

# CS440/ECE448 Lecture 17: Bayesian Inference

Slides by Svetlana Lazebnik, 10/2016

Modified by Mark Hasegawa-Johnson, 10/2017

# Review: Probability

- Random variables, events
- Axioms of probability
- Atomic events
- Joint and marginal probability distributions
- Conditional probability distributions
- Product rule, chain rule
- Independence and conditional independence

# Outline: Bayesian Inference

- Bayes Rule
- Law of Total Probability
- Misdiagnosis
- $\text{MAP} = \text{MPE}$
- The “Naïve Bayesian” Assumption
- Bag of Words (BoW)
- Parameter Estimation for the BoW model

# Bayes Rule



Rev. Thomas Bayes  
(1702-1761)

- The product rule gives us two ways to factor a joint probability:

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

- Therefore, 
$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$
- Why is this useful?
  - Can update our beliefs about A based on evidence B
    - $P(A)$  is the *prior* and  $P(A | B)$  is the *posterior*
  - Key tool for probabilistic inference: can get *diagnostic probability* from *causal probability*
    - E.g.,  $P(\text{Cavity} = \text{true} | \text{Toothache} = \text{true})$  from  $P(\text{Toothache} = \text{true} | \text{Cavity} = \text{true})$

# Bayes Rule example

Dan & Dana are getting married tomorrow, at an outdoor ceremony in the desert. In recent years, it has rained only 5 days each year ( $5/365 = 0.014$ ). Unfortunately, the weatherman has predicted rain for tomorrow. When it actually rains, the weatherman correctly forecasts rain 90% of the time. When it doesn't rain, he incorrectly forecasts rain 10% of the time. What is the probability that it will rain on their wedding?

# Bayes Rule example

Dan & Dana are getting married tomorrow, at an outdoor ceremony in the desert. In recent years, it has rained only 5 days each year ( $5/365 = 0.014$ ). Unfortunately, the weatherman has predicted rain for tomorrow. When it actually rains, the weatherman correctly forecasts rain 90% of the time. When it doesn't rain, he incorrectly forecasts rain 10% of the time. What is the probability that it will rain on their wedding?

$$P(\text{rain}|\text{forecast}) = \frac{0.9 * 0.014}{P(\text{forecast})}$$

# Bayes Rule example

Dan & Dana is getting married tomorrow, at an outdoor ceremony in the desert. In recent years, it has rained only 5 days each year ( $5/365 = 0.014$ ). Unfortunately, the weatherman has predicted rain for tomorrow. When it actually rains, the weatherman correctly forecasts rain 90% of the time. When it doesn't rain, he incorrectly forecasts rain 10% of the time. What is the probability that it will rain on their wedding?

$$P(\text{rain}|\text{forecast}) = \frac{0.9 \times 0.014}{0.9 * 0.014 + 0.1 * 0.986}$$

# Outline: Bayesian Inference

- Bayes Rule
- Law of Total Probability
- Misdiagnosis
- MAP = MPE
- The “Naïve Bayesian” Assumption
- Bag of Words (BoW)
- Parameter Estimation for the BoW model



## Law of total probability

$$\begin{aligned} P(X = x) &= \sum_{i=1}^n P(X = x, Y = y_i) \\ &= \sum_{i=1}^n P(X = x | Y = y_i) P(Y = y_i) \end{aligned}$$

# Bayes Rule example

Dan & Dana is getting married tomorrow, at an outdoor ceremony in the desert. In recent years, it has rained only 5 days each year ( $5/365 = 0.014$ ). Unfortunately, the weatherman has predicted rain for tomorrow. When it actually rains, the weatherman correctly forecasts rain 90% of the time. When it doesn't rain, he incorrectly forecasts rain 10% of the time. What is the probability that it will rain on their wedding?

