

Probability for Learning

- Probability for classification and modeling concepts.
- Bayesian probability
 - Notion of probability interpreted as partial belief
- Bayesian Estimation
 - Calculate the validity of a proposition
 - Based on prior estimate of its probability
 - and New relevant evidence

Bayes Theorem

- **Goal:** To determine the most probable hypothesis, given the data D plus any initial knowledge about the prior probabilities of the various hypotheses in H .

Bayes Theorem

Bayes Rule:
$$P(h | D) = \frac{P(D | h)P(h)}{P(D)}$$

- $P(h)$ = prior probability of hypothesis h
- $P(D)$ = prior probability of training data D
- $P(h | D)$ = probability of h given D (posterior density)
- $P(D | h)$ = probability of D given h (likelihood of D given h)

Dilemma at the movies

This person dropped their ticket in the hallway.

Do you call out

“Excuse me, ma’am!”

or

“Excuse me, sir!”

You have to make a guess.



Dilemma at the movies

What if they're standing in line for the men's restroom?

Bayesian inference is a way to capture common sense.

It helps you use what you know to make better guesses.

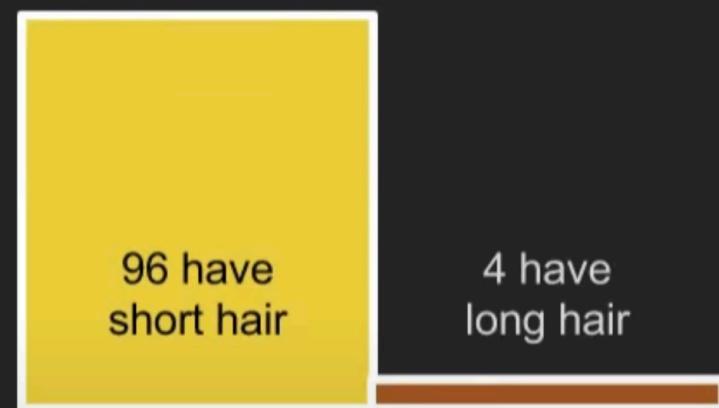


Put numbers to our dilemma

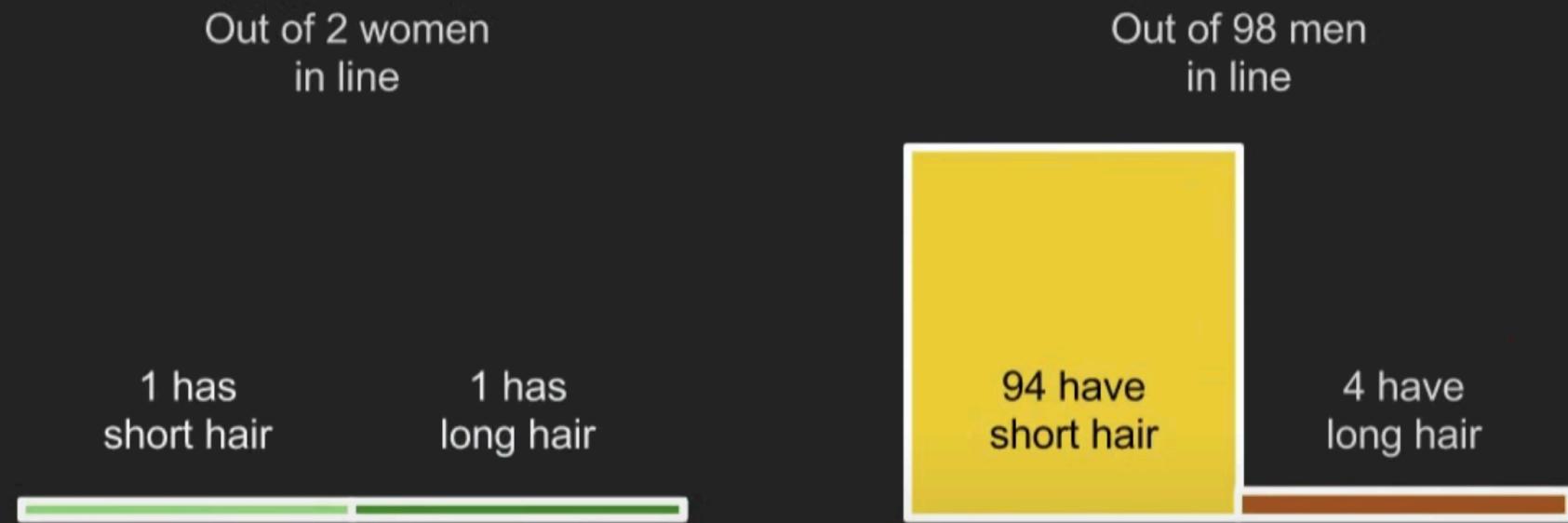
Out of 100 women
at the movies



Out of 100 men
at the movies



Put numbers to our dilemma



But there are 98 men and 2 women in line for the men's restroom.



