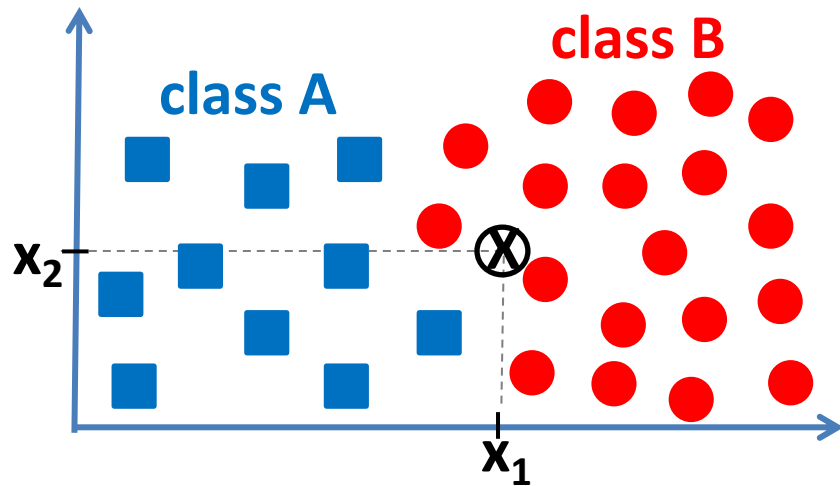


Machine Learning

Lecture 07

Naïve-Bayes model

- Which class does \mathbf{x} belong to?



Class (prior) probabilities tell us how dominant each individual class is over the whole dataset

$$p(\text{classA}) = 10/30 = 33.3\%$$

$$p(\text{classB}) = 20/30 = 66.7\%$$

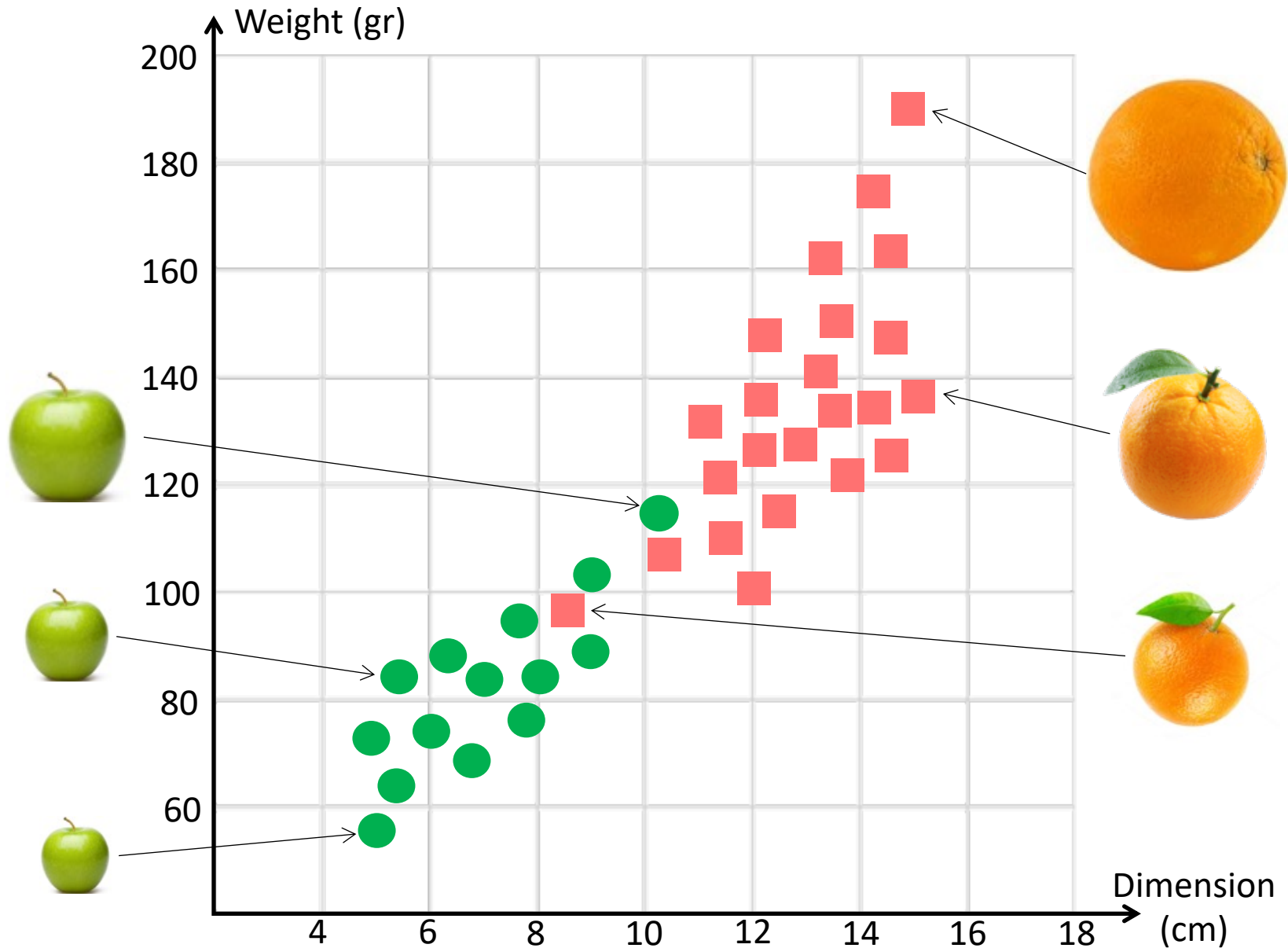
- By examining prior probabilities, we can conclude:

```
if P(classB) > P(classA)
  then classB
else classA
```

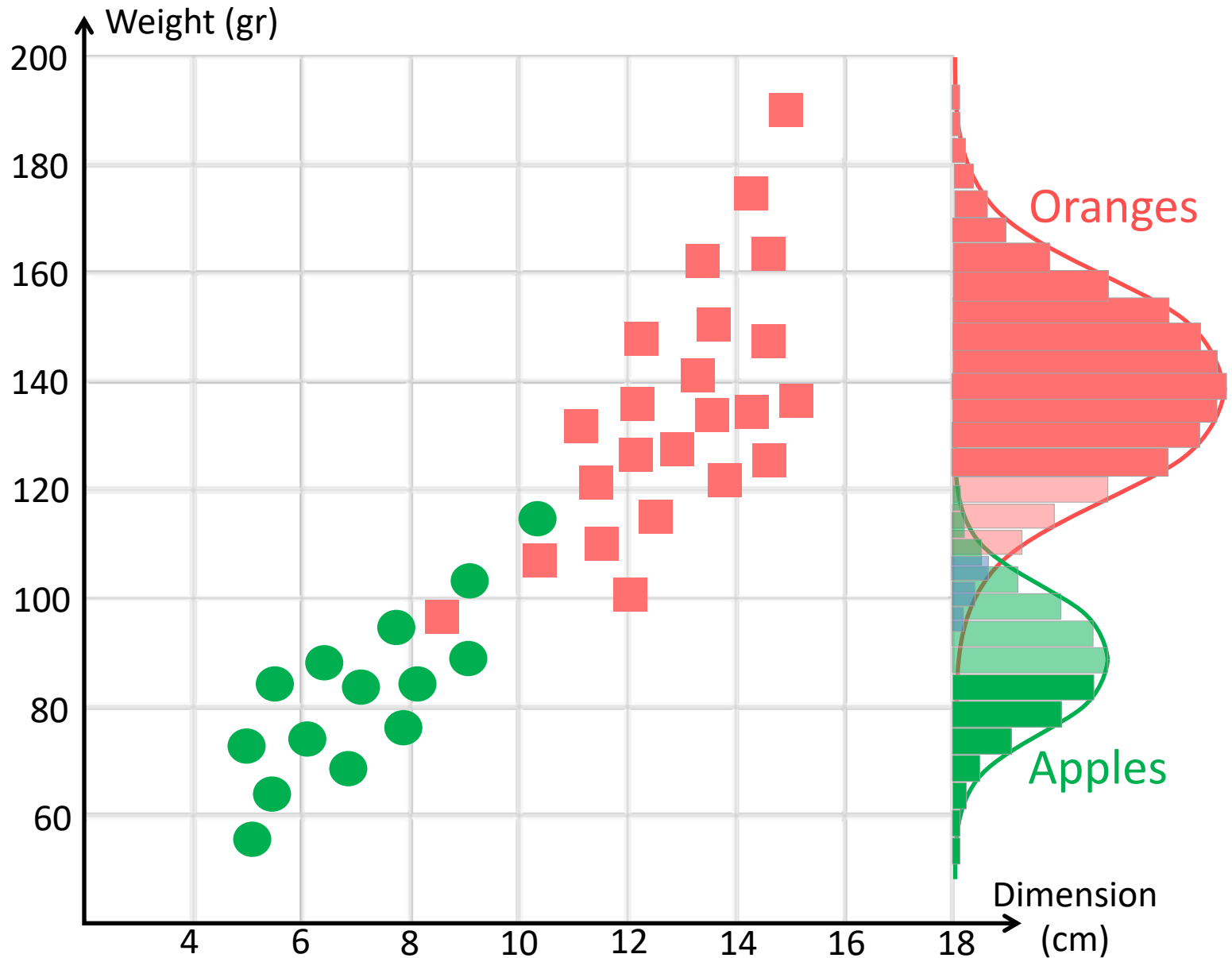
These rules don't tell the whole story though. What we really need is a comparison between conditional probabilities. Given the point $X=(x_1, x_2)$, what is the probability that X belongs to A (or B)?

