

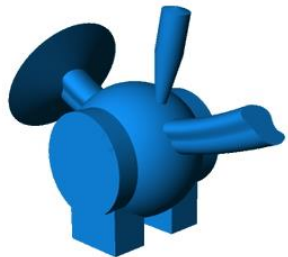
Naïve Bayes

Semester 1, 2023

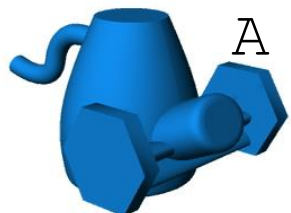
Kris Ehinger

Outline

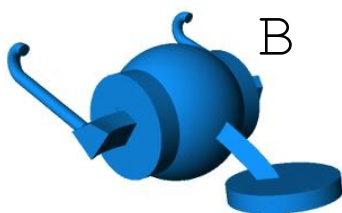
- Probabilistic learning
- Bayes' Rule
- Naïve Bayes classifier
- Practical issues, assumptions, etc.



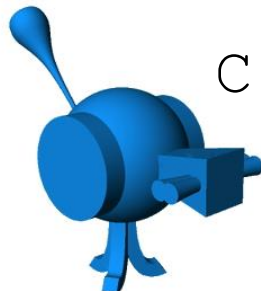
This is a “tufa”
Are there any more tufas
in the set below?



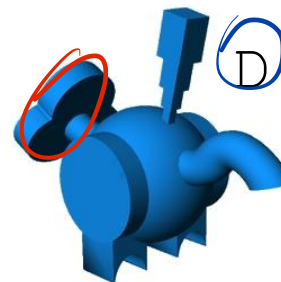
A



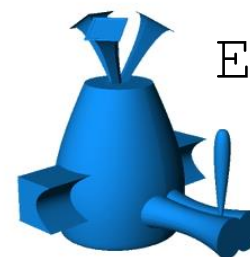
B



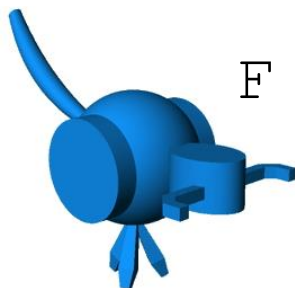
C



D



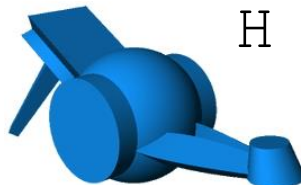
E



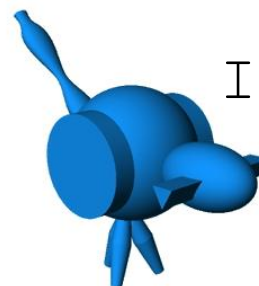
F



G



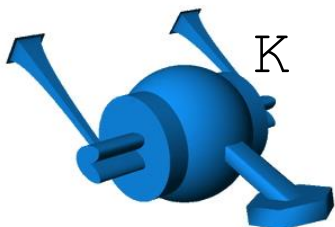
H



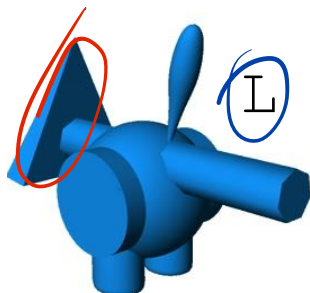
I



J



K



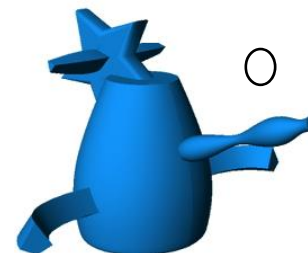
L



M



N



O

Concept learning

- People can grasp a “class” concept from very limited data:
 - One (or few) examples
 - Noisy features
 - Diverse data set
 - Ambiguity about what features might be important
- How to get this generalisation performance from machines?

Why probabilistic models?