Out of 100 people
at the movies

50 are women

50 are men

25 women
have
short hair

25 women
have
long hair

48 men
have
short hair

2 men have long hair





2 are
women

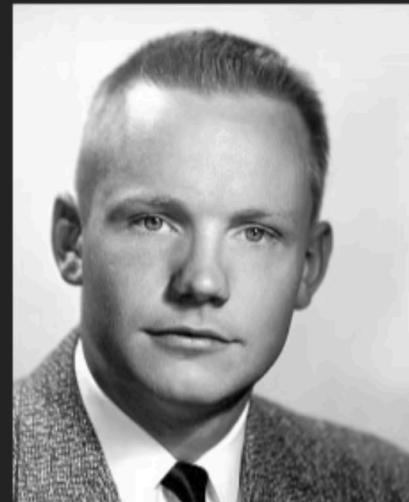
One
woman
has
short
hair

One
woman
has
long
hair

Out of 100 people
In line for the men's
restroom
98 are men

94 men
have
short hair

4 men have long hair



Translate to math

$P(\text{something}) = \# \text{ something} / \# \text{ everything}$

$P(\text{woman}) = \text{Probability that a person is a woman}$

$$= \# \text{ women} / \# \text{ people}$$

$$= 50 / 100 = .5$$

$P(\text{man}) = \text{Probability that a person is a man}$

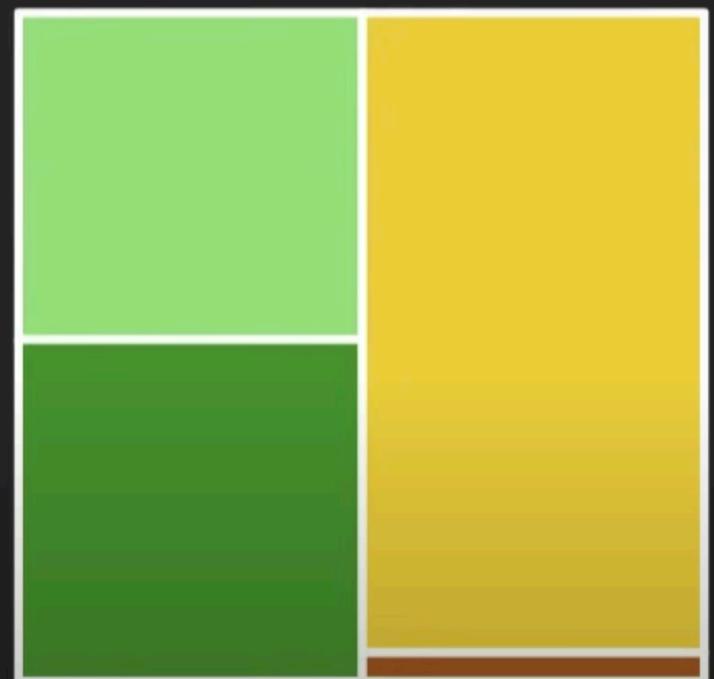
$$= \# \text{ men} / \# \text{ people}$$

$$= 50 / 100 = .5$$

Out of 100 people
at the movies

50 are women

50 are men



Translate to math

$P(\text{something}) = \# \text{ something} / \# \text{ everything}$

2 are
women

Out of 100 people
In line for the men's
restroom

$P(\text{woman}) = \text{Probability that a person is a woman}$

98 are men

$$= \# \text{ women} / \# \text{ people}$$

$$= 2 / 100 = \mathbf{.02}$$

$P(\text{man}) = \text{Probability that a person is a man}$

$$= \# \text{ men} / \# \text{ people}$$

$$= 98 / 100 = \mathbf{.98}$$



Conditional probabilities

$P(\text{long hair} | \text{woman})$

If I know that a person is a woman, what is the probability that person has long hair?

$P(\text{long hair} | \text{woman})$

$$= \# \text{ women with long hair} / \# \text{ women}$$

$$= 25 / 50 = .5$$

Out of 100 people
at the movies

50 are women

25 women
have
short hair

25 women
have
long hair



Conditional probabilities

This doesn't change when we consider people in line.

$$P(\text{long hair} \mid \text{woman})$$

$$= \# \text{ women with long hair} / \# \text{ women}$$

$$= 1 / 2 = .5$$

2 are
women

One
woman
has
short
hair

One
woman
has
long
hair

Out of 100 people
In line for the men's
restroom



Conditional probabilities

If I know that a person is a man, what is the probability that person has long hair?

$$P(\text{long hair} \mid \text{man})$$

$$= \# \text{ men with long hair} / \# \text{ men}$$

$$= 2 / 50 = .04$$

Whether in line or not.

Out of 100 people
at the movies

50 are men



Conditional probabilities

$P(A | B)$ is the probability of A, given B.

“If I know B is the case, what is the probability that A is also the case?”

$P(A | B)$ is not the same as $P(B | A)$.



$P(\text{cute} | \text{puppy})$ is not the same as $P(\text{puppy} | \text{cute})$

If I know the thing I'm holding is a puppy, what is the probability that it is cute?

If I know the the thing I'm holding is cute, what is the probability that it is a puppy?

Joint probabilities

What is the probability that a person is both a woman and has short hair?