$$p(\text{classB}|\mathbf{X}) > p(\text{classA}|\mathbf{X})?$$

Intuition on NB



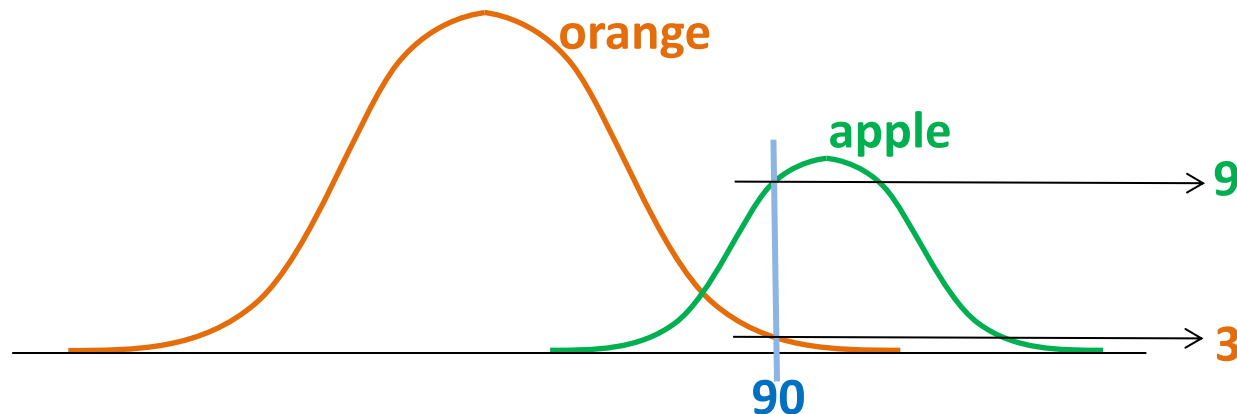
Intuition on NB



Intuition on NB

- We want to classify a fruit based on its weight (90 grams)?
- Given the probability distributions of weight for oranges and apples, we can find out which class is more likely?

$p(\text{class} | w=90)$ = probability of a class given an observed w



$$p(\text{orange} | \text{weight}=90) = \frac{3}{(9+3)} = 0.25$$

$$p(\text{apple} | \text{weight}=90) = \frac{9}{(9+3)} = 0.75$$

- For classifications based on dimension and weight, you need to construct a joint probability distribution based on weights and dimensions.

- Naive-Bayes algorithm for predictive modeling
 - A statistical classifier that can predict class membership probabilities (a generative method)
 - Instead of mapping the instance to a class, assign the class with the maximum probability given the instance
 - An eager learner (ideal for real-time predictions)
 - Despite its simplicity, it's surprisingly accurate
 - Somewhat immune to overfitting
 - Suffers from high bias and exhibits low variance
 - Underlying probabilistic model: The Bayes Theorem
 - Basic assumption: Given a class, every attribute is independent of all other attributes (the value of one feature has no effect on that of another – thus the word "Naive"... more on this later). This is hardly ever true, but it works pretty well in most cases.

Characteristics of NB

- Learning method: Supervised learning (generative, parametric and a linear model)
- Handling of mixed data types: Yes (both categorical and continuous data)
- Requires linearly separable data? No
- Classification ability: Binary and Multi-class
- Computationally demanding? No
- Ability for online classification: Yes
- Interpretability: High
- When to consider: Text classification, spam detection, target marketing, medical diagnosis, credit approval, recommender systems, emotion detection, face recognition, sentiment analysis

Basic rules of probability

Event: A possible outcome of an observation

$P(E)$ is probability of an event $E \Rightarrow E$ is impossible : $P(E)=0$

E is certain : $P(E)=1$

For any event E : $0 \leq P(E) \leq 1$

All events that are not E : $E^c \Rightarrow P(E) + P(E^c) = 1$

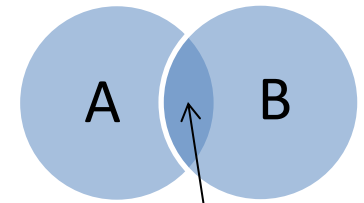
Union of two events A and B : $A \cup B$

Probability of all events in A or B or both: $P(A \cup B)$

Probability of all events common to both A & B : $P(A \cap B)$

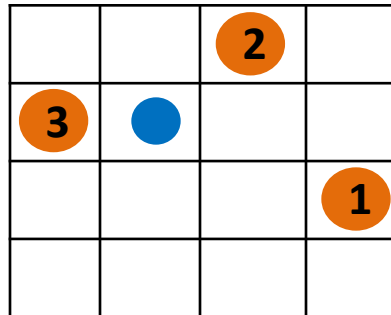
Probability of A or B : $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

Zero if A and B have no elements in common (A and B are mutually exclusive)



Intuition behind the Bayes' Theorem

- we can never be fully certain of the world, as it is constantly changing. Fundamental principle behind this theorem, is: update and improve our knowledge of reality (prior) as we get more and more data or evidence.
- Suppose we have an invisible blue ball in the field below



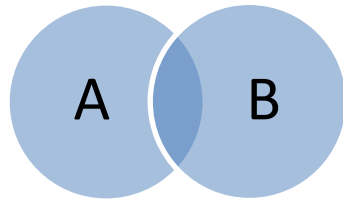
Someone is placing a visible yellow ball and tells you where the blue ball is located with respect to the yellow one, and you update your knowledge for the probability of the location of the blue ball:

Prior gets updated
as new evidence is
gathered

0. \Rightarrow Prior : $P(x) = 1/16$
1. "Left" \Rightarrow Posterior: $P(x|1) = 1/8$
2. "Behind" \Rightarrow Posterior: $P(x|1,2) = 1/4$
3. "Front" \Rightarrow Posterior: $P(x|1,2,3) = 1/2$

Conditional Probability

- Probability of A given that B has occurred:



$$P(A | B) = \frac{P(A \cap B)}{P(B)} \quad \text{1} \quad (\text{given that } P(B) \neq 0)$$

Similarly: $P(B | A) = \frac{P(B \cap A)}{P(A)} \quad \text{2} \quad (\text{given that } P(A) \neq 0)$

Conditional probability (from 1 & 2):

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

- Bayes' theorem:

$$P(\text{hypothesis} | \text{data}) = \frac{P(\text{data} | \text{hypothesis})P(\text{hypothesis})}{P(\text{data})}$$

Likelihood of the evidence if the hypothesis is true

Prior probability

Posterior probability of hypothesis given the evidence

Prior probability that the evidence itself is true

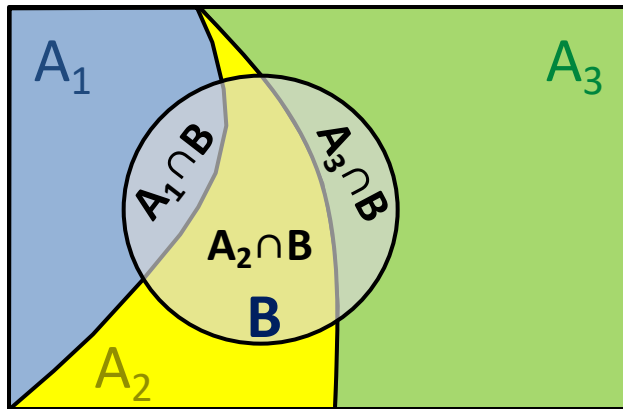
Conditional Probability – cont'd

- Expanded form:

$$\begin{aligned} P(B) &= P[(B \cap A) \cup (B \cap A^c)] = P(B \cap A) + P(B \cap A^c) \\ &= P(B | A)P(A) + P(B | A^c)P(A^c) \end{aligned}$$