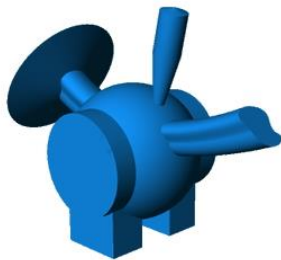
- Framework for modelling systems that are noisy/uncertain
- Rules that let you generalise from limited observations
- Based on laws of probability, so gives an optimal prediction given the available data
 - and the assumptions built into the model

Probabilistic learner

- Goal is classification, so we'll build a **supervised** model
- We need to build a probabilistic model of the training data, and then use that to predict the class of the test data
- In probability terms: given an instance T, which class c is most likely?
- $\hat{c} = \arg \max_{c \in C} P(c | T)$

Probabilistic learner

- How to do this?
- For each class c :
 - Find all examples of c in the training data
 - Count the number of times T has been observed
- Choose class \hat{c} with the greatest frequency of observed T



What if you've never seen this specific instance before?

Probabilistic learner

- Unfortunately, this requires a massive amount of data!
- We need to have seen every possible combination of attributes in the training set, ideally at least a few times per class, to have a good estimate of frequency
- For m attributes with k possible values and C classes, this mean $O(Ck^m)$ instances
 - 2 classes, 20 binary attributes: >2 million
 - 2 classes, 10 attributes with 10 values: 20 trillion

Naïve Bayes solution

- Assume different attributes are statistically independent, conditional on class
- Compute $P(c|T)$ from $P(T|c)$ using Bayes rule

$$P(c_j|T) = \frac{P(T|c_j)P(c_j)}{P(T)} \rightarrow \text{all same for } c_j$$

$$\max P(c_j|T) \rightarrow \max P(T|c_j)P(c_j)$$

Bayes' Rule

Bayes' Rule example

- Consider this situation:

You go to a shop that you're pretty sure is open today. When you get close, however, you notice all the lights are off inside and the windows are dark.

- What do you conclude?

Bayesian mental model

- When a shop is open, it usually doesn't turn the lights off:

$$P(\text{off} \mid \text{open}) = 0.01$$

- When a shop is closed, it usually does turn the lights off:

$$P(\text{off} \mid \text{closed}) = 0.85$$

$$P(\text{off}) = \overset{P(\text{off} \mid \text{open})}{P(\text{off} \mid \text{open})} \overset{P(\text{open})}{P(\text{open})} + \overset{P(\text{off} \mid \text{closed})}{P(\text{off} \mid \text{closed})} \overset{P(\text{closed})}{P(\text{closed})}$$

- So what are the chances the shop is open given that the lights are off?

$$P(\text{open} \mid \text{off})$$

Bayes' Rule

- For hypothesis H and evidence x
 - $P(H)$, the prior, is the initial degree of belief in H
 - $P(H|x)$, the posterior, is the degree of belief in H , given x
 - $P(x|H)$ is the likelihood of observing evidence x given that H is true

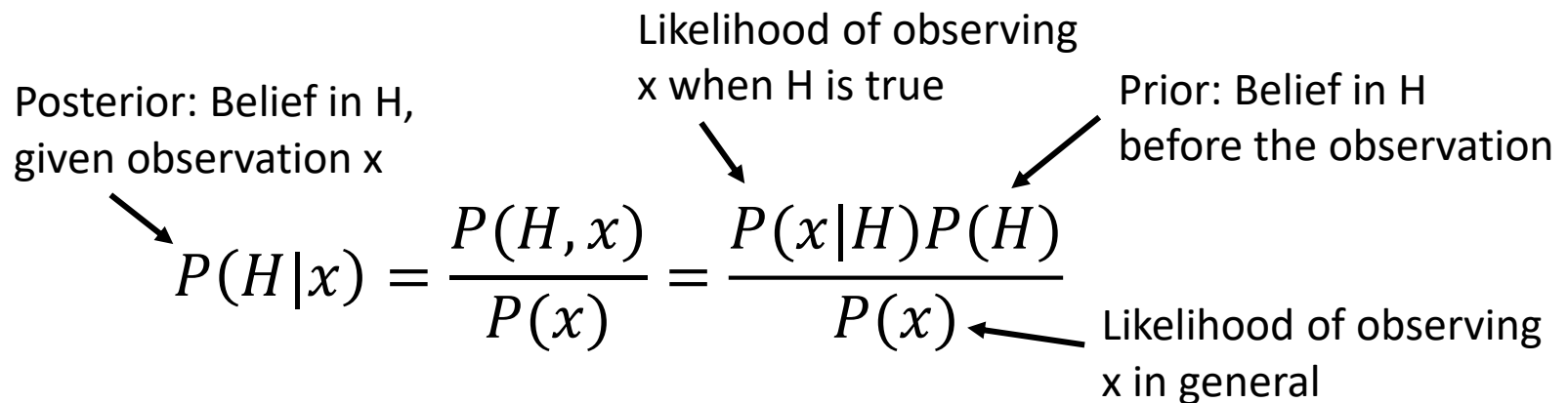


Diagram illustrating Bayes' Rule with annotations:

Posterior: Belief in H , given observation x → $P(H|x)$

Likelihood of observing x when H is true → $P(x|H)$

Prior: Belief in H before the observation → $P(H)$

Likelihood of observing x in general → $P(x)$

$$P(H|x) = \frac{P(H, x)}{P(x)} = \frac{P(x|H)P(H)}{P(x)}$$

Updating beliefs

- What are the chances the shop is open given that the lights are off?

$$P(H|x) = \frac{P(x|H)P(H)}{P(x)}$$

- Depends on prior P(open)...

$$P(\text{open} \mid \text{off}) = \frac{\overset{0.01}{P(\text{off} \mid \text{open})} \overset{0.95}{P(\text{open})}}{\underset{0.01}{P(\text{off} \mid \text{open})} \underset{0.95}{P(\text{open})} + \overset{0.85}{P(\text{off} \mid \text{closed})} \underset{0.05}{P(\text{closed})}}$$

$P(\text{open}) = 0.95$

$$P(\text{open} \mid \text{off}) = \frac{0.01 \times 0.95}{0.001 + 0.043} = 0.183$$

Bayes' Rule

- Bayes' Rule allows us to compute $P(H|x)$ when $P(x|H)$ can be estimated
 - In many situations, it's fairly easy to estimate $P(x|H)$ from real world observations or a theoretical model
- Powerful tool to represent a model's (or person's) beliefs, and how those beliefs are updated through experience with the world

Naïve Bayes classifier

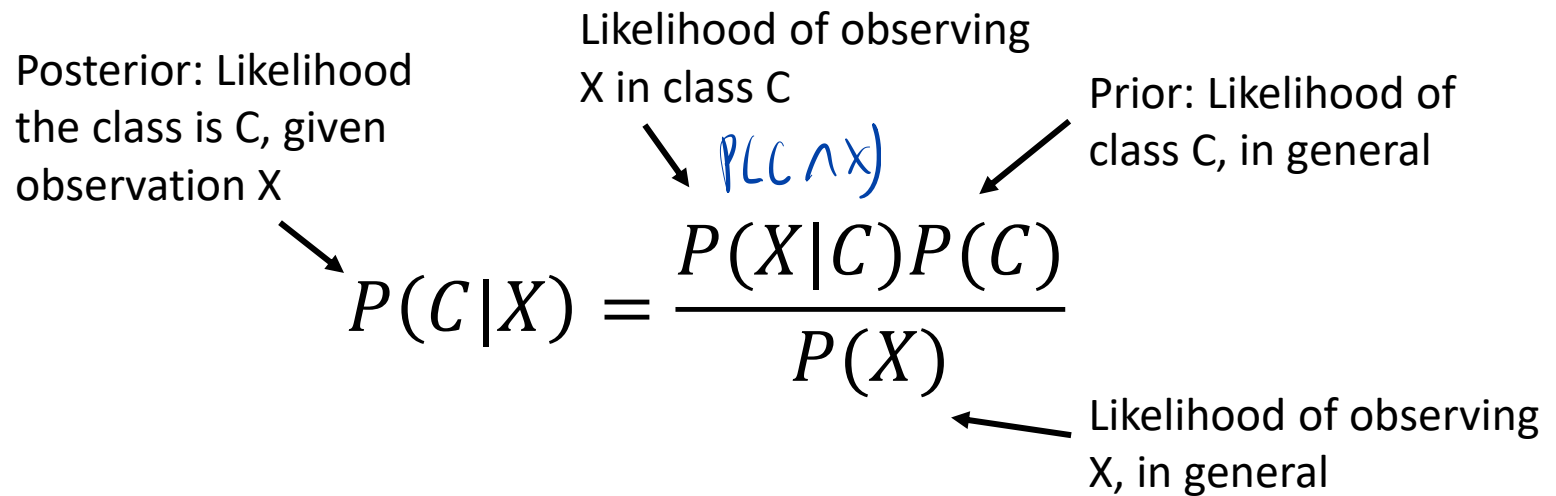
Bayes' Rule

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

Posterior: Likelihood
the class is C, given
observation X

Likelihood of observing
X in class C

Prior: Likelihood of
class C, in general


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

Likelihood of observing
X, in general

Naïve Bayes learner

- Task: classify an instance T into one of the possible classes $c_j \in C$

$$\begin{aligned}\hat{c} &= \arg \max_{c_j \in C} P(c_j | T) \\ &= \arg \max_{c_j \in C} \frac{P(T | c_j) P(c_j)}{\cancel{P(T)}}\end{aligned}$$