$$\begin{aligned} P(\text{rain} \mid \text{predict}) &= \frac{P(\text{predict} \mid \text{rain})P(\text{rain})}{P(\text{predict})} \\ &= \frac{P(\text{predict} \mid \text{rain})P(\text{rain})}{P(\text{predict} \mid \text{rain})P(\text{rain}) + P(\text{predict} \mid \neg\text{rain})P(\neg\text{rain})} \\ &= \frac{0.9 \times 0.014}{0.9 \times 0.014 + 0.1 \times 0.986} = \frac{0.0126}{0.0126 + 0.0986} = 0.111 \end{aligned}$$

# Outline: Bayesian Inference

- Bayes Rule
- Law of Total Probability
- **Misdiagnosis**
- $\text{MAP} = \text{MPE}$
- The “Naïve Bayesian” Assumption
- Bag of Words (BoW)
- Parameter Estimation for the BoW model

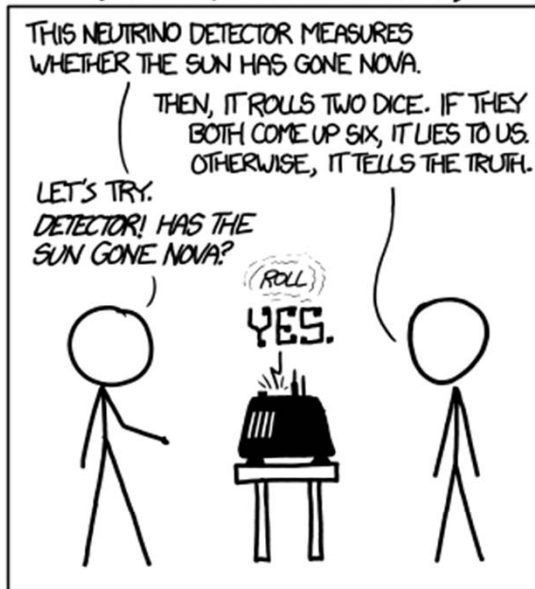
# Bayes rule: Example

1% of women at age forty who participate in routine screening have breast cancer. 80% of women with breast cancer will get positive mammographies. 9.6% of women without breast cancer will also get positive mammographies. A woman in this age group had a positive mammography in a routine screening. What is the probability that she actually has breast cancer?

$$\begin{aligned}P(\text{cancer} \mid \text{positive}) &= \frac{P(\text{positive} \mid \text{cancer})P(\text{cancer})}{P(\text{positive})} \\&= \frac{P(\text{positive} \mid \text{cancer})P(\text{cancer})}{P(\text{positive} \mid \text{cancer})P(\text{cancer}) + P(\text{positive} \mid \neg \text{cancer})P(\neg \text{Cancer})} \\&= \frac{0.8 \times 0.01}{0.8 \times 0.01 + 0.096 \times 0.99} = \frac{0.008}{0.008 + 0.095} = 0.0776\end{aligned}$$

<https://www.youtube.com/watch?v=BcvLAW-JRss>

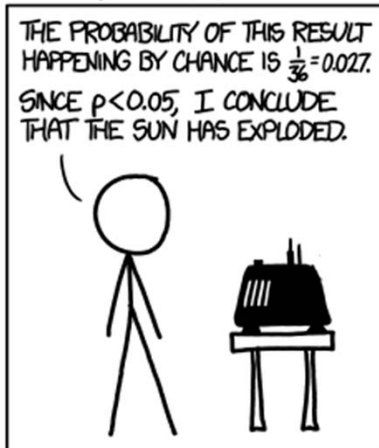
DID THE SUN JUST EXPLODE?  
(IT'S NIGHT, SO WE'RE NOT SURE.)



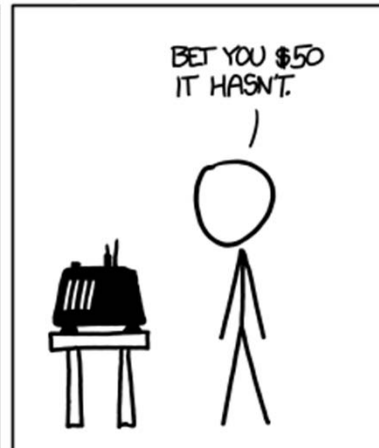
$$P(\text{nova} | \text{yes}) = \frac{P(\text{yes} | \text{nova})P(\text{nova})}{P(\text{yes})}$$

$$P(\neg \text{nova} | \text{yes}) = \frac{P(\text{yes} | \neg \text{nova})P(\neg \text{nova})}{P(\text{yes})}$$