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)} = \frac{P(B | A)P(A)}{\underbrace{P(B | A)P(A) + P(B | A^c)P(A^c)}_{\sum_j P(B | A_j)P(A_j)}}$$

- Total probability theorem:**

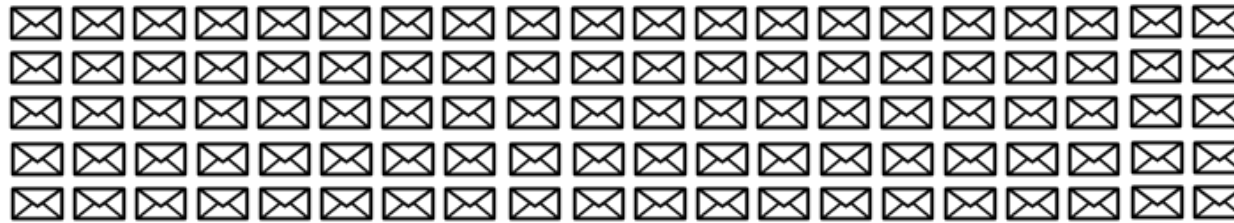


$$P(B) = P(A_1 \cap B) + P(A_2 \cap B) + P(A_3 \cap B)$$

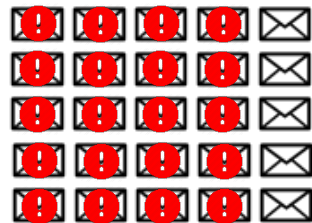
$$\begin{aligned} P(B) &= P(B | A_1)P(A_1) + P(B | A_2)P(A_2) + \\ &\quad P(B | A_3)P(A_3) \end{aligned}$$

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

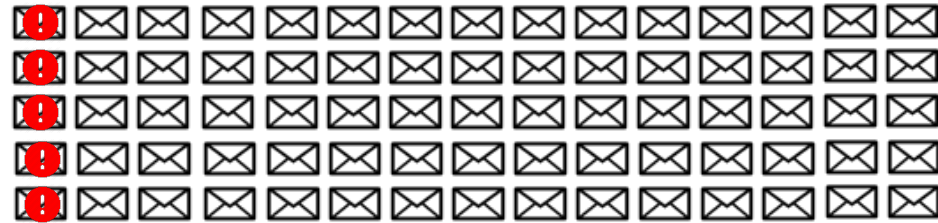
100 emails



25 spam (25%)

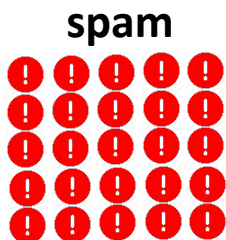


75 ham (75%)



! : Messages containing the word "Buy"

- Suppose 20 of the messages labeled "spam" contains the word "Buy"
- And 5 of the messages labeled "ham" contains the word "Buy"



If a mail contains the Word "Buy", what is the probability that it is **spam**?

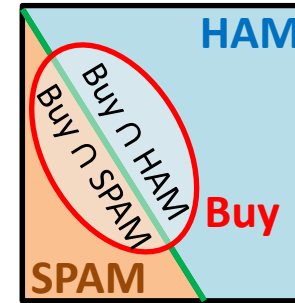
20 spam, 5 ham => $p = 20/(20+5) = 80\%$

Example I – cont'd

- Let's solve this by using the Bayes' theorem:

$$P(\text{spam} \mid \text{Buy}) = ?$$

$$P(\text{spam} \mid \text{Buy}) = \frac{P(\text{Buy} \mid \text{spam})P(\text{spam})}{P(\text{Buy})}$$



$$P(\text{Buy}) = P(\text{Buy} \cap \text{spam}) + P(\text{Buy} \cap \text{ham})$$
$$P(\text{Buy}) = P(\text{Buy} \mid \text{spam})P(\text{spam}) + P(\text{Buy} \mid \text{ham})P(\text{ham})$$

$$P(\text{Buy} \mid \text{spam}) = \frac{20}{25} \qquad P(\text{spam}) = 0.75$$

$$P(\text{Buy} \mid \text{ham}) = \frac{5}{75} \qquad P(\text{ham}) = 0.25$$

$$P(\text{spam} \mid \text{Buy}) = \frac{\frac{20}{25} * 0.25}{\frac{20}{25} * 0.25 + \frac{5}{75} * 0.75} = \frac{1/5}{1/5 + 5/100} = \frac{0.2}{0.25} = 0.8$$

Example 2

- A new method (VBI: Vision-based Inspection) is deployed to identify defective products. Here is how VBI works:
 - 1% of products are defective (statistical fact)
 - 80% of the VBI identify the defect when present
 - 9.6% of the VBI identify the defect when it's not there
- What's the chance the product is defective when it tests positive?

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

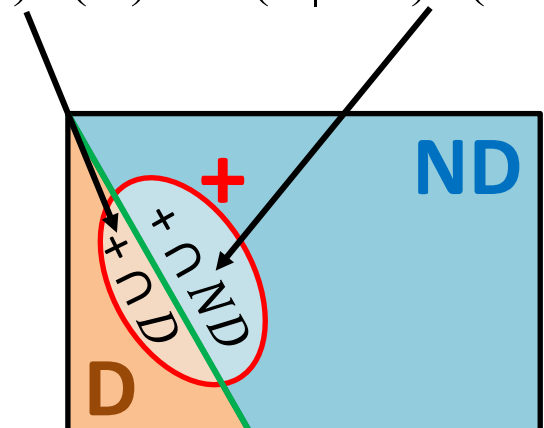
Total probability theorem

$$P(+) = P(+ \cap D) + P(+ \cap ND)$$

$$= P(+|D)P(D) + P(+|ND)P(ND)$$

$$P(+|D) = \frac{P(+ \cap D)}{P(D)}$$

$$P(D|+) = \frac{0.80 \cdot 0.01}{0.80 \cdot 0.01 + 0.096 \cdot 0.99} = 7.8\%$$



NB is a generative classifier

- Discriminative classifiers learn the probability of a class given the feature set (remember Logistic regression)

$p(y/x)$ where y is the class label and x is the feature vector

– Probability of a certain class given the feature set?

- Generative classifiers learn what kind of features are plausible given the class label

$p(x/y)$ probability of the feature set given the class label

- Together with the $p(x/y)$, we also have the $p(y)$, which is the **class prior** (independent of the features, this is the probability of having the class label in the population)
- Example: $p(y=1)$: probability of carrying the Ebola virus which is supposed to be small. So, the probability of having Ebola virus given the features (set of symptoms):

$$p(y=1|x) = \frac{p(x|y=1)p(y=1)}{p(x)}$$

$$\begin{aligned} p(x) &= \sum p(x|y_i)p(y_i) \\ &= p(x|y=1)p(y=1) + p(x|y=0)p(y=0) \end{aligned}$$

Why is NB so popular?

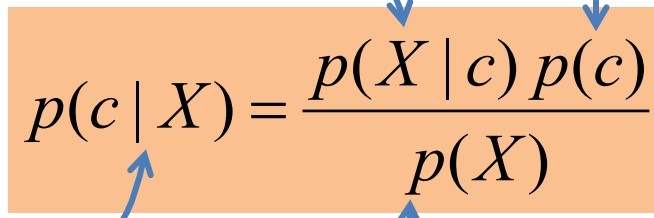
- **Classification problem:** What is the probability of class c given the set of features by X ? It follows from the Bayes' theorem:

Likelihood (probability of predictor given the class)

Prior probability of class

Posterior probability of class given the predictor

Prior probability of predictor (normalizer)

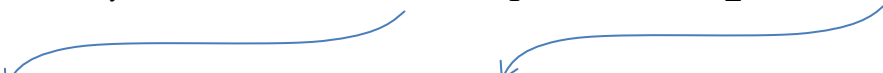
$$p(c | X) = \frac{p(X | c) p(c)}{p(X)}$$
A diagram illustrating the components of Bayes' theorem. The equation $p(c | X) = \frac{p(X | c) p(c)}{p(X)}$ is centered in an orange box. Four blue arrows point from descriptive text to parts of the equation: one from 'Likelihood (probability of predictor given the class)' to $p(X | c)$, one from 'Prior probability of class' to $p(c)$, one from 'Posterior probability of class given the predictor' to $p(c | X)$, and one from 'Prior probability of predictor (normalizer)' to $p(X)$.