$P(c_j \cap T) = P(c_j | T) P(T) = P(T | c_j) P(c_j)$ ✓

Same for all $c_j \in C$, so we can ignore it

Naïve Bayes learner

- Task: classify an instance T into one of the possible classes $c_j \in C$

$$\hat{c} = \arg \max_{c_j \in C} P(c_j | T)$$

$$= \arg \max_{c_j \in C} P(T | c_j) P(c_j)$$

- Each class generates instances
- Each class could have generated this instance, with some likelihood
- Which class is most likely to have generated this instance?

Naïve Bayes learner

- Task: classify an instance $T = \langle x_1, x_2, \dots, x_n \rangle$ into one of the possible classes $c_j \in C$
attributes

$$\hat{c} = \arg \max_{c_j \in C} P(c_j | x_1, x_2, \dots, x_n)$$

$$= \arg \max_{c_j \in C} \frac{P(x_1, x_2, \dots, x_n | c_j) P(c_j)}{\cancel{P(x_1, x_2, \dots, x_n)}}$$

Same for all $c_j \in C$, so we can ignore it

$$= \arg \max_{c_j \in C} P(x_1, x_2, \dots, x_n | c_j) P(c_j)$$

?
hard to compute

What makes it “naïve”?

- Naïve Bayes assumes all attributes are independent, conditional on the class:

$$P(x_1, x_2, \dots, x_n | c_j) \approx P(x_1 | c_j) P(x_2 | c_j) \dots P(x_n | c_j)$$

↑
AND \cap

$$= \prod_i P(x_i | c_j)$$

- Makes the problem tractable – easy to compute these probabilities
- But it's almost always untrue in real-world data
- But naïve Bayes (usually) works pretty well anyway

Complete naïve Bayes learner

- Task: classify an instance $T = \langle x_1, x_2, \dots, x_n \rangle$ into one of the possible classes $c_j \in C$

$$\hat{c} = \arg \max_{c_j \in C} P(c_j | x_1, x_2, \dots, x_n) \quad \begin{matrix} \frac{P(T|c_j)P(c_j)}{P(T)} \end{matrix}$$

$$= \arg \max_{c_j \in C} P(x_1, x_2, \dots, x_n | c_j) P(c_j) \quad \begin{matrix} \downarrow \text{bayes rule + ignore denominator} \end{matrix}$$

Complete naïve Bayes learner

- Task: classify an instance $T = \langle x_1, x_2, \dots, x_n \rangle$ into one of the possible classes $c_j \in C$

$$\hat{c} = \arg \max_{c_j \in C} P(c_j | x_1, x_2, \dots, x_n)$$

$$= \arg \max_{c_j \in C} P(c_j) \prod_i P(x_i | c_j)$$

- Each class generates attributes, with some likelihood
- Which class is most likely to have generated these attributes?

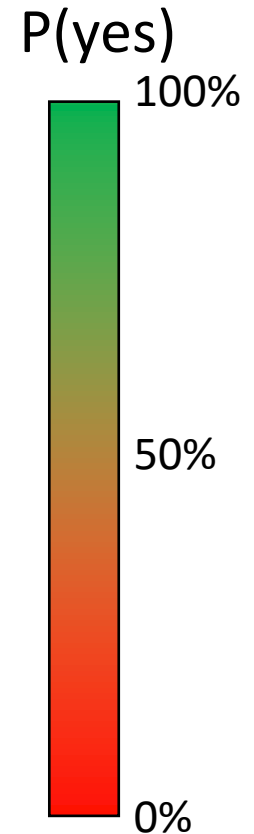
Bayesian prior

- The prior $P(c_j)$ can be estimated from the frequency of classes in the training set (maximum likelihood estimate)
- Naïve Bayes learns the priors from the training set and uses them in prediction
 - Good if the training set correctly reflects the real-world / test set distribution of classes (or is close)
 - But potentially a problem if you want to apply the classifier to a new situation with new priors