FREQUENTIST STATISTICIAN:



BAYESIAN STATISTICIAN:



<https://xkcd.com/1132/>

See also: <https://xkcd.com/882/>

# Probabilistic inference

- Suppose the agent has to make a decision about the value of an unobserved *query variable*  $X$  given some observed *evidence variable(s)*  $E = e$ 
  - Partially observable, stochastic, episodic environment
  - Examples:  $X = \{\text{spam, not spam}\}$ ,  $e = \text{email message}$   
 $X = \{\text{zebra, giraffe, hippo}\}$ ,  $e = \text{image features}$

✗

Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...

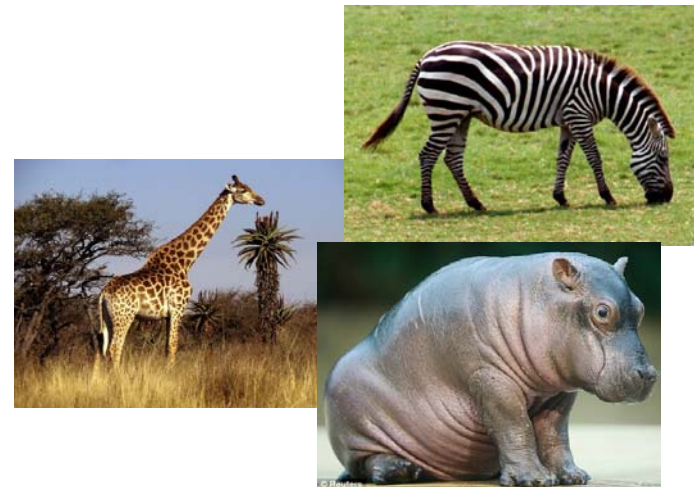
✗

TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES FOR ONLY \$99

✓

Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.



# Outline: Bayesian Inference

- Bayes Rule
- Law of Total Probability
- Misdiagnosis
- **MAP = MPE**
- The “Naïve Bayesian” Assumption
- Bag of Words (BoW)
- Parameter Estimation for the BoW model

# Bayesian decision theory

- Let  $x$  be the value predicted by the agent (known) and  $x^*$  be the true value of  $X$  (unknown).
- The agent has a **loss function**, which is 0 if  $x = x^*$  and 1 otherwise
- Expected loss for predicting  $x$ :

$$\sum_{x^*} L(x, x^*) P(x^* | o) = \sum_{x \neq x^*} P(x^* | o) = 1 - P(x | o)$$

- What is the estimate of  $X$  that minimizes the expected loss?
  - The one that has the greatest posterior probability  $P(x|o)$
  - This is called the **Maximum a Posteriori (MAP)** decision



# MAP decision

- Value  $x$  of  $X$  that has the highest posterior probability given the observation  $O=o$ :

$$\begin{aligned}x^* &= \operatorname{argmax} P(x|o) = \operatorname{argmax} \frac{P(o|x)P(x)}{P(o)} \\ &= \operatorname{argmax} P(o|x)P(x)\end{aligned}$$

$$\underset{\text{posterior}}{P(x|e)} \propto \underset{\text{likelihood}}{P(e|x)} \underset{\text{prior}}{P(x)}$$

- Maximum Likelihood (ML) decision:

$$x^* = \operatorname{argmax} P(o|x)$$

# Outline: Bayesian Inference

- Bayes Rule
- Law of Total Probability
- Misdiagnosis
- $\text{MAP} = \text{MPE}$
- The “Naïve Bayesian” Assumption
- Bag of Words (BoW)
- Parameter Estimation for the BoW model

# Naïve Bayes model

- Suppose we have many different types of observations (symptoms, features)  $E_1, \dots, E_n$  that we want to use to obtain evidence about an underlying hypothesis  $X$
- MAP decision:

$$\begin{aligned} P(X = x \mid E_1 = e_1, \dots, E_n = e_n) \\ \propto P(X = x)P(E_1 = e_1, \dots, E_n = e_n \mid X = x) \end{aligned}$$

- If each feature  $E_i$  can take on  $k$  values, how many entries are in the (conditional) joint probability table  $P(E_1, \dots, E_n \mid X = x)$ ?