$P(\text{woman with short hair})$

$$= P(\text{woman}) * P(\text{short hair} | \text{woman})$$

$$= .5 * .5 = .25$$

Out of probability of 1

$$P(\text{woman}) = .5$$

$$P(\text{man}) = .5$$



Joint probabilities

Out of probability of 1

$P(\text{woman with long hair})$

$$= P(\text{woman}) * P(\text{long hair} | \text{woman})$$

$$= .5 * .5 = \mathbf{.25}$$

$$P(\text{woman}) = .5$$

$$P(\text{man}) = .5$$



Joint probabilities

Out of probability of 1

$P(\text{man with short hair})$

$$= P(\text{man}) * P(\text{short hair} | \text{man})$$

$$= .5 * .96 = .48$$

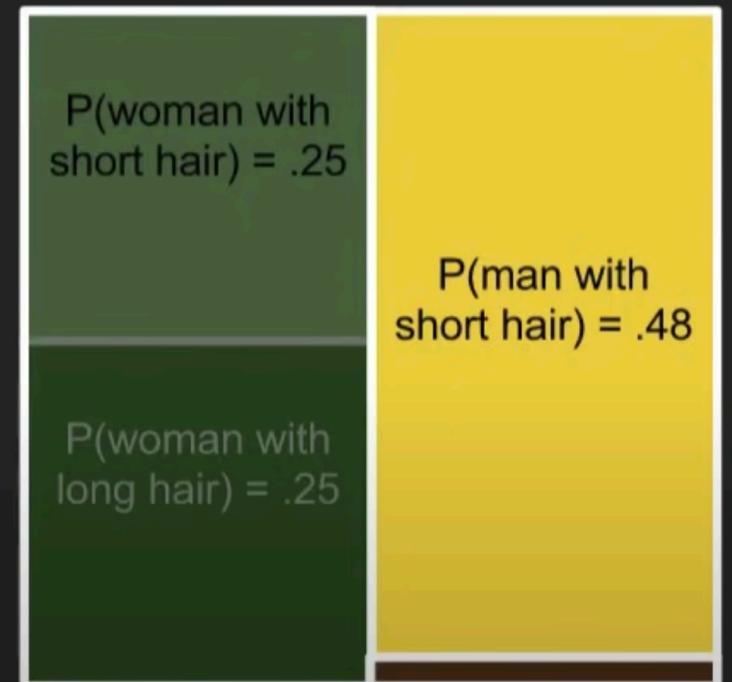
$$P(\text{woman}) = .5$$

$$P(\text{man}) = .5$$

$$P(\text{woman with short hair}) = .25$$

$$P(\text{man with short hair}) = .48$$

$$P(\text{woman with long hair}) = .25$$



Joint probabilities

Out of probability of 1

$P(\text{man with long hair})$

$$= P(\text{man}) * P(\text{long hair} | \text{man})$$

$$= .5 * .04 = \mathbf{.02}$$

$$P(\text{woman}) = .5$$

$$P(\text{man}) = .5$$



Joint probabilities

Out of probability of 1

If $P(\text{man}) = .98$ and $P(\text{woman}) = .02$,
then the answers change.

$P(\text{man with long hair})$

$$\begin{aligned} &= P(\text{man}) * P(\text{long hair} | \text{man}) \\ &= .96 * .04 = \mathbf{.04} \end{aligned}$$

$P(\text{woman}) = .02$

$P(\text{man}) = .98$

$P(\text{woman with short hair}) = .01$

$P(\text{woman with long hair}) = .01$

$P(\text{man with short hair}) = .94$

$P(\text{man with long hair}) = \mathbf{.04}$

Joint probabilities

Out of probability of 1

$P(\text{woman with long hair})$

$P(\text{woman}) = .02$

$P(\text{man}) = .98$

$$= P(\text{woman}) * P(\text{long hair} | \text{woman})$$

$$= .02 * .5 = \mathbf{.01}$$

$P(\text{woman with short hair}) = .01$

$P(\text{man with short hair}) = .94$

$P(\text{woman with long hair}) = .01$

$P(\text{man with long hair}) = .04$

Joint probabilities

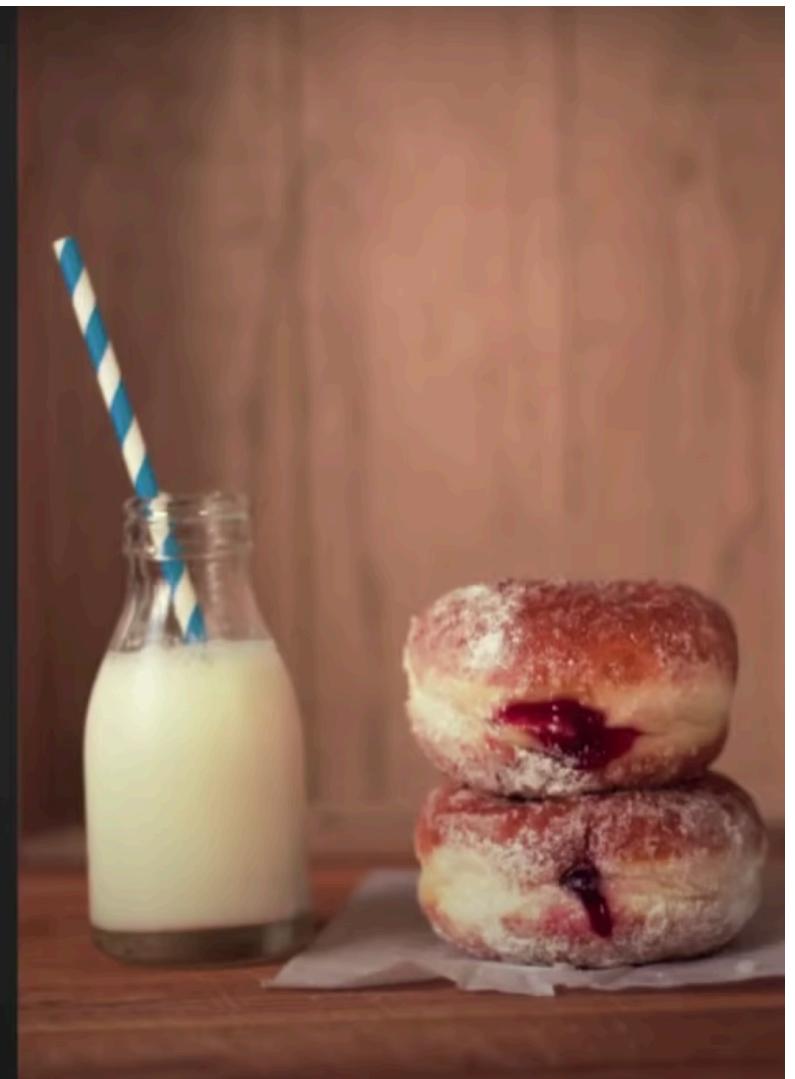
$P(A \text{ and } B)$ is the probability that both A and B are the case.