- The numerator is computed as follows:

$$\begin{aligned} p(X | c)p(c) &= p(c)p(x_1, x_2, x_3, \dots, x_n | c) \\ &= p(c)p(x_1 | c)p(x_2, x_3, \dots, x_n | c, x_1) \\ &= p(c)p(x_1 | c)p(x_2 | c, x_1)p(x_3, \dots, x_n | c, x_1, x_2) \\ &= \dots \end{aligned}$$

Why is NB so popular?

- You have to keep track of all probabilities for each possible value of x_i . In the simplest case, x_i is just binary (0 or 1) yielding 2^n probabilities to estimate.
- It is clear that this is way too complex and a simplification is needed. By assuming all features are **conditionally independent**, last equation in the previous slide can be further simplified into the following form:

$$p(X | c)p(c) = p(c) \prod_{i=1}^n p(x_i | c) = p(c) p(x_1 | c) p(x_2 | c) \dots p(x_n | c)$$


Probability of class c generating the observed value for feature **1**

Probability of class c generating the observed value for feature **n**

- It all comes down to counting and multiplication which can all be pre-computed. Classification then becomes easy and fast. Assign the class label of whichever has the highest likelihood and it's done!

Types of NB models

- Based on the assumption for distribution of $p(X/c)$
- **Gaussian NB**
 - Used for continuous data (real numbers), features are assumed to follow a Normal (Gaussian) distribution
- **Bernoulli NB**
 - Used when the feature vectors are binary (male/female) or boolean (True/False) – the binomial model
- **CategoricalNB**
 - Used for categorical features with discrete values (e.g., ordinal features with 1-2-3 used for cool-warm-hot)
- **Multinomial NB**
 - Used for features containing counts and frequencies (useful for counting how often a word occurs in a text document)

Class	Long	Sweet	Yellow	Total
Banana	400	350	450	500
Orange	0	150	300	300
Other	100	150	50	200
Total	500	650	800	1000

We have a training set of 1000 fruits

Based on the features of “long”, “sweet” and “yellow”, the fruit can belong to one of the following classes: “Banana”, “Orange” or “Other”

- New data: {**long, sweet, yellow**} => What is the likelihood of this fruit being any one of the classes?
- **Step 1:** Let's calculate the probability that the fruit is a banana

$$P(\text{Banana} \mid \text{Long, Sweet, Yellow}) = P(B \mid LSY) = \frac{P(LSY \mid B)P(B)}{P(LSY)} = ?$$

- Assume the features are conditionally independent (given B):

$$P(LSY \mid B) = P(L \mid B)P(S \mid B)P(Y \mid B) \quad (\text{Naive assumption of NB})$$

- **Step 2:** Let's work out the equations above and plug them in

Example – cont'd

Class	Long	Sweet	Yellow	Total
Banana	400	350	450	500
Orange	0	150	300	300
Other	100	150	50	200
Total	500	650	800	1000

$$\begin{aligned}
 P(\text{Long} \mid \text{Banana}) &= 400 / 500 = 0.8 \\
 P(\text{Sweet} \mid \text{Banana}) &= 350 / 500 = 0.7 \\
 P(\text{Yellow} \mid \text{Banana}) &= 450 / 500 = 0.9 \\
 P(\text{Banana}) &= 500 / 1000 = 0.5
 \end{aligned}$$

$$P(\text{Banana} \mid \text{Long, Sweet, Yellow}) = P(B \mid \text{LSY}) = \frac{P(\text{LSY} \mid B)P(B)}{P(\text{LSY})} = ?$$

$$P(\text{LSY} \mid B) = 0.8 * 0.7 * 0.9 = 0.504$$

$$P(B) = 0.5$$

$$P(\text{Banana} \mid \text{LSY}) = \frac{0.504 * 0.5}{P(\text{LSY})} = \frac{0.252}{P(\text{LSY})} = \frac{0.252}{0.27075} = 93\%$$

$$\begin{aligned}
 P(\text{LSY}) &= P(\text{LSY} \mid B)P(B) + P(\text{LSY} \mid \text{Orange})P(\text{Orange}) \\
 &\quad + P(\text{LSY} \mid \text{Other})P(\text{Other})
 \end{aligned}$$

- Carrying out the same computation for the other two classes:

$$P(\text{Orange} \mid \text{LSY}) = 0$$

$$P(\text{Other} \mid \text{LSY}) = \frac{0.019}{P(\text{LSY})} = 7\%$$

} Given the current attributes, Naive-Bayes classifies the fruit as Banana

A classification example (2)

- Suppose we're building a tweet classifier which decides whether a tweet is "**positive**" or "**negative**"
- Here is our training data composed of 5 tweets:

tweet	Text	Class
1	I loved the movie	+
2	hated the movie	-
3	great movie, great acting	+
4	a poor script	-
5	bad script, bad movie	-

- Question: Which class does the tweet "**hated the poor acting**" belong to? (+) or (-)?

A classification example (2) – cont'd

- Since NB is a probabilistic classifier, we need to calculate the probability that the tweet "**hated the poor acting**" is "+" or "-", i.e.,
$$P(+ \mid \text{hated the poor acting}) = ?$$
$$P(- \mid \text{hated the poor acting}) = ?$$
- We need to somehow convert this text into numbers that we can do calculations on. We will
 - treat each tweet as a set of words it contains
 - ignore the word order and sentence construction and use "**word frequencies**"
- Features: Words
- Data: Word frequency (count) for each tweet (row)
- Too simplistic? Yes, but it works surprisingly well!

A classification example (2) – cont'd

- Number of unique words (unigrams) forming the vocabulary **V** : 11
{ I, loved, the, movie, hated, a, great, poor, acting, bad, script }
- We'll now convert the text documents into numerical feature vectors, where features are possible words.
- This process of turning text documents into a feature vector is called "**vectorization**".
- This specific strategy (tokenization, counting and normalization) is called the "**Bag of Words**" representation.

#	I	loved	the	movie	hated	a	great	poor	acting	bad	script	class
1	1	1	1	1		(values are the number of times a word occurs in the given document)						+
2			1	1	1							-
3				1			2		1			+
4						1		1			1	-
5				1						2	1	-

A classification example (2) – cont'd

- In terms of conditional probabilities:

$$P(+ \mid \textit{hated the poor acting}) = \frac{P(\textit{hated the poor acting} \mid +)P(+)}{P(\textit{hated the poor acting})}$$

$$P(- \mid \textit{hated the poor acting}) = \frac{P(\textit{hated the poor acting} \mid -)P(-)}{P(\textit{hated the poor acting})}$$

- We just need to compare the numerators for both and take the one with the greater probability for concluding the class of this new text.
- **Problem:** We don't have "*hated the poor acting*" in our training data.
- **Solution:** Assume every word in a sentence is independent of the other words. Instead of looking at the entire sentence, we'll look at individual words.
- A strong but super useful assumption (thus Naiveness)

A classification example (2) – cont'd

- So we write $P(\text{hated the poor acting})$ as:

$$P(\text{hated the poor acting}) = P(\text{hated}) P(\text{the}) P(\text{poor}) P(\text{acting})$$

- Applying this assumption to what we had before:

$$P(\text{hated the poor acting} | -) = P(\text{hated} | -) \times P(\text{the} | -) \times P(\text{poor} | -) \times P(\text{acting} | -)$$

- All of these individual words actually show up multiple times in our training data and they can be calculated:

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

$$P(\text{hated} | -) = \frac{\text{how many times "hated" appears in "-"}}{\text{total number of words in "-"}} = \frac{P(\text{hated} \cap -)}{P(-)} \xrightarrow{\text{blue}} \frac{1}{10}$$

$$P(\text{hated the poor acting} | -) = P(\text{hated} | -) \times P(\text{the} | -) \times P(\text{poor} | -) \times P(\text{acting} | -)$$

$$P(\text{hated the poor acting} | -) = 0$$

?

doesn't appear in "-"
So, $P(\text{acting} | -) = 0$

If one of the terms is zero, the whole expression is nullified!

