbers are much higher for the case of missing values. But we care only about the ratio of yes and no.

2. Missing values: in the evidence set

- With missing value:

Outlook	Temp.	Humidity	Windy	Play
?	Cool	High	True	?

$$P(\text{yes} \mid E) = 0.0238 / P(E)$$

$$P(\text{no} \mid E) = 0.0343 / P(E)$$

After normalization: $P(\text{yes} \mid E) = \mathbf{41\%}$, $P(\text{no} \mid E) = \mathbf{59\%}$

- Without missing value:

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

$$P(\text{yes} \mid E) = 0.0053 / P(E)$$

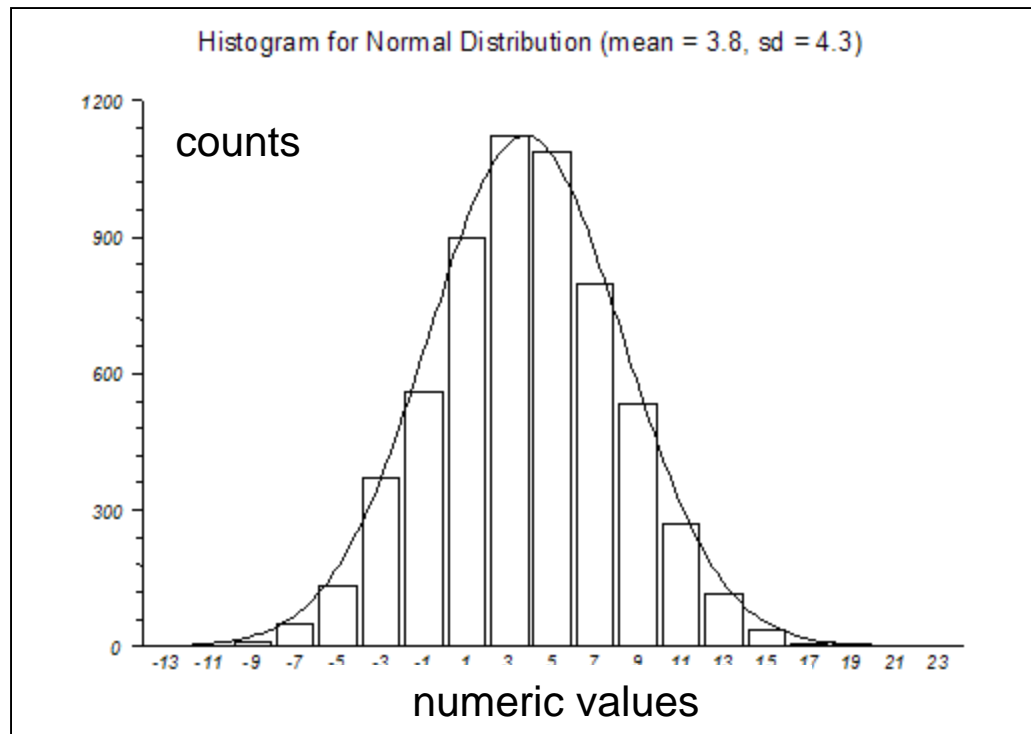
$$P(\text{no} \mid E) = 0.0206 / P(E)$$

After normalization: $P(\text{yes} \mid E) = \mathbf{21\%}$, $P(\text{no} \mid E) = \mathbf{79\%}$

Of course, this is a very small dataset where each count matters, but the prediction is still the same: most probably – no play