# Naïve Bayes model

- Suppose we have many different types of observations (symptoms, features)  $E_1, \dots, E_n$  that we want to use to obtain evidence about an underlying hypothesis  $X$

- MAP decision:

$$P(X = x \mid E_1 = e_1, \dots, E_n = e_n) \\ \propto P(X = x)P(E_1 = e_1, \dots, E_n = e_n \mid X = x)$$

- We can make the simplifying assumption that the different features are **conditionally independent given the hypothesis**:

$$P(E_1 = e_1, \dots, E_n = e_n \mid X = x) = \prod_{i=1}^n P(E_i = e_i \mid X = x)$$

- If each observation and the hypothesis can take on  $k$  values, what is the complexity of storing the resulting distributions?
- W.o naïve Bayes:  $k(k^n - 1)$
- With naïve Bayes: each  $p(e|x)$  requires  $(k-1)*k$  ( $k$  values of  $x$ ,  $k-1$  of  $e$ )
- There are  $n$  of them  $\rightarrow n*(k-1)*k$

# Naïve Bayes model

- Posterior:

$$P(X = x \mid E_1 = e_1, \dots, E_n = e_n)$$

- MAP decision:

$$x^* = \arg \max_x \underbrace{P(x \mid e)}_{\text{posterior}} \propto \underbrace{P(x)}_{\text{prior}} \underbrace{\prod_{i=1}^n P(e_i \mid x)}_{\text{likelihood}}$$

# Case study: Text document classification

- **MAP decision:** assign a document to the class with the highest posterior  $P(\text{class} \mid \text{document})$
- Example: spam classification
  - Classify a message as spam if  $P(\text{spam} \mid \text{message}) > P(\neg\text{spam} \mid \text{message})$



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...



TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES  
FOR ONLY \$99



Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

# Case study:

## Text document classification

- **MAP decision:** assign a document to the class with the highest posterior  $P(\text{class} \mid \text{document})$
- We have  $P(\text{class} \mid \text{document}) \propto P(\text{document} \mid \text{class})P(\text{class})$
- To enable classification, we need to be able to estimate the **likelihoods**  $P(\text{document} \mid \text{class})$  for all classes and **priors**  $P(\text{class})$

# Outline: Bayesian Inference

- Bayes Rule
- Law of Total Probability
- Misdiagnosis
- $\text{MAP} = \text{MPE}$
- The “Naïve Bayesian” Assumption
- **Bag of Words (BoW)**
- Parameter Estimation for the BoW model



# Naïve Bayes Representation

- Goal: estimate likelihoods  $P(\text{document} \mid \text{class})$  and priors  $P(\text{class})$
- Likelihood: ***bag of words*** representation
  - The document is a sequence of words  $(w_1, \dots, w_n)$
  - The order of the words in the document is not important
  - Each word is conditionally independent of the others given document class



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...



TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES  
FOR ONLY \$99



Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

# Naïve Bayes Representation

- Goal: estimate likelihoods  $P(\text{document} \mid \text{class})$  and priors  $P(\text{class})$
- Likelihood: **bag of words** representation
  - The document is a sequence of words  $(w_1, \dots, w_n)$
  - The order of the words in the document is not important
  - Each word is conditionally independent of the others given document class