A classification example (2) – cont'd

- Solution? Laplace smoothing:** Add "1" to every count so that none of the terms is zero. So, add "1" to all possible words in the divisor as well (to ensure a ratio not greater than 1)
- The possible words (entire dictionary of 11) are:
 $V = \{ \overset{1}{\text{'I'}}, \text{'loved'}, \text{'the'}, \text{'movie'}, \text{'hated'}, \text{'great'}, \text{'a'}, \text{'acting'}, \text{'poor'}, \text{'script'}, \overset{11}{\text{'bad'}} \}$
- With Laplace Smoothing we get:

word	P(word +)	P(word -)
hated	$\frac{0 + 1}{8 + 11}$	$\frac{1 + 1}{10 + 11}$
the	$\frac{1 + 1}{8 + 11}$	$\frac{1 + 1}{10 + 11}$
poor	$\frac{0 + 1}{8 + 11}$	$\frac{1 + 1}{10 + 11}$
acting	$\frac{1 + 1}{8 + 11}$	$\frac{0 + 1}{10 + 11}$

$(\text{poor} \cap +) \rightarrow 0 + 1$ +1 for "poor"
 $p(\text{poor} | +) = \frac{0 + 1}{8 + 11} = 0.0526$
number of all words in class '+' $\rightarrow 8 + 11$ +11 for every word in vocabulary

$$p(w_k | +) = \frac{\text{count}(w_k | +)}{\sum_{w \in V} \text{count}(w_k, +)}$$

$$p(w_k | +) = \frac{\text{count}(w_k | +) + 1}{\sum_{w \in V} [\text{count}(w_k, +) + 1]}$$

number of times word w_k occurs in class "+"

$$p(w_k | +) = \frac{\text{count}(w_k | +) + 1}{\sum_{w \in V} \text{count}(w_k, +) + |V|}$$

number of all words in class "+" $\rightarrow \sum_{w \in V} \text{count}(w_k, +)$
number of unique words in V $\rightarrow |V|$

A classification example (2) – cont'd

- We can now multiply all the probabilities:

$$P(\text{hated the poor acting} | +) = P(\text{hated} | +) \times P(\text{the} | +) \times P(\text{poor} | +) \\ \times P(\text{acting} | +) = 3.07 \cdot 10^{-5}$$

$$P(\underbrace{\text{hated the poor acting}}_{\text{tweet}} | -) = P(\text{hated} | -) \times P(\text{the} | -) \times P(\text{poor} | -) \\ \times P(\text{acting} | -) = 4.11 \cdot 10^{-5}$$

$$P(+ | \text{tweet}) = \frac{P(\text{tweet} | +) P(+)}{P(\text{tweet})} = \frac{P(\text{tweet} | +) P(+)}{P(\text{tweet} | +) P(+)+P(\text{tweet} | -) P(-)} \quad \leftarrow \frac{2}{5}$$
$$= 33.2\%$$

$$P(- | \text{tweet}) = \frac{P(\text{tweet} | -) P(-)}{P(\text{tweet})} = \frac{P(\text{tweet} | -) P(-)}{P(\text{tweet} | +) P(+)+P(\text{tweet} | -) P(-)} \quad \leftarrow \frac{3}{5}$$

= 66.8% ✓

NB classifies "hated the poor acting" as "-"

NB for continuous data

- NB is designed to work primarily with categorical attributes. In many practical cases, however, variables are often numerical. So how does NB handle numerical attributes?

Temperature	Humidity	Play
85	86	no
89	85	no
80	87	no
90	90	no
86	83	no
83	76	yes
70	88	yes
68	80	yes
64	65	yes
69	70	yes
75	80	yes
75	70	yes
72	90	yes
81	75	yes

Famous "tennis play" data

Method 1: Discretize the continuous variables into different categories using intervals (binning) and turn them into nominal values.

- e.g. temp. smaller than 70 as **cool**, 71-79 as **medium**, above 80 as **hot**
- You can then construct their frequency tables to compute the desired probabilities.
- This can be very subjective and cause loss of information, thus may not be suitable for what you need.

NB for continuous data – cont'd

- **Method 2:** Use probability density functions preserving the continuous values as given. One common practice is to assume Normal distributions for numerical variables

Temp.	Humidity	Play
85	85	no
80	90	no
65	70	no
72	95	no
71	80	no
83	78	yes
70	96	yes
68	80	yes
64	65	yes
69	70	yes
75	80	yes
75	70	yes
72	90	yes
81	75	yes

- Compute the mean and variance from the data set:

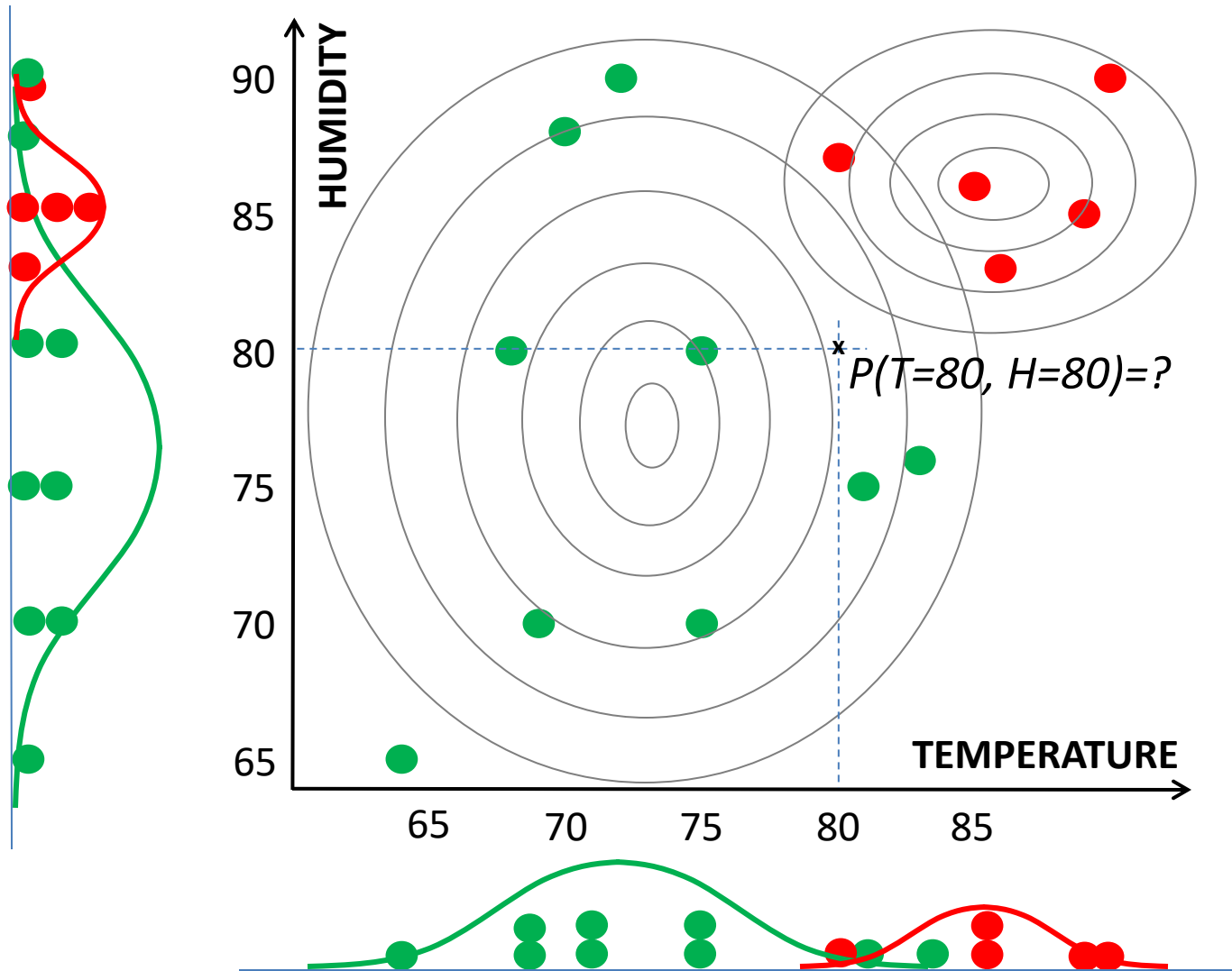
μ_T and σ_T for both “yes”
 μ_H and σ_H and “no” classes