Also written $P(A, B)$ or $P(A \cap B)$

$P(A \text{ and } B)$ is the same as $P(B \text{ and } A)$

The probability that I am having a jelly donut with my milk is the same as the probability that I am having milk with my jelly donut.

$P(\text{donut and milk}) = P(\text{milk and donut})$



Marginal probabilities

Out of probability of 1

$$P(\text{long hair}) = P(\text{woman with long hair}) +$$

$$P(\text{woman}) = .02$$

$$P(\text{man}) = .98$$

$$P(\text{man with long hair})$$

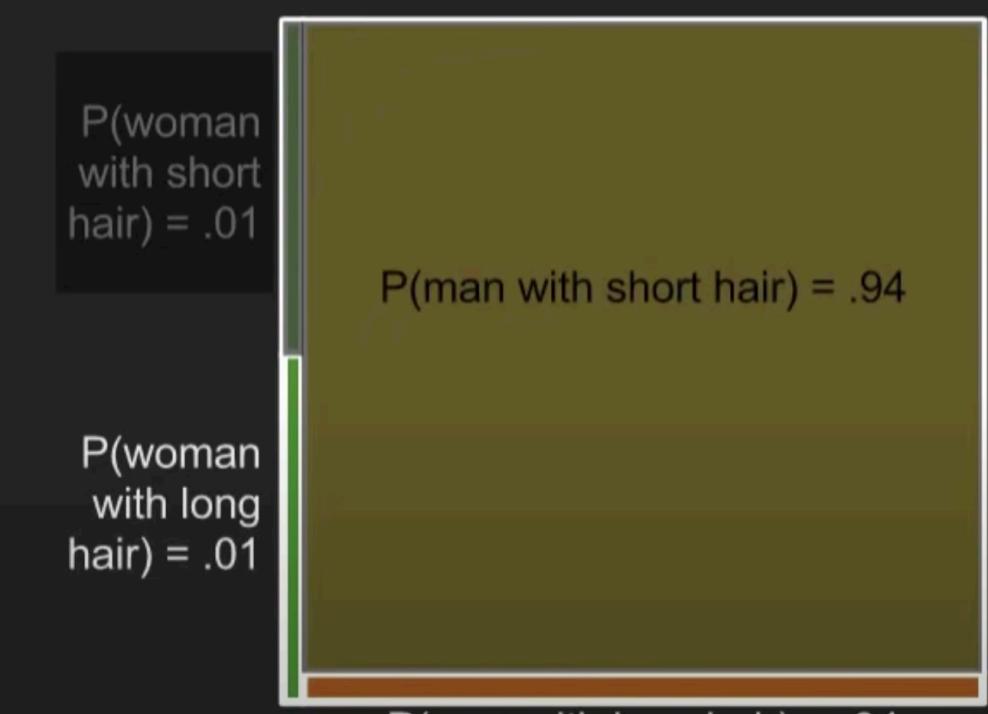
$$= .01 + .04 = \mathbf{.05}$$

$$P(\text{woman with short hair}) = .01$$

$$P(\text{man with short hair}) = .94$$

$$P(\text{woman with long hair}) = .01$$

$$P(\text{man with long hair}) = .04$$



Marginal probabilities

Out of probability of 1

$$P(\text{short hair}) = P(\text{woman with short hair}) + P(\text{man with short hair})$$