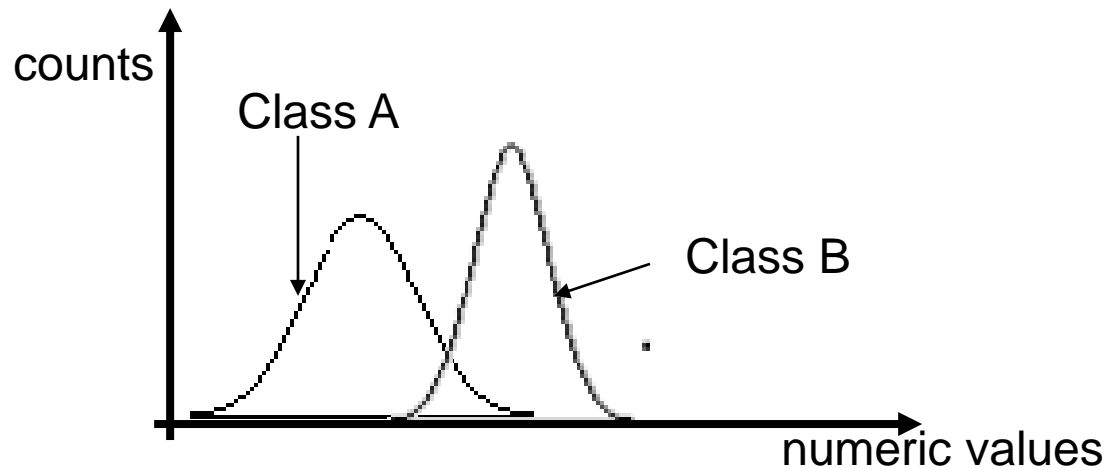
Normal distribution

- Usual assumption: attributes have a normal or Gaussian probability distribution.



Two classes have different distributions

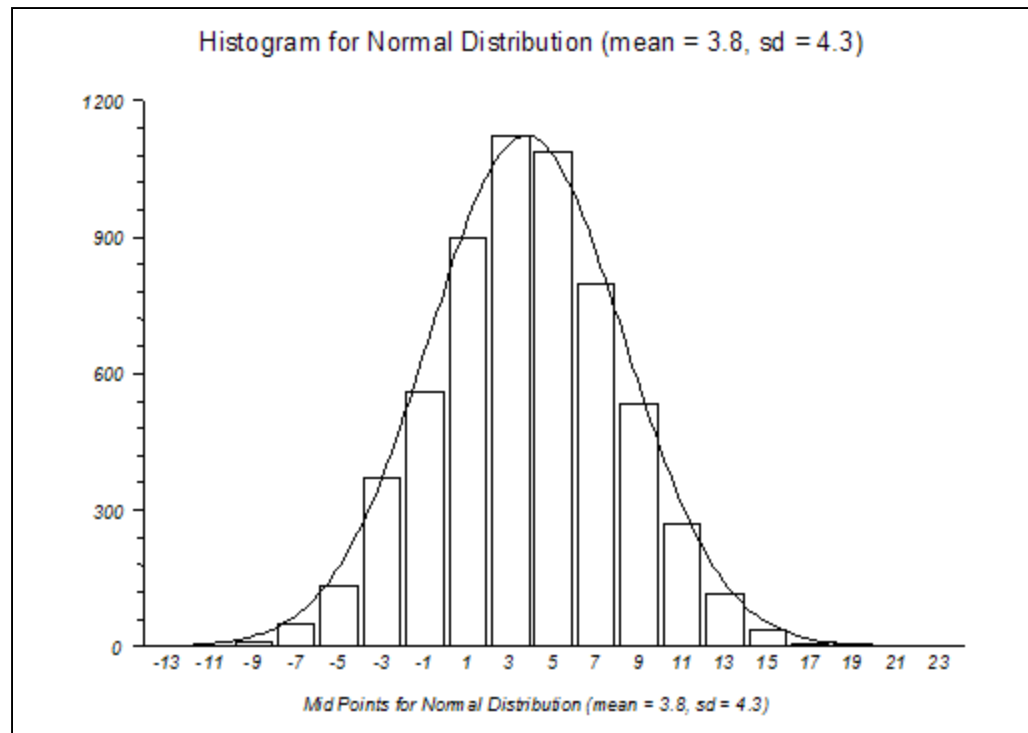
- Class A is normally distributed around its mean with its standard deviation. Class B is normally distributed around the different mean and with a different std



Probability density function

- Probability density function (PDF) for the normal distribution:

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



For a given x – evaluates its probability according to the distribution of probabilities in a given class