- Use the normal distribution to compute probabilities:

$$p(temp = t \mid play = yes) = \frac{1}{\sqrt{2\pi}\sigma_{T,y}} e^{-\frac{(t-\mu_{T,y})^2}{2\sigma_{T,y}^2}}$$

$$p(hum = h \mid play = yes) = \frac{1}{\sqrt{2\pi}\sigma_{H,y}} e^{-\frac{(h-\mu_{H,y})^2}{2\sigma_{H,y}^2}}$$

NB for continuous data – cont'd



NB for continuous data – cont'd

- What is the class (probability) for playing tennis given input x where $x = x(\text{temp}=T_x, \text{humidity}=H_x)$?

$$p(H_x | \text{yes}) = \frac{1}{\sqrt{2\pi}\sigma_{H,y}} e^{-\frac{(H_x - \mu_{H,y})^2}{2\sigma_{H,y}^2}} \quad p(H_x | \text{no}) = \frac{1}{\sqrt{2\pi}\sigma_{H,n}} e^{-\frac{(H_x - \mu_{H,n})^2}{2\sigma_{H,n}^2}}$$

$$p(T_x | \text{yes}) = \frac{1}{\sqrt{2\pi}\sigma_{T,y}} e^{-\frac{(T_x - \mu_{T,y})^2}{2\sigma_{T,y}^2}} \quad p(T_x | \text{no}) = \frac{1}{\sqrt{2\pi}\sigma_{T,n}} e^{-\frac{(T_x - \mu_{T,n})^2}{2\sigma_{T,n}^2}}$$

$$p(x | \text{yes}) = p(H_x | \text{yes}) p(T_x | \text{yes}) \quad p(x | \text{no}) = p(H_x | \text{no}) p(T_x | \text{no})$$

$$p(\text{yes} | x) = \frac{p(x | \text{yes}) p(\text{yes})}{p(x | \text{yes}) p(\text{yes}) + p(x | \text{no}) p(\text{no})}$$

$$p(\text{no} | x) = \frac{p(x | \text{no}) p(\text{no})}{p(x | \text{yes}) p(\text{yes}) + p(x | \text{no}) p(\text{no})}$$

NB and missing values

- Handling of missing values with NB is easy: Just skip it!
- We can simply ignore the instance for which the value is missing. Remember the expression for NB:

$$p(X_1, X_2, \dots, X_j, \dots, X_n | y) = \prod_{i \neq j}^n P(X_i | y)$$

Product goes over all attributes except "j", the one that's skipped (missing)

Tossing a coin three times: $E = \{X_1 = \text{Head}, X_2 = \text{missing}, X_3 = \text{Tail}\}$
 $p(E) = p(\text{H}, \text{H}, \text{T}) + p(\text{H}, \text{T}, \text{T})$ as the missing value is either **H** or **T**
 $P(E) = 1/8 + 1/8 = 1/4$ (same as $p(\text{H}, \text{T}) = 1/4$)

$$p(X_1, \dots, X_j, \dots, X_n | y) = P(X_1 | y) \dots P(X_j | y) \dots P(X_n | y)$$

$$\sum_{X_j} p(X_1, \dots, X_j, \dots, X_n | y) = \sum_{X_j} P(X_1 | y) \dots P(X_j | y) \dots P(X_n | y)$$

over all possible values of X_j

$$= P(X_1 | y) \dots \sum_{X_j} P(X_j | y) \dots P(X_n | y) = \prod_{i \neq j}^n P(X_i | y)$$

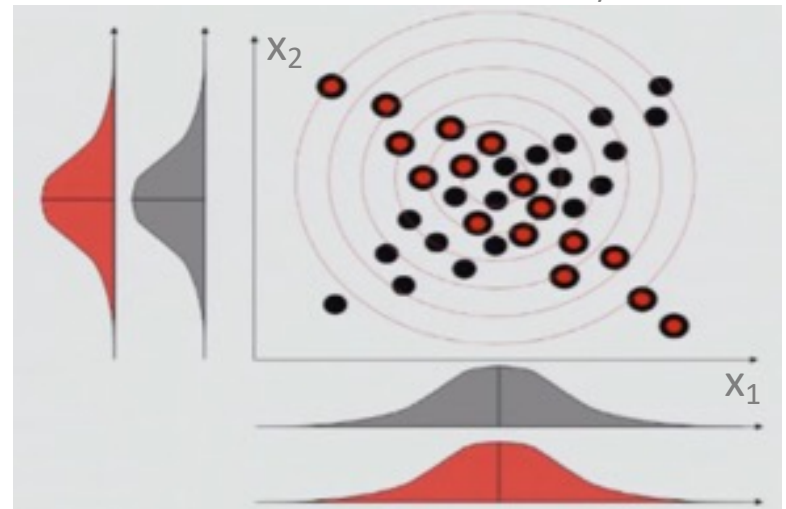
1

- **Zero Frequency (zero observation) problem**
 - When a categorical feature has a value in the test set, but not in the training set, the model assigns a 0 probability and is unable to make a prediction.
 - To fix it, we apply a smoothing technique (Laplace correction) and add 1 to every feature. This is essentially assigning an arbitrarily low probability to that class to make it different from zero.
- **Not much "fine-tuning"**
 - NB classifiers have very limited options for parameter tuning. So, it's best to focus on data pre-processing and feature selection.
 - In case overfitting occurs, one can tweak the priors.

- **No combination of classifiers**
 - Combination of classifiers (like ensembling, bagging and boosting) won't help with NB as the purpose with these techniques is to reduce the variance. NB has no variance to minimize (no loss function to minimize).
 - NB employs a simple hypothesis function with no dependencies in features (unable to model feature interactions unlike many other classifiers), hence, it's a high-bias/low-variance model that is unable to handle complex problems.
- **Deviations from normality**
 - For continuous features that don't follow a Normal distribution, try variable transformation.

- **Inversely correlated attributes**
 - Once you project all data points on both axes, distributions for 2 classes are on top of each other (almost identical) providing no extra prediction power. So, NB can't distinguish between the classes. What differentiates these classes isn't the means or variances. It's the covariance. For black class, x_1 and x_2 are positively correlated, for red class, they are negatively correlated. This is the only difference between the classes. NB cannot deal with correlation. It assumes conditional independence (naive). If the only thing that separates the 2 classes is the way in which attributes are correlated, you're out of luck.

Victor Lavrenko IAML5.9 - Gaussian Naive Bayes classifier



- **Double-counting effect (highly-correlated features)**
 - When both features encode essentially the same information, the classifier "double-counts" the effect of highly correlated features (voting them twice overinflates their importance), and gets pushed closer to a given label than is justified. This is particularly true for large and sparse data (text mining).
 - Of course, we don't usually build NB classifiers that contain two identical features. However, we do build classifiers that contain features which are dependent on one another. For example, the features ends-with **(a)** and ends-with **(vowel)** are dependent on each other. Because if an input value has the first feature, then it must also have the second feature. For features like these, the duplicated information may be given more weight than is justified by the training set.

- **Presence of mixed data (continuous + categorical)**
 - GaussianNB (continuous), BernoulliNB (binary), and MultinomialNB (nominal) are employed for different types of data. How to proceed when we have mixed data?
 - A few approaches (see the Notebook for some examples):
 - Use binning with intervals that contain roughly equal number of data points for continuous variables and turn them into categorical variables (loss of information)
 - Compute the probabilities for each feature type independently (due to nature of the underlying assumption of NB) and multiply them together (by taking care of normalization)
 - Fit different models for different features and then transform the data set by taking the class assignment probabilities as new features, and then refitting a new model (say GaussianNB) on the new features.