↓
new class

Computing probability

Outlook	Temp	Humidity	Windy	Play?
sunny	cool	normal	false	yes
sunny	cool	normal	false	yes
sunny	cool	normal	false	yes
sunny	cool	normal	false	yes
sunny	cool	normal	false	yes
sunny	cool	normal	false	yes
sunny	cool	normal	false	yes
overcast	hot	high	true	no



Today it's sunny, cool, normal, false. What's the class?

$$P(c_j | x_1, x_2, x_3, x_4) \Rightarrow \text{yes}, P=1$$

Computing probability

Outlook	Temp	Humidity	Windy	Play?
rainy	hot	normal	true	yes
rainy	hot	normal	true	no
rainy	hot	normal	true	yes
rainy	hot	normal	true	no
rainy	hot	normal	true	yes
rainy	hot	normal	true	no
sunny	cool	normal	false	yes
sunny	mild	high	false	no
overcast	cool	high	true	no

P(yes)

100%

50%

0%

Today it's rainy, hot, normal, true. What's the class?

$P(c_j | \text{rainy, hot, normal, true})$, yes : $P=0.5$
no : $P=0.5$

Computing probability

Outlook	Temp	Humidity	Windy	Play?
<u>overcast</u>	<u>mild</u>	normal	true	yes
sunny	<u>mild</u>	normal	<u>false</u>	yes
<u>overcast</u>	hot	<u>high</u>	true	yes
sunny	cool	<u>high</u>	<u>false</u>	yes
rainy	cool	normal	true	no
<u>overcast</u>	hot	normal	true	no
sunny	hot	normal	<u>false</u>	no
sunny	<u>mild</u>	normal	true	no
overcast	cool	<u>high</u>	true	no

P(yes)

100%

50%

0%

Today it's overcast, mild, high, false. What's the class?

$\max P(x_1 | c_j) P(x_2 | c_j) P(c_j)$

each of features are more common in yes

This is naive bayes. consider each feature independently

Naïve Bayes example

- Given a training data set, what probabilities do we need to estimate?

Headache	Sore	Temperature	Cough	Diagnosis
severe	mild	high	yes	Flu
no	severe	normal	yes	Cold
mild	mild	normal	yes	Flu
mild	no	normal	no	Cold
severe	severe	normal	yes	Flu

We need $P(c_j)$, $P(x_i | c_j)$ for every x_i , c_j

Naïve Bayes example

Headache	Sore	Temperature	Cough	Diagnosis
severe	mild	high	yes	Flu
no	severe	normal	yes	Cold
mild	mild	normal	yes	Flu
mild	no	normal	no	Cold
severe	severe	normal	yes	Flu