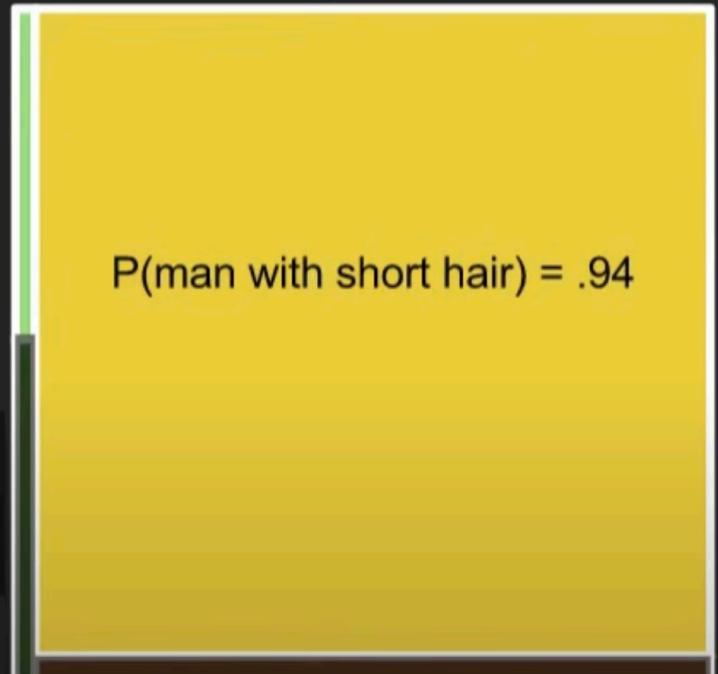
$$P(\text{woman}) = .02$$

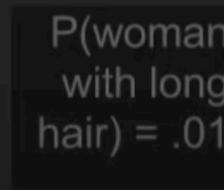
$$P(\text{man}) = .98$$

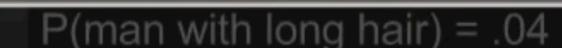
$$P(\text{man with short hair})$$

$$= .01 + .94 = \mathbf{.95}$$

$$P(\text{woman with short hair}) = .01$$


$$P(\text{man with short hair}) = .94$$


$$P(\text{woman with long hair}) = .01$$


$$P(\text{man with long hair}) = .04$$

What we really care about

We know the person has long hair.
Are they a man or a woman?

$P(\text{man} \mid \text{long hair})$

We don't know this answer yet.



Thomas Bayes noticed something cool

$$P(\text{man with long hair}) = P(\text{long hair}) * P(\text{man} | \text{long hair})$$

$$P(\text{long hair and man}) = P(\text{man}) * P(\text{long hair} | \text{man})$$

Because $P(\text{man and long hair}) = P(\text{long hair and man})$

Thomas Bayes noticed something cool

$$P(\text{man with long hair}) = P(\text{long hair}) * P(\text{man} \mid \text{long hair})$$

$$P(\text{long hair and man}) = P(\text{man}) * P(\text{long hair} \mid \text{man})$$

$$\text{Because } P(\text{man and long hair}) = P(\text{long hair and man})$$

$$P(\text{long hair}) * P(\text{man} \mid \text{long hair}) = P(\text{man}) * P(\text{long hair} \mid \text{man})$$

$$P(\text{man} \mid \text{long hair}) = P(\text{man}) * P(\text{long hair} \mid \text{man}) / P(\text{long hair})$$

$$P(A \mid B) = P(B \mid A) * P(A) / P(B)$$

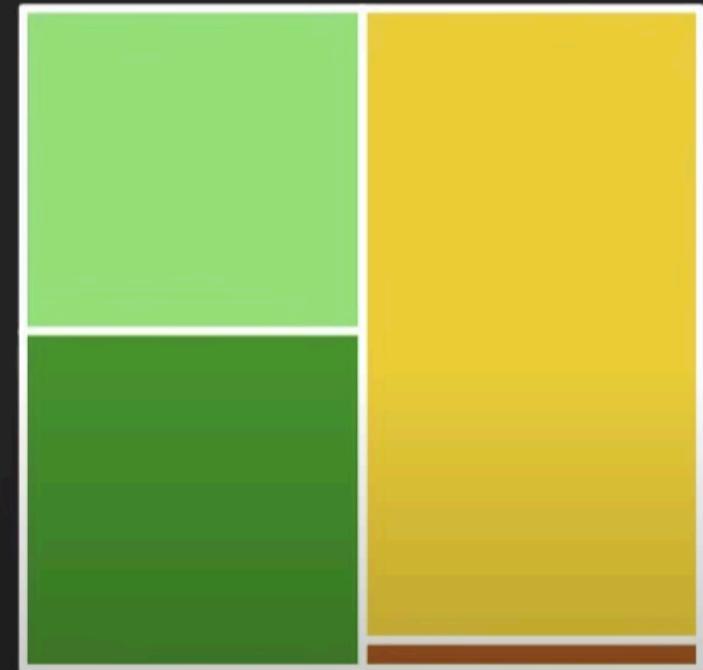
Back to the movie theater, this time with Bayes

$$P(\text{man} \mid \text{long hair}) = \frac{P(\text{man}) * P(\text{long hair} \mid \text{man})}{P(\text{long hair})}$$

$$= \frac{P(\text{man}) * P(\text{long hair} \mid \text{man})}{P(\text{woman with long hair}) + P(\text{man with long hair})}$$

$$P(\text{man} \mid \text{long hair}) = \frac{.5 * .04}{.25 + .02} = .02 / .27 = .07$$

$$P(\text{woman}) = .5 \quad P(\text{man}) = .5$$



$$P(\text{long hair} \mid \text{man}) = .04 \\ P(\text{long hair} \mid \text{woman}) = .5$$

Back to the movie theater, this time with Bayes

$$P(\text{man} \mid \text{long hair}) = \frac{P(\text{man}) * P(\text{long hair} \mid \text{man})}{P(\text{long hair})}$$

$$= \frac{P(\text{man}) * P(\text{long hair} \mid \text{man})}{P(\text{woman with long hair}) + P(\text{man with long hair})}$$

$$P(\text{man} \mid \text{long hair}) = \frac{.98 * .04}{.01 + .04} = .04 / .05 = .80$$



$$P(\text{long hair} \mid \text{man}) = .04$$

$$P(\text{long hair} \mid \text{woman}) = .5$$

Probability distributions

Probability is like a pot with just one cup of coffee left in it.