- **Use log probabilities**
 - Probabilities can be small numbers and often lead to inaccurate results (due to precision of floating-point values) when multiplied together. It is best to work in the logarithmic realm and add them up rather than multiply them. You can convert it back to original later.
- **Remove redundant features**
 - Due to possible effect of double-counting of correlated features, you can remove one of the features.
- **One vs Rest approach**
 - For features that have many categorical levels, you can try a two-class feature (category A vs all-else) and keep experimenting with different combinations

- **Pros**

- A **very good choice for a small amount of training data** (as the estimates are based on the joint density function)
- Reliably a high bias/low variance classifier. Try for problems with noisy data, lots of features, and few samples
- Despite its simplicity, surprisingly accurate
- Easy, **very fast to train** (single scan) and classify
- **Not sensitive to irrelevant features**
- Works with both categorical and continuous data (requires work)
- Handles multiple classes
- **Tolerant to missing values** (not implemented in scikit-learn)
- Can also be **used to generate instances of the problem**

- **Cons**

- Assumes independence of features (cannot learn interactions between features – fits feature weights independently)
- Suffers from multi-collinearity
- Sensitive to how input is prepared

using scikit-learn Library

```
from sklearn.naive_bayes import GaussianNB
from sklearn.naive_bayes import BernoulliNB
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import accuracy_score
```

```
Clf = GaussianNB()
Clf.fit(X_train, y_train)
accuracy_score(y_test, Clf.predict(X_test))
```


Addendum

Bias-Variance property of NB

- **Why is NB a high-bias/low variance classifier?**
- The Naive Bayes classifier employs a very simple (linear) hypothesis function to model data. It suffers from high bias, or error resulting from inaccuracies in its hypothesis class, because its hypothesis function is so simple it cannot accurately represent many complex situations.
- On the other hand, it exhibits low variance, or failure to generalize to unseen data based on its training set, because its hypothesis class' simplicity prevents it from overfitting to its training data.

Is NB a linear classifier?

- **Is NB a linear classifier? What is a decision boundary like?**
 - Suppose we have two populations whose distributions are Gaussian: $N(\mu_0, \sigma_0^2)$ and $N(\mu_1, \sigma_1^2)$
 - Want to classify whether a person is female ($c=0$) or male ($c=1$) based on height

$$p(x | c = 0) = \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} \quad p(x | c = 1) = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}}$$

- What we need is:

$$\hat{c}(x) = ?$$

$$\hat{c}(x) = 1 \text{ if } p(c = 1 | x) > 0.5; \quad \hat{c}(x) = 0 \text{ otherwise}$$

- We'll assume equal probability of being male or female

$$p(c = 1) = 0.5$$

- This is the condition that will form the decision boundary

Is NB a linear classifier?

- Probability of class=male given the feature:

$$p(c = 1 | x) \propto p(x | c = 1) p(c = 1)$$

- Let's assume $p(c = 1) = p_1$

$$\log[p(x | c = 1) p(c = 1)] = \log\left[\frac{1}{\sqrt{2\pi}\sigma_1}\right] - \frac{1}{2\sigma_1^2} (x - \mu_1)^2 + \log p_1$$

$$\log[p(x | c = 0) p(c = 0)] = \log\left[\frac{1}{\sqrt{2\pi}\sigma_0}\right] - \frac{1}{2\sigma_0^2} (x - \mu_0)^2 + \log p_0$$

- The decision boundary between any two classes is where the log of the ratio (log odds) is zero (why?)

$$\log \frac{p(x | c = 1) p(c = 1)}{p(x | c = 0) p(c = 0)} = -\frac{1}{2\sigma_1^2} (x^2 - 2\mu_1 x) + \frac{1}{2\sigma_0^2} (x^2 - 2\mu_0 x)$$

constant

$$-\frac{\mu_1^2}{2\sigma_1^2} + \frac{\mu_0^2}{2\sigma_0^2} + \log\left[\frac{1}{\sqrt{2\pi}\sigma_1}\right] - \log\left[\frac{1}{\sqrt{2\pi}\sigma_0}\right] + \log p_1 - \log p_0$$

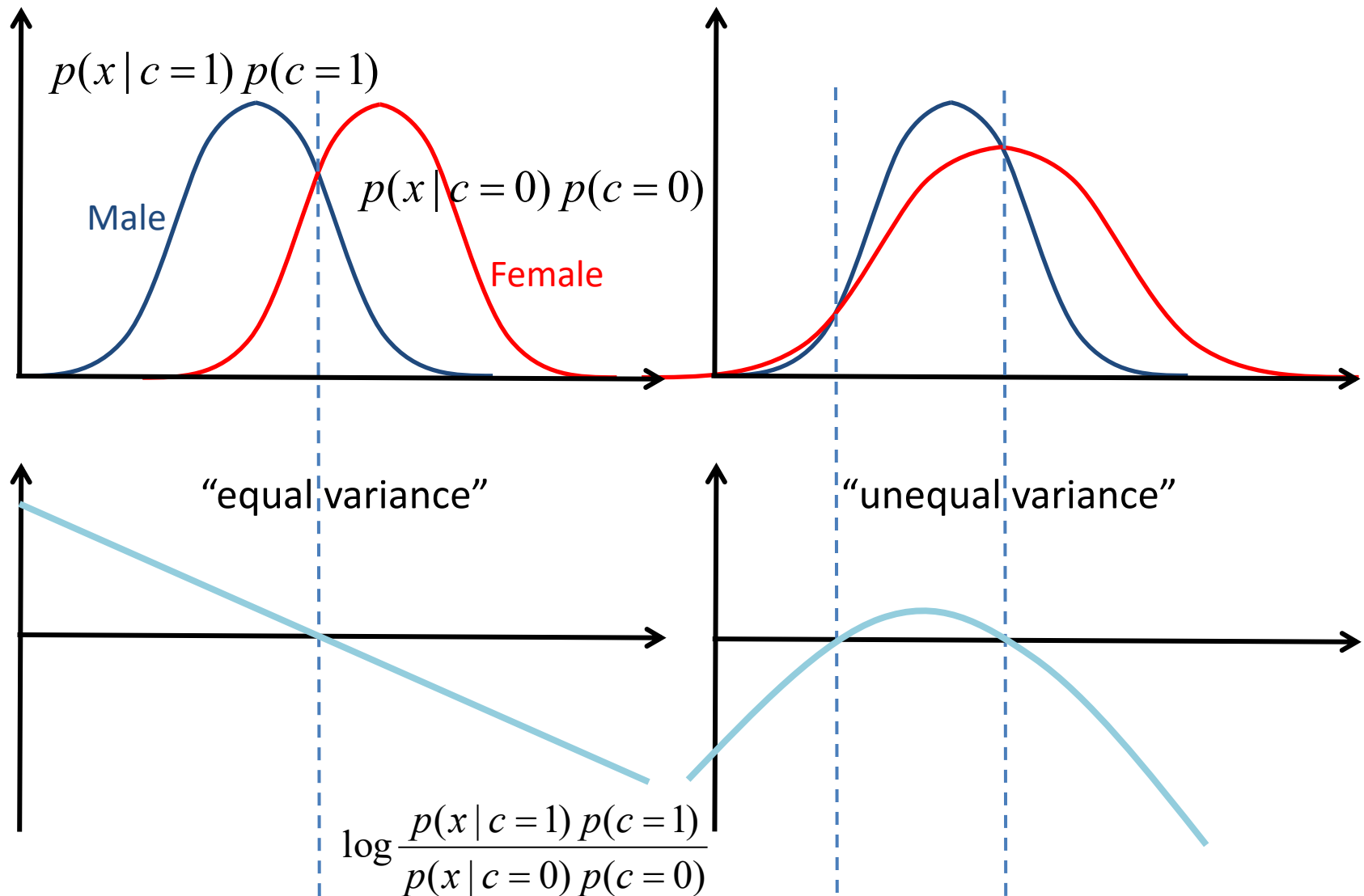
Is NB a linear classifier?