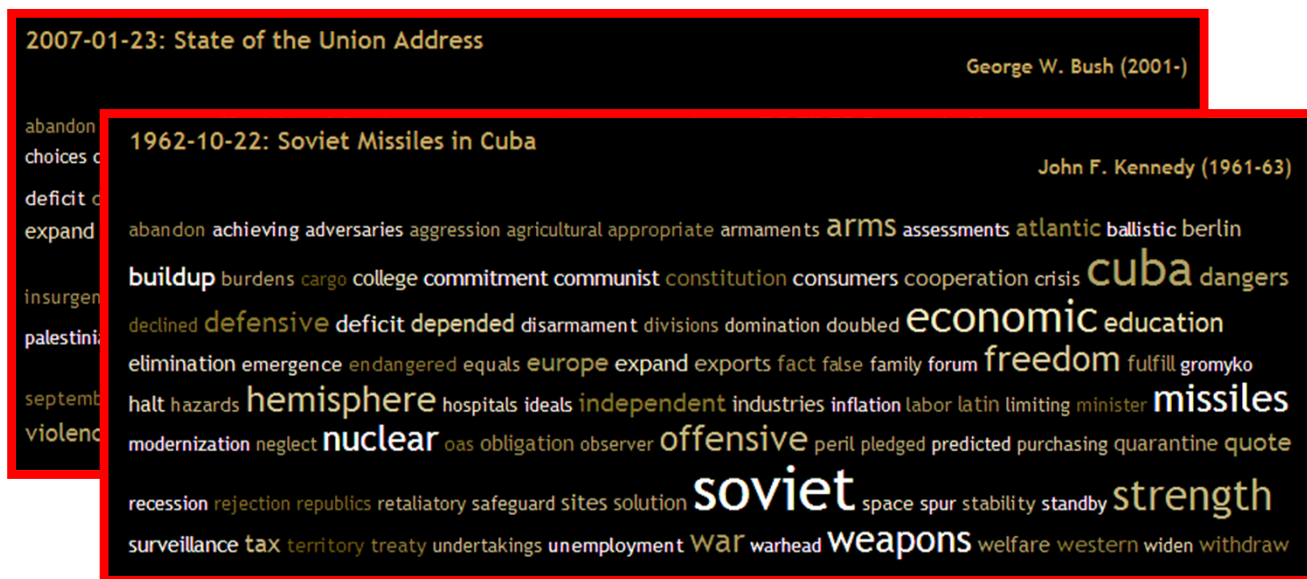
$$P(\text{document} \mid \text{class}) = P(w_1, \dots, w_n \mid \text{class}) = \prod_{i=1}^n P(w_i \mid \text{class})$$

# Bag of words illustration



US Presidential Speeches Tag Cloud  
<http://chir.ag/projects/preztags/>

# Bag of words illustration



US Presidential Speeches Tag Cloud  
<http://chir.ag/projects/preztags/>

# Bag of words illustration



US Presidential Speeches Tag Cloud  
<http://chir.ag/projects/preztags/>

# Naïve Bayes Representation

- Goal: estimate likelihoods  $P(\text{document} \mid \text{class})$  and  $P(\text{class})$
- Likelihood: **bag of words** representation
  - The document is a sequence of words  $(w_1, \dots, w_n)$
  - The order of the words in the document is not important
  - Each word is conditionally independent of the others given document class

$$P(\text{document} \mid \text{class}) = P(w_1, \dots, w_n \mid \text{class}) = \prod_{i=1}^n P(w_i \mid \text{class})$$

- Thus, the problem is reduced to estimating marginal likelihoods of individual words  $P(w_i \mid \text{class})$

# Parameter estimation

- Model parameters: feature likelihoods  $P(\text{word} \mid \text{class})$  and priors  $P(\text{class})$ 
  - How do we obtain the values of these parameters?

prior

spam:	0.33
¬spam:	0.67

$P(\text{word} \mid \text{spam})$

the :	0.0156
to :	0.0153
and :	0.0115
of :	0.0095
you :	0.0093
a :	0.0086
with:	0.0080
from:	0.0075
...	

$P(\text{word} \mid \neg\text{spam})$

the :	0.0210
to :	0.0133
of :	0.0119
2002:	0.0110
with:	0.0108
from:	0.0107
and :	0.0105
a :	0.0100
...	