Probability distributions

If you only have one cup, you can fill it completely.



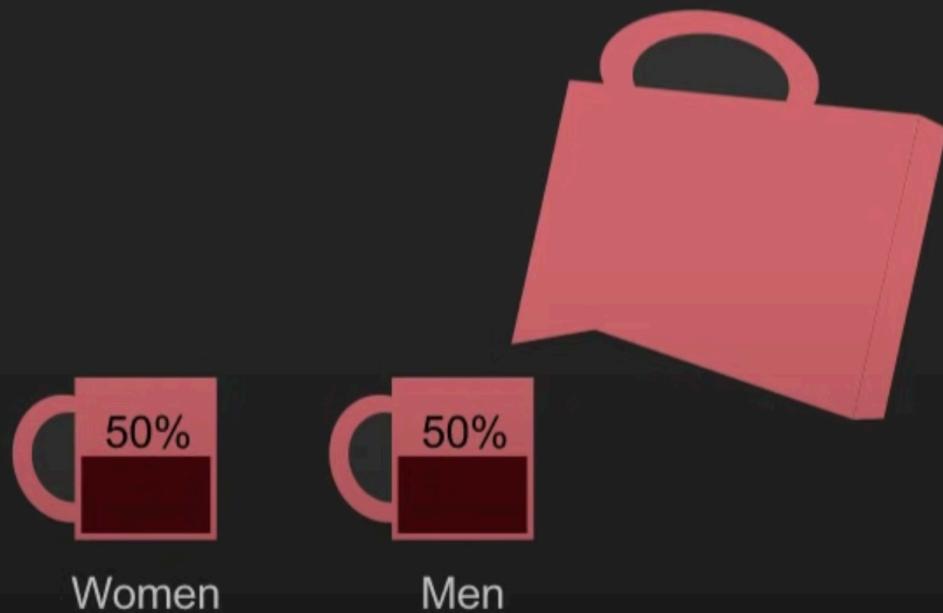
Probability distributions

If you have two cups, you have to decide how to share (distribute) it.



Probability distributions

Our people are distributed between two groups, women and men.



Probability distributions

We can distribute them more.



Women with
short hair



Men with
short hair



Women with
long hair



Men with
long hair

Probability distributions



Women with
short hair



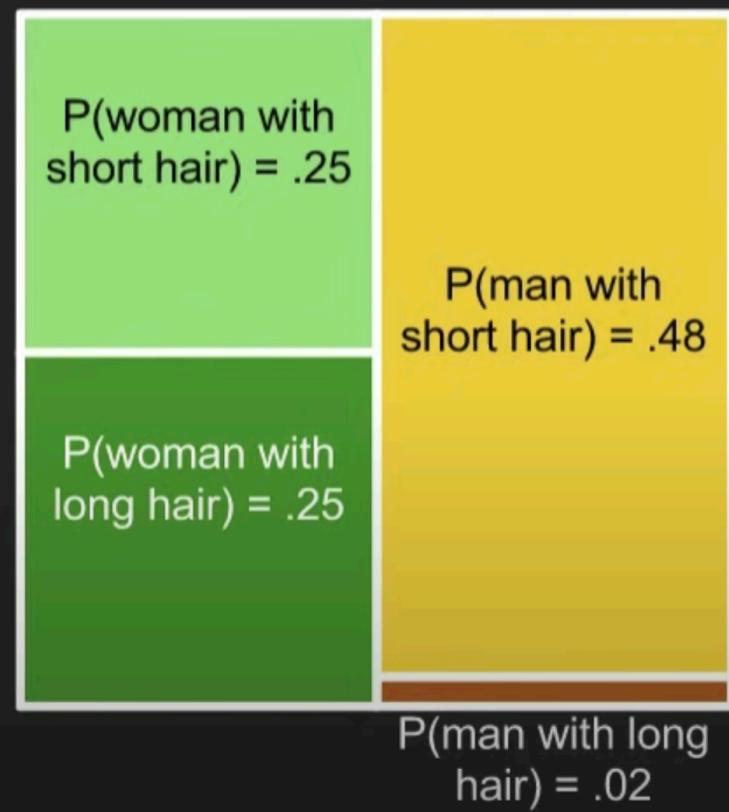
Men with
short hair



Women with
long hair



Men with
long hair



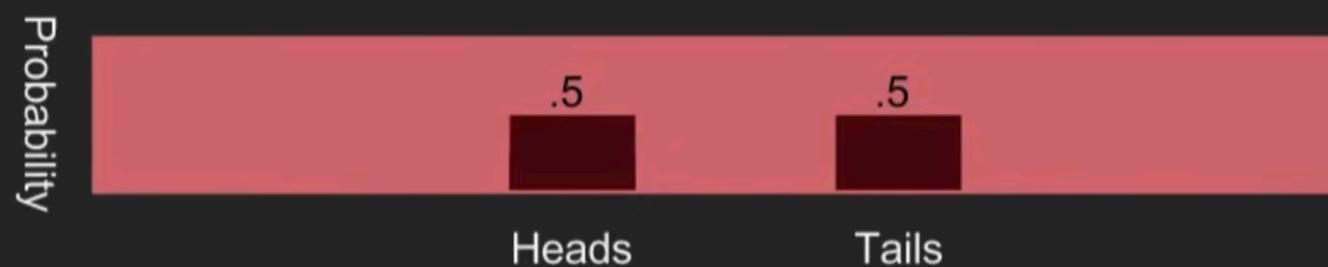
Probability distributions

It's helpful to think of probabilities as beliefs



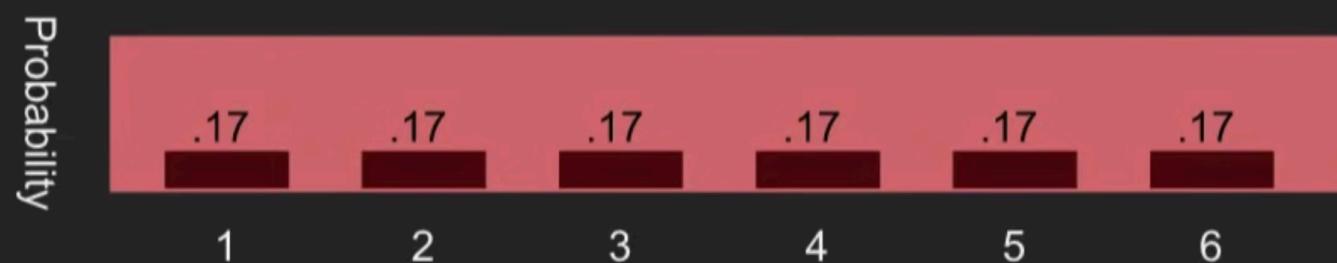
Probability distributions

Flipping a fair coin



Probability distributions

Rolling a fair die



Bayes Theorem

Bayes Rule:
$$P(h | D) = \frac{P(D | h)P(h)}{P(D)}$$

- $P(h)$ = prior probability of hypothesis h
- $P(D)$ = prior probability of training data D
- $P(h | D)$ = probability of h given D (posterior density)
- $P(D | h)$ = probability of D given h (likelihood of D given h)

Maximum A Posteriori (MAP) Hypothesis

$$P(h | D) = \frac{P(D | h)P(h)}{P(D)}$$

The Goal of Bayesian Learning: the most probable hypothesis given the training data (Maximum A Posteriori hypothesis)

$$\begin{aligned} h_{MAP} &= \arg \max_{h \in H} P(h | D) \\ &= \arg \max_{h \in H} \frac{P(D | h)P(h)}{P(D)} \\ &= \arg \max_{h \in H} P(D | h)P(h) \end{aligned}$$

Compute ML Hypo

$$h_{ML} = \arg \max_{h \in H} p(D | h)$$

$$= \arg \max_{h \in H} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{d_i - h(x_i)}{\sigma}\right)^2}$$

$$= \arg \max_{h \in H} \sum_{i=1}^m -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2} \left(\frac{d_i - h(x_i)}{\sigma}\right)^2$$

$$= \arg \min_{h \in H} \sum_{i=1}^m (d_i - h(x_i))^2$$

Bayes Optimal Classifier

Question: Given new instance x , what is its most probable classification?

- $h_{MAP}(x)$ is not the most probable classification!

Example: Let $P(h_1|D) = .4$, $P(h_2|D) = .3$, $P(h_3|D) = .3$

Given new data x , we have $h_1(x)=+$, $h_2(x) = -$, $h_3(x) = -$

What is the most probable classification of x ?

Bayes optimal classification:

$$\arg \max_{v_j \in V} \sum_{h_i \in H} P(v_j | h_i)P(h_i | D)$$

where V is the set of all the values a classification can take and v_j is one possible such classification.

Example:

$$\begin{array}{llll} P(h_1|D) = .4, & P(-|h_1)=0, & P(+|h_1)=1 & \sum_{h_i \in H} P(+|h_i)P(h_i|D) = .4 \\ P(h_2|D) = .3, & P(-|h_2)=1, & P(+|h_2)=0 & \sum_{h_i \in H} P(-|h_i)P(h_i|D) = .6 \\ P(h_3|D) = .3, & P(-|h_3)=1, & P(+|h_3)=0 & \end{array}$$

features: $X = (X_1, X_2, \dots, X_n)$

label: Y

$$P(Y=y|X=(x_1,x_2,...,x_n))$$

for what value of y

$$P(Y = y | X = (x_1, x_2, \dots, x_n))$$

is maximum

but... $P(Y|X)$ is hard to find!

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

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

Diagram illustrating the components of Bayes' Theorem:

- Likelihood**: $P(X|Y)$
- Prior**: $P(Y)$
- Evidence**: $P(X)$
- Posterior**: $P(Y|X)$

The diagram shows curved arrows pointing from the Prior and Evidence to the Likelihood term in the numerator, and from the Likelihood and Evidence terms to the Posterior term in the numerator.