- So the log of the ratio becomes:

$$\log \frac{p(x | c = 1) p(c = 1)}{p(x | c = 0) p(c = 0)} = \frac{1}{2} \left(\frac{1}{\sigma_0^2} - \frac{1}{\sigma_1^2} \right) x^2 + \left(\frac{\mu_1}{\sigma_1^2} - \frac{\mu_0}{\sigma_0^2} \right) x + \mathbf{constant}$$

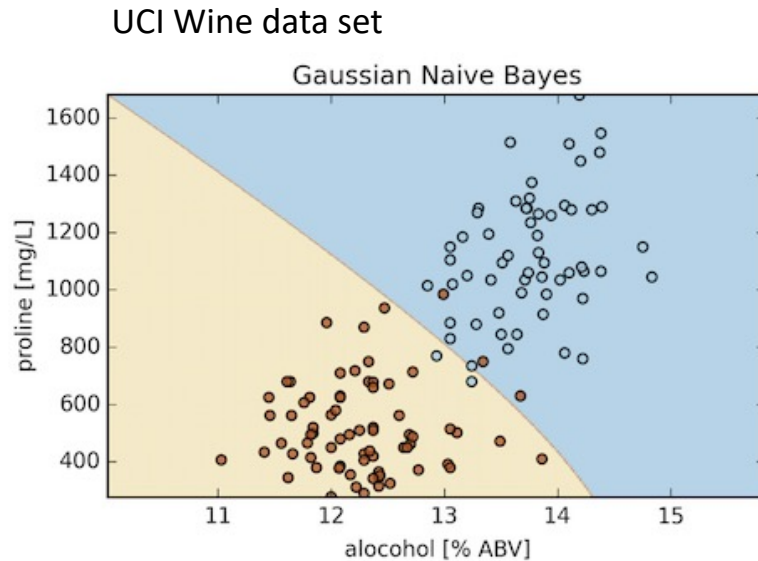
- If the data in each class has a Gaussian density with **unequal variance**, then the boundary between any two classes is a **quadratic function** of x .
- If the data in each class has a Gaussian density with **equal variance**, the coefficient of the quadratic term disappears and the boundary between any two classes becomes a **line**.
- Summary for the decision boundary:
 - General case: parabolic curve
 - Different means, same variance: Line (or plane)
 - Same mean, different variance: circle / ellipse
- See the next slide for graphical representations

Is NB a linear classifier?

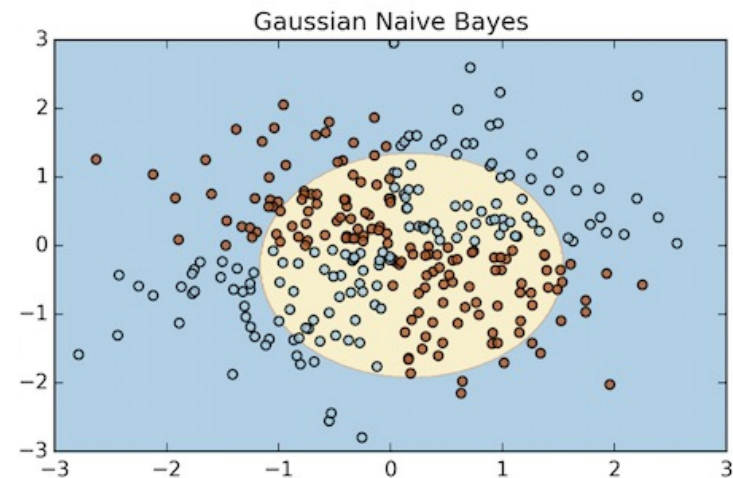


Is NB a linear classifier?

- So, it's a (piecewise) quadratic decision boundary for the Gaussian model. The multinomial model has a linear boundary. Here are some plots taken from S. Raschka:



An XOR toy data set



Ref: <https://www.quora.com/What-is-the-decision-boundary-for-Naive-Bayes>

- How do to deal with mixed (categorical + continuous) data in NB?
- There are 3 types of formulations for NB in scikit-learn:
 1. BernoulliNB (binary features)
 2. MultinomialNB (nominal features)
 3. GaussianNB (continuous features)
- In the presence of mixed data, you cannot mix the features. We have a few options:
- **1.** Convert your continuous data into a categorical representation by computing percentiles for each continuous variable and then bin those percentiles as categorical levels. For example, you can take the temperature for and create the following bins: "cold", "cool", "normal", "warm", "hot" ensuring that each bin contains approximately 20% of the population of your training set (not recommended as discretization will throw away useful information contained in a continuously varying data)

Mixed data in NB – cont'd

- **2.** You can independently fit a GaussianNB model on the continuous part of the data and a MultinomialNB (BernoulliNB depending on the nature of the categorical data) model on the categorical part. Then transform all the dataset by taking the class assignment probabilities (with `predict_proba` method) as new features as can be seen from the command down below. You can then easily refit a new model (e.g. a new GaussianNB) on these new features.

```
np.hstack((multinomial_probas, gaussian_probas))
```

Ref: <https://stackoverflow.com/questions/14254203/mixing-categorical-and-continuous-data-in-naive-bayes-classifier-using-scikit-learn>

- **3.** The last option is to find the individual probability values (as in 2) for both the continuous and categorical features. Let's assume that we have two features (feature1 is continuous, and feature2 is categorical). Here is the probability we need to calculate:

- Example data set

feature1	feature2	class
126	low	c1
220	low	c2
191	medium	c1
128	high	c1
125	low	c2
...

$$p(c_1 | f_1, f_2) = \frac{p(f_1, f_2 | c_1)p(c_1)}{p(f_1, f_2)}$$

$$p(c_1 | f_1, f_2) = \frac{p(f_1 | c_1)p(f_2 | c_1)p(c_1)}{p(f_1, f_2 | c_1)p(c_1) + p(f_1, f_2 | c_2)p(c_2)} = \frac{p(f_1 | c_1)p(f_2 | c_1)p(c_1)}{D}$$

Similarly:

$$p(c_2 | f_1, f_2) = \frac{p(f_1 | c_2)p(f_2 | c_2)p(c_2)}{D}$$

Mixed data in NB – cont'd

- By using the GaussianNB on *feature1* , and MultinomialNB on *feature2* based on the target values, we get the probabilities (predict_proba method) separately:

$$p(c_1 | f_1) = \frac{p(f_1 | c_1)p(c_1)}{p(f_1 | c_1)p(c_1) + p(f_1 | c_2)p(c_2)} = \frac{p(f_1 | c_1)p(c_1)}{P_1} = p_{11} \Rightarrow p(f_1 | c_1) = \frac{p_{11}}{p(c_1)} P_1$$

$$p(c_1 | f_2) = \frac{p(f_2 | c_1)p(c_1)}{p(f_2 | c_1)p(c_1) + p(f_2 | c_2)p(c_2)} = \frac{p(f_2 | c_1)p(c_1)}{P_2} = p_{12} \Rightarrow p(f_2 | c_1) = \frac{p_{12}}{p(c_1)} P_2$$

$$p(c_2 | f_1) = \frac{p(f_1 | c_2)p(c_2)}{p(f_1 | c_1)p(c_1) + p(f_1 | c_2)p(c_2)} = \frac{p(f_1 | c_2)p(c_2)}{P_1} = p_{21} \Rightarrow p(f_1 | c_2) = \frac{p_{21}}{p(c_2)} P_1$$

$$p(c_2 | f_2) = \frac{p(f_2 | c_2)p(c_2)}{p(f_2 | c_1)p(c_1) + p(f_2 | c_2)p(c_2)} = \frac{p(f_2 | c_2)p(c_2)}{P_2} = p_{22} \Rightarrow p(f_2 | c_2) = \frac{p_{22}}{p(c_2)} P_2$$

- Now going back to $p(c_1 | f_1, f_2)$ and $p(c_2 | f_1, f_2)$ we find:

$$p(c_1 | f_1, f_2) = \frac{p(f_1 | c_1)p(f_2 | c_1)p(c_1)}{D} = \frac{P_1 P_2}{D} \frac{p_{11} p_{12}}{p(c_1)} \quad , \text{ and}$$

$$p(c_2 | f_1, f_2) = \frac{p(f_1 | c_2)p(f_2 | c_2)p(c_2)}{D} = \frac{P_1 P_2}{D} \frac{p_{21} p_{22}}{p(c_2)}$$

- Normalizing probabilities we get:

$$p(c_1 | f_1, f_2) = \frac{p(c_1 | f_1, f_2)}{p(c_1 | f_1, f_2) + p(c_2 | f_1, f_2)} = \frac{\frac{p_{11}p_{12}}{p(c_1)}}{\frac{p_{11}p_{12}}{p(c_1)} + \frac{p_{21}p_{22}}{p(c_2)}}$$

$$p(c_2 | f_1, f_2) = \frac{p(c_2 | f_1, f_2)}{p(c_1 | f_1, f_2) + p(c_2 | f_1, f_2)} = \frac{\frac{p_{21}p_{22}}{p(c_2)}}{\frac{p_{11}p_{12}}{p(c_1)} + \frac{p_{21}p_{22}}{p(c_2)}}$$

- where p_1 and p_2 are the priors for the target variable and can easily be calculated from the training data. p_{11} , p_{12} , p_{21} , and p_{22} are the probabilities computed by fitting a model to each feature (continuous or categorical) using the relevant NB model from scikit-learn library and then executing the "predict_proba" command separately.