Probability and density

- Relationship between probability and density:

$$\Pr\left[c - \frac{\varepsilon}{2} < x < c + \frac{\varepsilon}{2}\right] \approx \varepsilon * f(c)$$

- But: to compare posteriori probabilities it is enough to calculate PDF, because ε cancels out
- Exact relationship:

$$\Pr[a \leq x \leq b] = \int_a^b f(t) dt$$

To compute probability $P(X=V \mid \text{class})$

- Gives \approx probability of $X=V$ of belonging to class A:

$$f(x \mid \text{class}) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- We approximate μ by the sample mean:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

- We approximate σ^2 by the sample variance:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Numeric weather data example

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

$$f(x | yes) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Compute the probability of temp=66 for class Yes:

$\sim\mu$ (mean) =
 $(83+70+68+64+69+75+75+72+81)/9 = 73$

$\sim\sigma^2$ (variance) = $((83-73)^2 + (70-73)^2 + (68-73)^2 + (64-73)^2 + (69-73)^2 + (75-73)^2 + (75-73)^2 + (72-73)^2 + (81-73)^2) / (9-1) = 38$

$$f(x | yes) = \frac{1}{\sqrt{38 * 2 * 3.14}} 2.7^{-\frac{(x-73)^2}{2 * 38}}$$

Density function for temp in class Yes

outlook	temperature	humidity	windy	play
sunny	85	85	FALSE	no
sunny	80	90	TRUE	no
overcast	83	86	FALSE	yes
rainy	70	96	FALSE	yes
rainy	68	80	FALSE	yes
rainy	65	70	TRUE	no
overcast	64	65	TRUE	yes
sunny	72	95	FALSE	no
sunny	69	70	FALSE	yes
rainy	75	80	FALSE	yes
sunny	75	70	TRUE	yes
overcast	72	90	TRUE	yes
overcast	81	75	FALSE	yes
rainy	71	91	TRUE	no

Substitute x=66:

$$f(x = 66 | yes) = \frac{1}{15.44} 2.7^{-\frac{(66-73)^2}{76}} = 0.034$$

P(temp=66|yes)=0.034

Numeric weather data example

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

$$f(x | yes) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Compute the probability of Humidity=90 for class Yes:

$\sim\mu$ (mean) =
 $(86+96+80+65+70+80+70+90+75)/9 = 79$

$\sim\sigma^2$ (variance) = $((86-79)^2 + (96-79)^2 + (80-79)^2 + (65-79)^2 + (70-79)^2 + (80-79)^2 + (70-79)^2 + (90-79)^2 + (75-79)^2) / (9-1) = 104$

$$f(x | yes) = \frac{1}{\sqrt{104 * 2 * 3.14}} 2.7^{-\frac{(x-79)^2}{2*104}}$$

Density function for humidity in class Yes

outlook	temperature	humidity	windy	play
sunny	85	85	FALSE	no
sunny	80	90	TRUE	no
overcast	83	86	FALSE	yes
rainy	70	96	FALSE	yes
rainy	68	80	FALSE	yes
rainy	65	70	TRUE	no
overcast	64	65	TRUE	yes
sunny	72	95	FALSE	no
sunny	69	70	FALSE	yes
rainy	75	80	FALSE	yes
sunny	75	70	TRUE	yes
overcast	72	90	TRUE	yes
overcast	81	75	FALSE	yes
rainy	71	91	TRUE	no

Substitute x=90:

$$f(x = 90 | yes) = \frac{1}{25.55} 2.7^{-\frac{(90-79)^2}{208}} = 0.022$$

P(humidity=90|yes)=0.022

Classifying a new day

- A new day E:

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

$$\begin{aligned} P(\text{play}=\text{yes} \mid E) &= \\ &P(\text{Outlook}=\text{Sunny} \mid \text{play}=\text{yes}) * \\ &P(\text{Temp}=66 \mid \text{play}=\text{yes}) * \\ &P(\text{Humidity}=90 \mid \text{play}=\text{yes}) * \\ &P(\text{Windy}=\text{True} \mid \text{play}=\text{yes}) * \\ &P(\text{play}=\text{yes}) / P(E) = \\ &= (2/9) * (0.034) * (0.022) * (3/9) \\ &\quad * (9/14) / P(E) = 0.000036 / \\ &P(E) \end{aligned}$$

$$\begin{aligned} P(\text{play}=\text{no} \mid E) &= \\ &P(\text{Outlook}=\text{Sunny} \mid \text{play}=\text{no}) * \\ &P(\text{Temp}=66 \mid \text{play}=\text{no}) * \\ &P(\text{Humidity}=90 \mid \text{play}=\text{no}) * \\ &P(\text{Windy}=\text{True} \mid \text{play}=\text{no}) * \\ &P(\text{play}=\text{no}) / P(E) = \\ &= (3/5) * (0.0291) * (0.038) * (3/5) \\ &\quad * (5/14) / P(E) = 0.000136 / \\ &P(E) \end{aligned}$$

After normalization: $P(\text{play}=\text{yes} \mid E) = 20.9\%$, $P(\text{play}=\text{no} \mid E) = 79.1\%$

Practicality

- Naïve Bayes works surprisingly well (even if independence assumption is clearly violated)
- Because classification doesn't require accurate probability estimates as long as maximum probability is assigned to correct class

Applications of Naïve Bayes

The best classifier for:

- Document classification
- Diagnostics
- Clinical trials
- Assessing risks

Text Categorization

- Text categorization is the task of assigning a given document to one of a fixed set of categories, on the basis of the words it contains.
- The class is the document category, and the evidence variables are the presence or absence of each word in the document.

Text Categorization

- The model consists of the prior probability $P(\text{Category})$ and the conditional probabilities $P(\text{Word}_i \mid \text{Category})$.
- For each category c , $P(\text{Category}=c)$ is estimated as the fraction of all the “training” documents that are of that category.
- Similarly, $P(\text{Word}_i = \text{true} \mid \text{Category} = c)$ is estimated as the fraction of documents of category that contain this word.
- Also, $P(\text{Word}_i = \text{true} \mid \text{Category} = \neg c)$ is estimated as the fraction of documents not of category that contain this word.

Text Categorization (cont'd)

- Now we can use naïve Bayes for classifying a new document with n words:

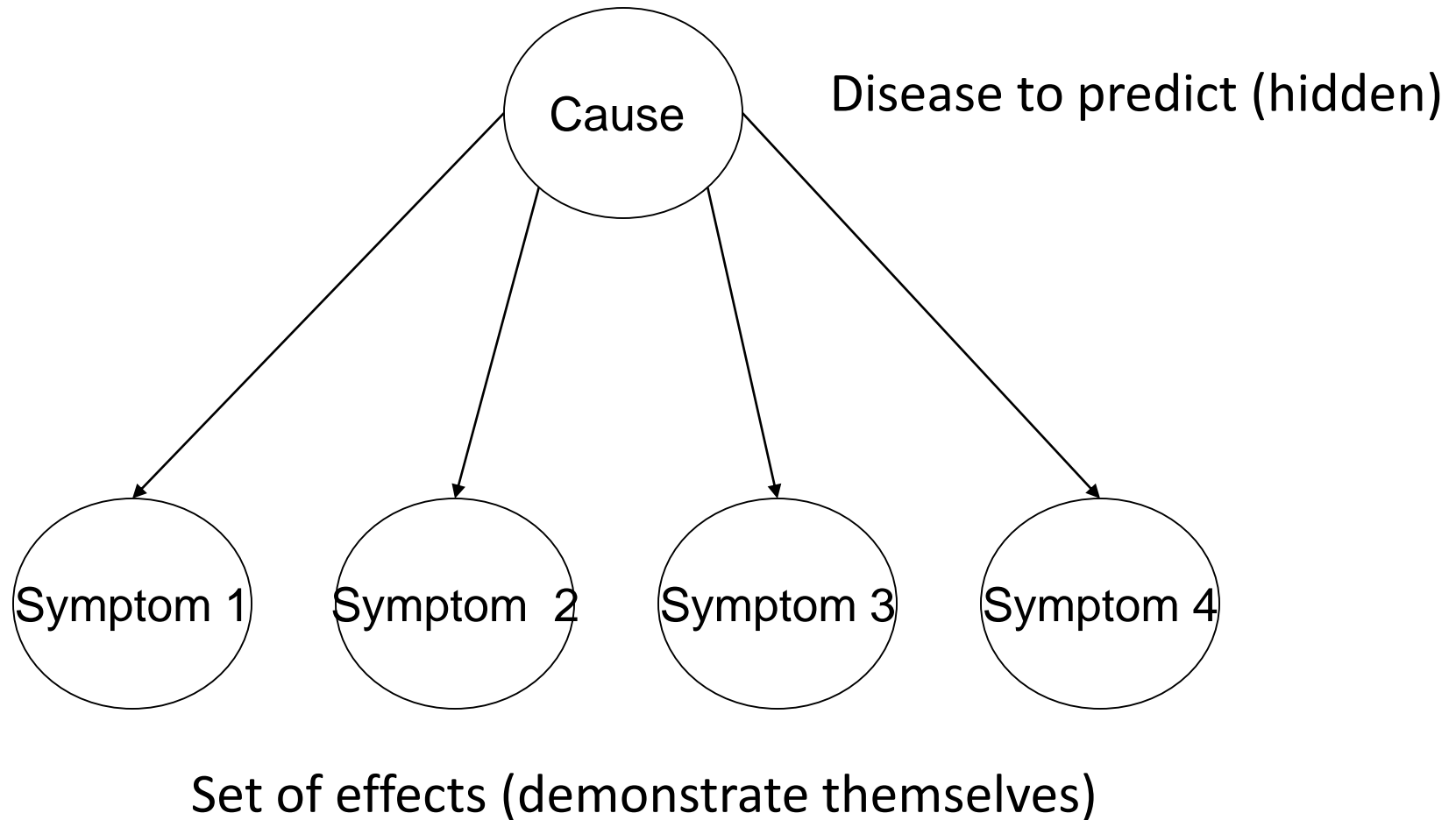
$$P(\text{Category} = c \mid \text{Word}_1 = \text{true}, \dots, \text{Word}_n = \text{true}) = \alpha * P(\text{Category} = c) \prod_{i=1}^n P(\text{Word}_i = \text{true} \mid \text{Category} = c)$$

$$P(\text{Category} = \neg c \mid \text{Word}_1 = \text{true}, \dots, \text{Word}_n = \text{true}) = \alpha * P(\text{Category} = \neg c) \prod_{i=1}^n P(\text{Word}_i = \text{true} \mid \text{Category} = \neg c)$$

$\text{Word}_1, \dots, \text{Word}_n$ are the words occurring in the new document
 α is the normalization constant.

- Observe that similarly with the “missing values” the new document doesn’t contain every word for which we computed the probabilities.

Diagnostics with Naïve Bayes



Example of diagnostic problem

- A doctor knows that **50%** of patients with a stiff neck were diagnosed with meningitis.
- The **doctor also knows some unconditional** facts (prior probabilities):
 - the prior probability that any patient has meningitis is **$1/50,000$**
 - the probability that he does not have a meningitis is **$49,999/50,000$**

Diagnostic problem

$$P(\text{StiffNeck}=\text{true} \mid \text{Meningitis}=\text{true}) = 0.5$$

$$P(\text{StiffNeck}=\text{true} \mid \text{Meningitis}=\text{false}) = 0.5$$

$$P(\text{Meningitis}=\text{true}) = 1/50000$$

$$P(\text{Meningitis}=\text{false}) = 49999/50000$$

$$P(\text{Meningitis}=\text{true} \mid \text{StiffNeck}=\text{true})$$

$$= P(\text{StiffNeck}=\text{true} \mid \text{Meningitis}=\text{true}) P(\text{Meningitis}=\text{true}) /$$