X_1	X_2	Y	
0	0	0	
0	1	1	
1	2	1	
0	0	1	
2	2	0	$X_1, X_2 \in \{0, 1, 2\}$
1	1	0	$Y \in \{0, 1\}$
0	2	1	
2	0	0	
2	1	0	
1	0	0	

Estimate the value of Y given that $X = (0, 2)$

X_1	X_2	Y
0	0	0
0	1	1
1	2	1
0	0	1
2	2	0
1	1	0
0	2	1
2	0	0
2	1	0
1	0	0

Let's compute $P(Y = 0|X = (0, 2))$ and $P(Y = 1|X = (0, 2)) \dots$

$$P(Y = 0) = \frac{\#Y = 0}{\#Y = 0 + \#Y = 1} = \frac{6}{10}$$

$$P(Y = 1) = \frac{\#Y = 1}{\#Y = 0 + \#Y = 1} = \frac{4}{10}$$

$$P(X = (0, 2)|Y = 1) = \frac{1}{4}$$

$$P(X = (0, 2)|Y = 0) = 0$$

$$P(X = (0, 2)|Y = 0) * P(Y = 0) = 0 * \frac{6}{10} = 0$$

$$P(X = (0, 2)|Y = 1) * P(Y = 1) = \frac{1}{4} * \frac{4}{10} = \frac{1}{10}$$

X_1	X_2	Y
0	0	0
0	1	1
1	2	1
0	0	1
2	2	0
1	1	0
0	2	1
2	0	0
2	1	0
1	0	0

Let's compute $P(Y = 0|X = (0, 2))$ and $P(Y = 1|X = (0, 2)) \dots$

$$P(Y = 0) = \frac{\#Y = 0}{\#Y = 0 + \#Y = 1} = \frac{6}{10}$$

$$P(Y = 1) = \frac{\#Y = 1}{\#Y = 0 + \#Y = 1} = \frac{4}{10}$$

$$P(X = (0, 2)|Y = 1) = \frac{1}{4}$$

$$P(X = (0, 2)|Y = 0) = 0$$

Estimated value of Y is 1 given $X = (0, 2)$

X_1	X_2	Y
0	0	0
0	1	1
1	2	1
0	0	1
2	2	0
1	1	0
0	2	1
2	0	0
2	1	0
1	0	0

Let's compute $P(Y = 0|X = (0, 2))$ and $P(Y = 1|X = (0, 2)) \dots$

$$P(Y = 0) = \frac{\#Y = 0}{\#Y = 0 + \#Y = 1} = \frac{6}{10}$$

$$P(Y = 1) = \frac{\#Y = 1}{\#Y = 0 + \#Y = 1} = \frac{4}{10}$$

$$P(X = (0, 2)|Y = 1) = \frac{1}{4}$$

$$P(X = (0, 2)|Y = 0) = 0$$

X_1	X_2	Y
0	0	0
0	1	1
1	2	1
0	0	1
2	2	0
1	1	0
0	2	1
2	0	0
2	1	0
1	0	0

Let's assume X_1, X_2 are independant!

Let's compute $P(Y = 0|X = (0, 2))$ and $P(Y = 1|X = (0, 2)) \dots$

$$P(Y = 0) = \frac{\#Y = 0}{\#Y = 0 + \#Y = 1} = \frac{6}{10}$$

$$P(Y = 1) = \frac{\#Y = 1}{\#Y = 0 + \#Y = 1} = \frac{4}{10}$$

$$P(X = (0, 2)|Y = 1) = P(X_1 = 0|Y = 1) * P(X_2 = 2|Y = 1)$$

$$P(X = (0, 2)|Y = 0) = P(X_1 = 0|Y = 0) * P(X_2 = 2|Y = 0)$$

X_1	X_2	Y
0	0	0
0	1	1
1	2	1
0	0	1
2	2	0
1	1	0
0	2	1
2	0	0
2	1	0
1	0	0

Let's compute $P(Y = 0|X = (0, 2))$ and $P(Y = 1|X = (0, 2)) \dots$

$$P(Y = 0) = \frac{\#Y = 0}{\#Y = 0 + \#Y = 1} = \frac{6}{10}$$

$$P(Y = 1) = \frac{\#Y = 1}{\#Y = 0 + \#Y = 1} = \frac{4}{10}$$

$$P(X = (0, 2)|Y = 1) = P(X_1 = 0|Y = 1) * P(X_2 = 2|Y = 1) = \frac{3}{4} \cdot \frac{2}{4}$$

$$P(X = (0, 2)|Y = 0) = P(X_1 = 0|Y = 0) * P(X_2 = 2|Y = 0) = \frac{1}{6} \cdot \frac{1}{6}$$

$$4/10 * \frac{3}{4} * \frac{2}{4} > \frac{1}{6} * \frac{1}{6} * 6/10$$

**don't forget to multiply the priors
i.e. $P(Y=1)$ and $P(Y=0)$**

X_1	X_2	Y
0	0	0
0	1	1
1	2	1
0	0	1
2	2	0
1	1	0
0	2	1
2	0	0
2	1	0
1	0	0

Naïve Bayes

Bayes rule:

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k) P(X_1 \dots X_n | Y = y_k)}{\sum_j P(Y = y_j) P(X_1 \dots X_n | Y = y_j)}$$

Assuming conditional independence among X_i 's:

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

So, classification rule for $X^{new} = < X_1, \dots, X_n >$ is:

$$Y^{new} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

Naïve Bayes Algorithm – discrete X_i

- Train Naïve Bayes (examples)
for each* value y_k
estimate $\pi_k \equiv P(Y = y_k)$
for each* value x_{ij} of each attribute X_i
estimate $\theta_{ijk} \equiv P(X_i = x_{ij}|Y = y_k)$

- Classify (X^{new})

$$Y^{new} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{new}|Y = y_k)$$

$$Y^{new} \leftarrow \arg \max_{y_k} \pi_k \prod_i \theta_{ijk}$$

* probabilities must sum to 1, so need estimate only n-1 parameters...