$$P(\text{Flu}) = 3/5$$

$$P(\text{Headache} = \text{severe} \mid \text{Flu}) = 2/3$$

$$P(\text{Headache} = \text{mild} \mid \text{Flu}) = 1/3$$

$$P(\text{Headache} = \text{no} \mid \text{Flu}) = 0/3$$

$$P(\text{Cold}) = 2/5$$

$$P(\text{Headache} = \text{severe} \mid \text{Cold}) = 0/2$$

$$P(\text{Headache} = \text{mild} \mid \text{Cold}) = 1/2$$

$$P(\text{Headache} = \text{no} \mid \text{Cold}) = 1/2$$

Naïve Bayes example

$$P(\text{Flu}) = 3/5$$

$$P(\text{Headache} = \text{severe} | \text{Flu}) = 2/3$$

$$P(\text{Headache} = \text{mild} | \text{Flu}) = 1/3$$

$$P(\text{Headache} = \text{no} | \text{Flu}) = 0/3$$

$$P(\text{Sore} = \text{severe} | \text{Flu}) = 1/3$$

$$P(\text{Sore} = \text{mild} | \text{Flu}) = 2/3$$

$$P(\text{Sore} = \text{no} | \text{Flu}) = 0/3$$

$$P(\text{Temp} = \text{high} | \text{Flu}) = 1/3$$

$$P(\text{Temp} = \text{normal} | \text{Flu}) = 2/3$$

$$P(\text{Cough} = \text{yes} | \text{Flu}) = 3/3$$

$$P(\text{Cough} = \text{no} | \text{Flu}) = 0/3$$

$$P(\text{Cold}) = 2/5$$

$$P(\text{Headache} = \text{severe} | \text{Cold}) = 0/2$$

$$P(\text{Headache} = \text{mild} | \text{Cold}) = 1/2$$

$$P(\text{Headache} = \text{no} | \text{Cold}) = 1/2$$

$$P(\text{Sore} = \text{severe} | \text{Cold}) = 1/2$$

$$P(\text{Sore} = \text{mild} | \text{Cold}) = 0/2$$

$$P(\text{Sore} = \text{no} | \text{Cold}) = 1/2$$

$$P(\text{Temp} = \text{high} | \text{Cold}) = 0/2$$

$$P(\text{Temp} = \text{normal} | \text{Cold}) = 2/2$$

$$P(\text{Cough} = \text{yes} | \text{Cold}) = 1/2$$

$$P(\text{Cough} = \text{no} | \text{Cold}) = 1/2$$

Naïve Bayes example

- A patient comes to the clinic with mild headache, severe soreness, normal temperature, and no cough. Are they more likely to have cold or flu?

$$P(x_1|c_j) P(x_2|c_j) P(x_3|c_j) P(c_j) \approx P(c_j) P(x_1, x_2, x_3|c_j) = \max_{c_j} \frac{P(c_j \wedge T)}{P(T)} \\ = \max_{c_j} P(c_j | T)$$

Cold: $P(C)P(H = m|C)P(S = s|C)P(T = n|C)P(C = n|C)$

$$\frac{2}{5} \left(\frac{1}{2}\right) \left(\frac{1}{2}\right) \left(\frac{2}{2}\right) \left(\frac{1}{2}\right) = 0.05$$

not P(patient having cold)
just relative probability

Flu: $P(F)P(H = m|F)P(S = s|F)P(T = n|F)P(C = n|F)$

$$\frac{3}{5} \left(\frac{1}{3}\right) \left(\frac{1}{3}\right) \left(\frac{2}{3}\right) \left(\frac{0}{3}\right) = 0$$

Naïve Bayes example

- A patient comes to the clinic with severe headache, mild soreness, high temperature, and no cough. Are they more likely to have cold or flu?

Cold: $P(C)P(H = s|C)P(S = m|C)P(T = h|C)P(C = n|C)$

$$\frac{2}{5} \left(\frac{0}{2}\right) \left(\frac{0}{2}\right) \left(\frac{0}{2}\right) \left(\frac{1}{2}\right) = \underline{0}$$

training dataset is too small
no all combination

Flu: $P(F)P(H = s|F)P(S = m|F)P(T = h|F)P(C = n|F)$

$$\frac{3}{5} \left(\frac{2}{3}\right) \left(\frac{2}{3}\right) \left(\frac{1}{3}\right) \left(\frac{0}{3}\right) = \underline{0}$$

equal probability

Zero values

- Problem: if any $P(x_i | c_j) = 0$, the final value will be zero
 - This means we need to see every possible pairing of attributes to class, which means a lot of data; likely to miss some pairings in a real-world data set
 - The 0s are actually informative – the fact that we never observed a pairing means it's probably rare
- Solution: treat unobserved events as possible but unlikely (all probabilities > 0)

Probabilistic smoothing I

- Simple option: whenever you encounter $P(x_i | c_j) = 0$, just replace 0 with a small positive constant ϵ
- ϵ should be very small, much less than $(1/N)$ ($N =$ number of instances), *given the class Flu: $N=3$, Cold: $N=2$*
- Since ϵ is tiny, assume all probabilities still sum to 1 (no need to scale other values)
- In practice, tends to reduce to comparisons to the cardinality of ϵ (meaning, whichever class has fewest ϵ s wins) *underestimate the rare event*

Probabilistic smoothing example

- A patient comes to the clinic with severe headache, mild soreness, high temperature, and no cough. Are they more likely to have cold or flu?

Cold: $P(C)P(H = s|C)P(S = m|C)P(T = h|C)P(C = n|C)$

$$\frac{2}{5}(\varepsilon)(\varepsilon)(\varepsilon)\left(\frac{1}{2}\right) = \frac{\varepsilon^3}{5}$$

Flu: $P(F)P(H = s|F)P(S = m|F)P(T = h|F)P(C = n|F)$

$$\frac{3}{5}\left(\frac{2}{3}\right)\left(\frac{2}{3}\right)\left(\frac{1}{3}\right)(\varepsilon) = \frac{4\varepsilon}{45}$$

Probabilistic smoothing II

- Slightly more complicated option: increase all counts by 1 (Laplace smoothing)
 - Unseen events get a count of 1
 - Events seen once become 2, twice 3, etc. ^{$|t| =$}
 - A more general version: increase all counts by α , which is a value between 0 and 1
 - Formula for an attribute X with d values:

Unsmoothed:

$$P_i = \frac{x_i}{N}$$

Smoothed:

$$P_i = \frac{x_i + \alpha}{N + \alpha d}$$

\Rightarrow level headache: $d=3$

Probabilistic smoothing example

Headache	Sore	Temperature	Cough	Diagnosis
severe	mild	high	yes	Flu
no	severe	normal	yes	Cold
mild	mild	normal	yes	Flu
mild	no	normal	no	Cold
severe	severe	normal	yes	Flu

$$P(\text{Headache} = \text{severe} | \text{Flu}) = \frac{1+2}{3+3} = 3/6$$

$$P(\text{Headache} = \text{mild} | \text{Flu}) = \frac{1+1}{3+3} = 2/6$$

$$P(\text{Headache} = \text{no} | \text{Flu}) = \frac{1+0}{3+3} = 1/6$$

$$P(\text{Headache} = \text{severe} | \text{Cold}) = \frac{1+0}{3+2} = 1/5$$

$$P(\text{Headache} = \text{mild} | \text{Cold}) = \frac{1+1}{3+2} = 2/5$$

$$P(\text{Headache} = \text{no} | \text{Cold}) = \frac{1+1}{3+2} = 2/5$$

Probabilistic smoothing II

- Most common method of Laplace smoothing is add-one smoothing ($\alpha=1$, all counts increase by 1)
- Probabilities are changed drastically when there are few instances, but the changes are smaller with more instances (mimics confidence)
- Add-one smoothing is known to overestimate the likelihood of rare events
 - But ϵ or lower values of α can underestimate
 - Hard to choose the right α in practice

Laplace
smoothing:
 N is large
denominator
is large.
probability
don't change
much

N is small

Probabilistic smoothing III

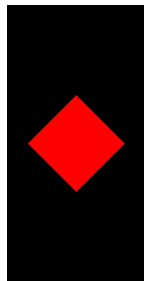
- Even more options:
 - Add-k smoothing: like Laplace smoothing, but adds a value $k > 1$ *overestimate likelihood of rare events*
 - Good-Turing estimation: uses the observed counts of different events to estimate how likely you are to see a never-before-seen event *use observed data to predict unobserved data*
 - Regression

Missing values

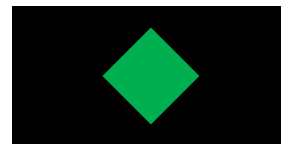
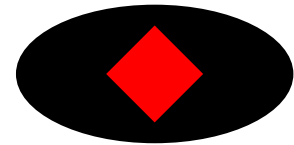
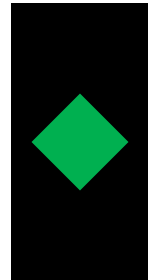
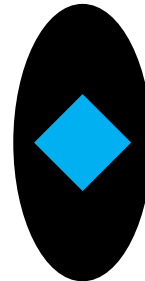
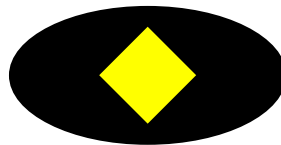
- What if an instance is missing some attribute?
- Missing values at test can simply be ignored – compute the likelihood of each class from the non-missing values
- Missing values in training can also be ignored – don't include them in the attribute-class counts, and the probabilities will be based on the non-missing values

Naïve Bayes example

- Simple shapes dataset:
 - Tall or Wide
 - Oval or Rectangle
 - Red, Yellow, Green, or Blue

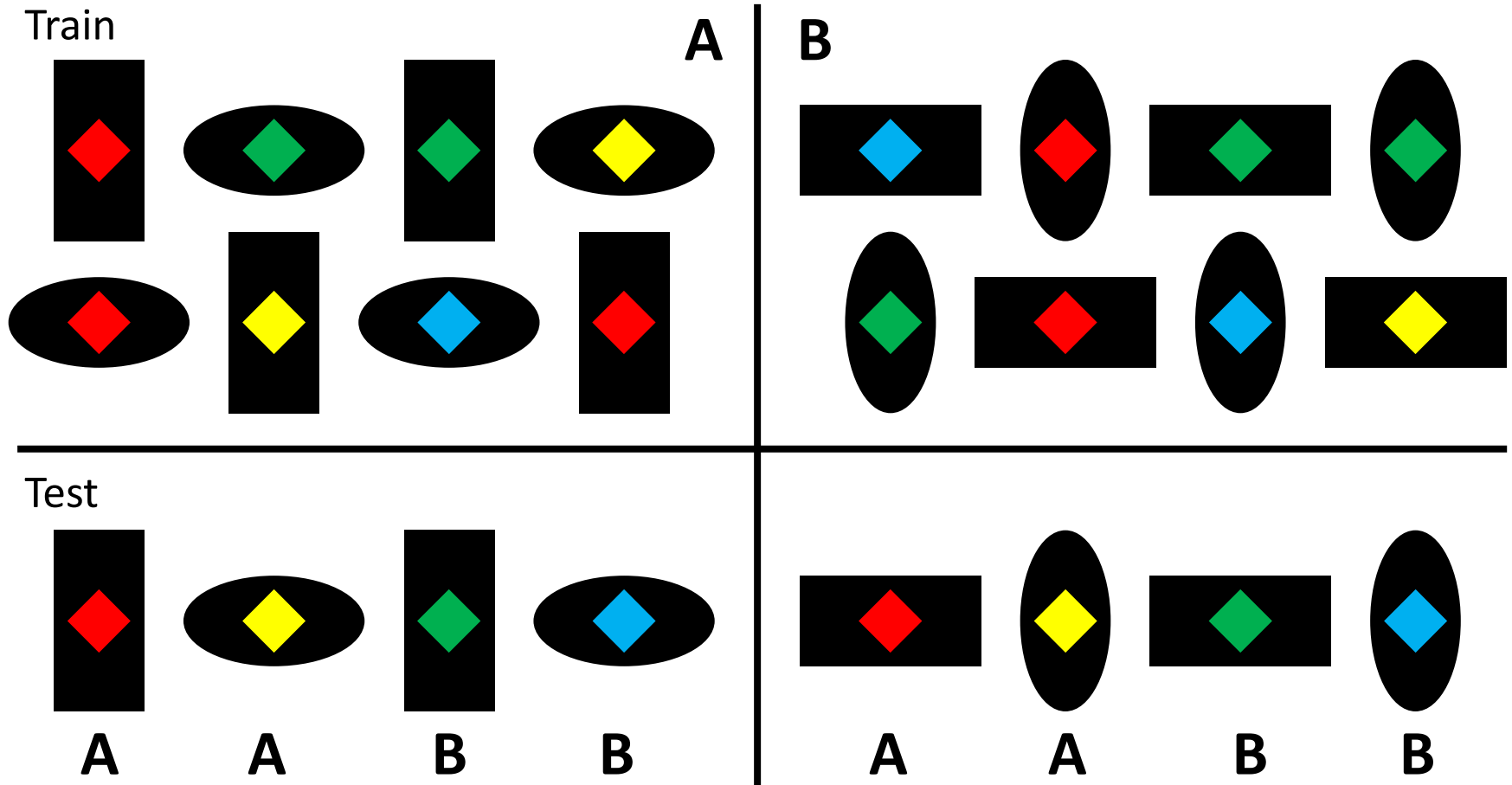


Tall, Rectangle, Red



Wide, Rectangle, Green

Naïve Bayes example



Naïve Bayes summary

- Why does it (usually) work, despite the wrong assumption of conditional independence?
 - We don't need a perfect estimate of $P(c|T)$ for every class – we just need to know which class is most likely
 - Ignoring the fact that some attributes are correlated tends to make all the class probabilities higher, but doesn't typically change their rank
 - Naïve Bayes is also robust to small errors in estimating $P(x_i|c_j)$ – these can change class probabilities but typically don't change class rank

Naïve Bayes

- Strengths:
 - Simple to build, fast
 - Computations scale well to high-dimensional datasets (1000s of attributes)
 - Explainable – generally easy to understand why the model makes the decision it does
- Weaknesses:
 - Inaccurate when there are many missing $P(x_i | c_j)$ values
 - Conditional independence assumption becomes problematic for complex systems