# Outline: Bayesian Inference

- Bayes Rule
- Law of Total Probability
- Misdiagnosis
- $\text{MAP} = \text{MPE}$
- The “Naïve Bayesian” Assumption
- Bag of Words (BoW)
- Parameter Estimation for the BoW model



# Parameter estimation

- Model parameters: feature likelihoods  $P(\text{word} \mid \text{class})$  and priors  $P(\text{class})$ 
  - How do we obtain the values of these parameters?
  - Need *training set* of labeled samples from both classes

$$P(\text{word} \mid \text{class}) = \frac{\text{\# of occurrences of this word in docs from this class}}{\text{total \# of words in docs from this class}}$$

- This is the *maximum likelihood* (ML) estimate, or estimate that maximizes the likelihood of the training data:

$$\prod_{d=1}^D \prod_{i=1}^{n_d} P(w_{d,i} \mid \text{class}_{d,i})$$

$d$ : index of training document,  $i$ : index of a word

# Parameter estimation

- ML (Maximum Likelihood) parameter estimate:

$$P(\text{word} \mid \text{class}) = \frac{\text{\# of occurrences of this word in docs from this class}}{\text{total \# of words in docs from this class}}$$

- Laplacian Smoothing estimate

- How can you estimate the probability of a word you never saw in the training set? (Hint: what happens if you give it probability 0, then it actually occurs in a test document?)
- **Laplacian smoothing:** pretend you have seen every vocabulary word one more time than you actually did

$$P(\text{word} \mid \text{class}) = \frac{\text{\# of occurrences of this word in docs from this class} + 1}{\text{total \# of words in docs from this class} + V}$$

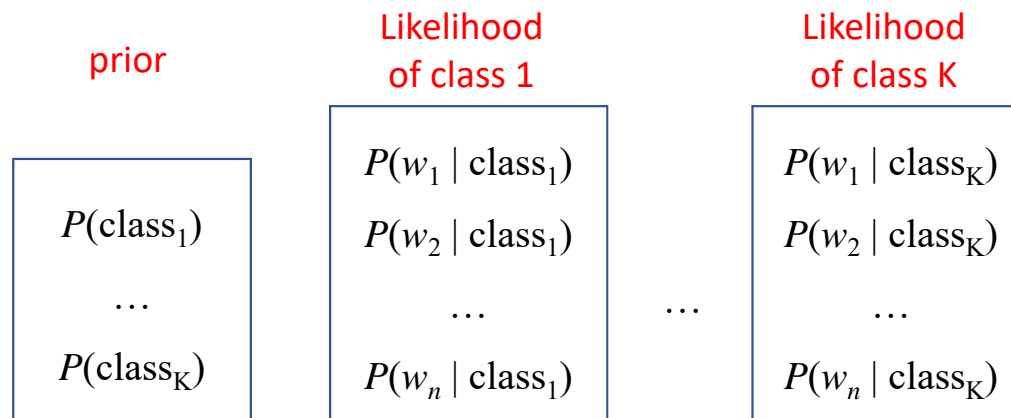
(V: total number of unique words)

# Summary: Naïve Bayes for Document Classification

- Naïve Bayes model: assign the document to the class with the highest posterior

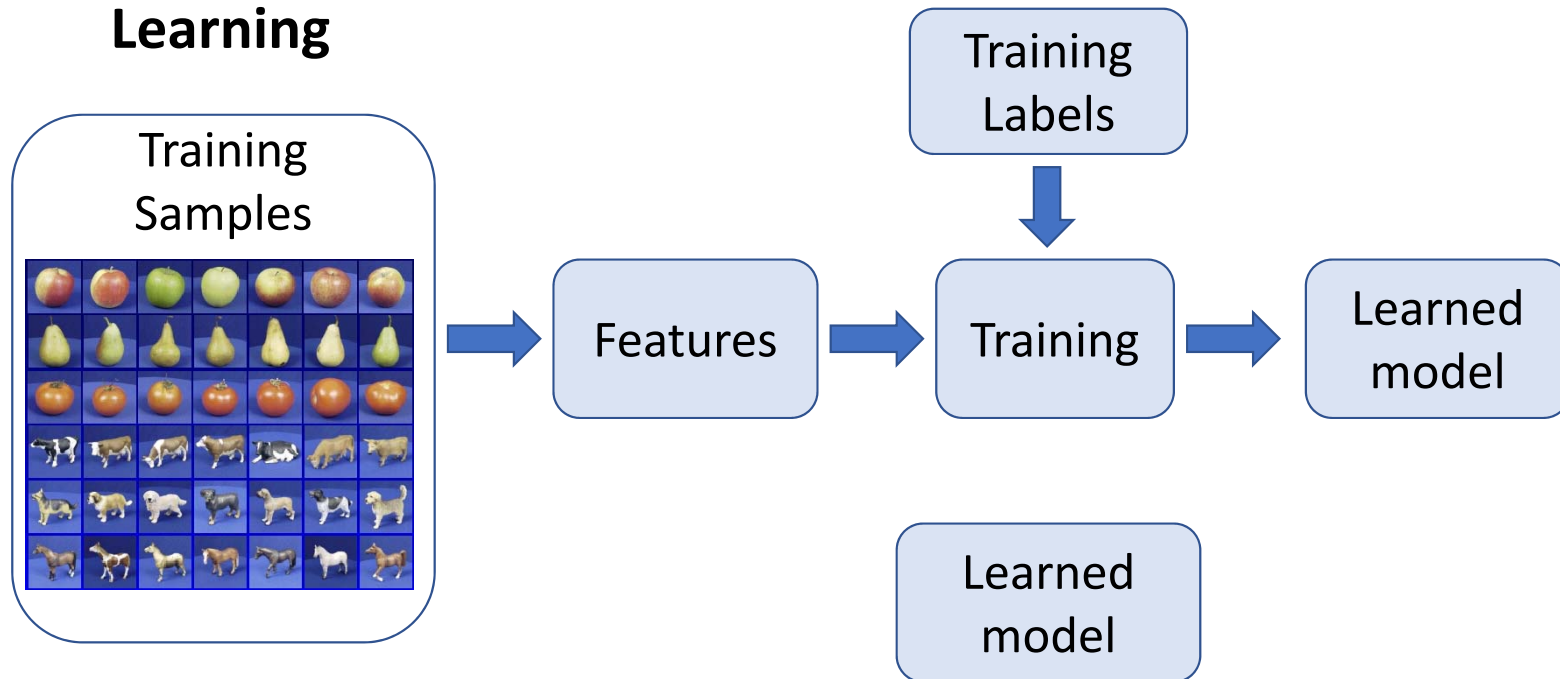
$$P(\text{class} \mid \text{document}) \propto P(\text{class}) \prod_{i=1}^n P(w_i \mid \text{class})$$

- Model parameters:

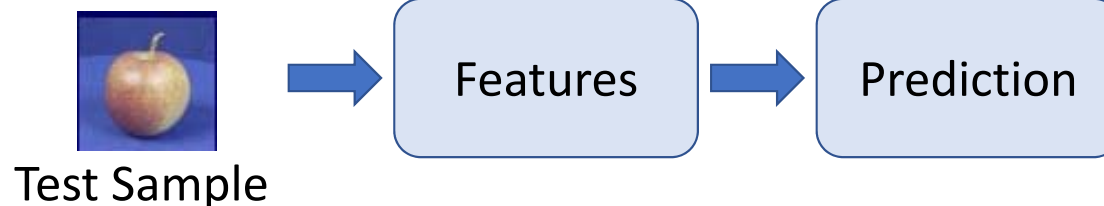


# Learning and inference pipeline

## Learning



## Inference



# Review: Bayesian decision making

- Suppose the agent has to make decisions about the value of an unobserved *query variable*  $X$  based on the values of an observed *evidence variable*  $E$
- **Inference problem:** given some observation  $E = e$ , what is  $P(X \mid e)$ ?
- **Learning problem:** estimate the parameters of the probabilistic model  $P(X \mid E)$  given a *training sample*  $\{(x_1, e_1), \dots, (x_n, e_n)\}$