$$P(\text{StiffNeck}=\text{true})$$

$$= (0.5) \times (1/50000) / P(\text{StiffNeck}=\text{true}) = 0.5 * 0.00002 / P(\text{StiffNeck}=\text{true}) =$$

$$0.00010 / P(\text{StiffNeck}=\text{true})$$

$$P(\text{Meningitis}=\text{false} \mid \text{StiffNeck}=\text{true})$$

$$= P(\text{StiffNeck}=\text{true} \mid \text{Meningitis}=\text{false}) P(\text{Meningitis}=\text{false}) /$$

$$P(\text{StiffNeck}=\text{true})$$

$$= (0.5) * (49999/50000) / P(\text{StiffNeck}=\text{true}) = 0.49999 / P(\text{StiffNeck}=\text{true})$$

1/5000 chance that the patient with a stiff neck has meningitis (due to the very low prior probability)

Bayes' rule critics: prior probabilities

- The doctor has the above quantitative information in the diagnostic direction from symptoms (evidences, effects) to causes.
- The problem is that prior probabilities are hard to estimate and they may fluctuate. Imagine, there is sudden epidemic of meningitis. The prior probability, $P(\text{Meningitis}=\text{true})$, will go up.
- Clearly, $P(\text{StiffNeck}=\text{true} \mid \text{Meningitis}=\text{true})$ is unaffected by the epidemic. It simply reflects the way meningitis works.
- The estimation of $P(\text{Meningitis}=\text{true} \mid \text{StiffNeck}=\text{true})$ will be incorrect until new data about $P(\text{Meningitis}=\text{true})$ are collected

Tax Data – Naive Bayes

Classify: (_, No, Married, 95K, ?)

(Apply also the Laplace normalization)

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Tax Data – Naive Bayes

Classify: (_, No, Married, 95K, ?)

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

$$P(\text{Yes}) = 3/10 = 0.3$$

$$P(\text{Refund}=\text{No} | \text{Yes}) = (3+1)/(3+2) = 0.8$$

$$P(\text{Status}=\text{Married} | \text{Yes}) = (0+1)/(3+3) = 0.17$$

$$f(\text{income} | \text{Yes}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Approximate μ with: $(95+85+90)/3 = 90$

Approximate σ^2 with:

$$\frac{((95-90)^2 + (85-90)^2 + (90-90)^2)}{(3-1)} = 25$$

$$f(\text{income}=95 | \text{Yes}) =$$

$$e^{-((95-90)^2 / (2*25))} / \sqrt{2*3.14*25} = .048$$

$$P(\text{Yes} | E) = \alpha * .8 * .17 * .048 * .3 = \alpha * .0019584$$

Tax Data

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Classify: (_, No, Married, 95K, ?)

$$P(\text{No}) = 7/10 = .7$$

$$P(\text{Refund}=\text{No} | \text{No}) = (4+1)/(7+2) = .556$$

$$P(\text{Status}=\text{Married} | \text{No}) = (4+1)/(7+3) = .5$$

$$f(\text{income} | \text{No}) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Approximate μ with:

$$(125+100+70+120+60+220+75)/7 = 110$$

Approximate σ^2 with:

$$((125-110)^2 + (100-110)^2 + (70-110)^2 + (120-110)^2 + (60-110)^2 + (220-110)^2 + (75-110)^2)/(7-1) = 2975$$

$$f(\text{income}=95 | \text{No}) =$$

$$e^{-((95-110)^2 / (2*2975))} / \sqrt{2*3.14*2975} = .00704$$

$$P(\text{No} | E) = \alpha * .556 * .5 * .00704 * 0.7 = \alpha * .00137$$

Tax Data

Classify: (_, No, Married, 95K, ?)

$$P(\text{Yes} | E) = \alpha * .0019584$$

$$P(\text{No} | E) = \alpha * .00137$$

$$\alpha = 1 / (.0019584 + .00137) = 300.44$$

$$P(\text{Yes} | E) = 300.44 * .0019584 = 0.59$$

$$P(\text{No} | E) = 300.44 * .00137 = 0.41$$

We predict "Yes."

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes