

Applied Artificial Intelligence

Session 7: Probability and Statistics for AI and Machine Learning II

Fall 2018

NC State University

Instructor: Dr. Behnam Kia

Course Website: <https://appliedai.wordpress.ncsu.edu/>

Random Experiment



Sample Space: $\Omega = \{Head, Tail\}$

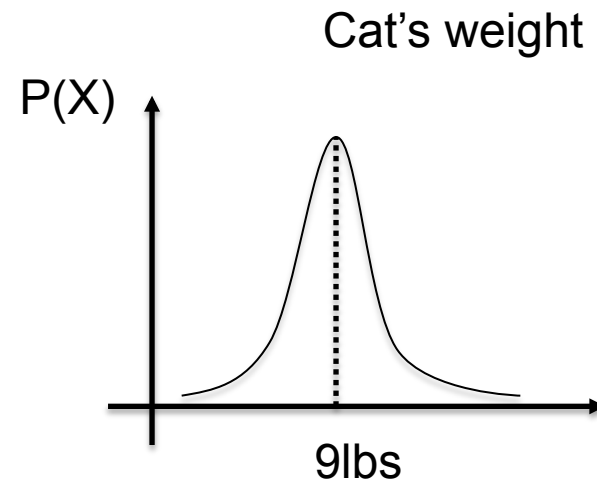
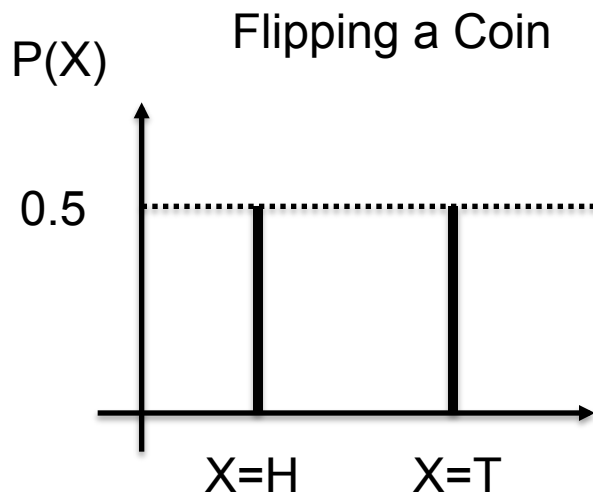
Random Variable

- A random variable is a variable that takes on different values based on outcomes of a random experiment.
- Random variable can be discrete or continuous.



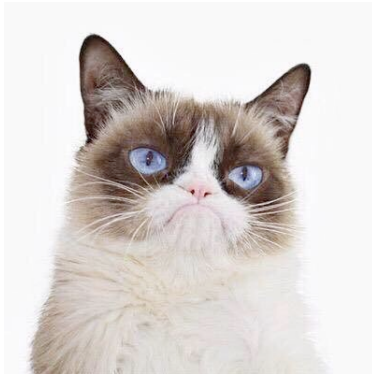
Probability Distribution Function

- Probability distribution function is a description of how likely a random variable or a set of variables is to take on each of its possible states.



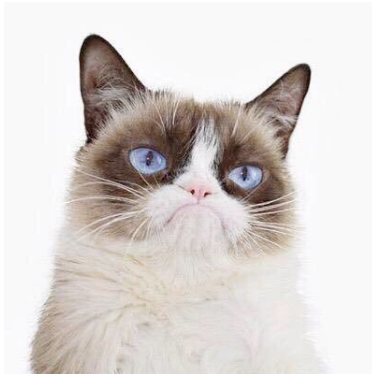
Hypothetical Scenario

- We are installing cameras and sensors in a neighborhood to record the presence of feral animals.
- There are two types of animals in the neighborhood; Dogs and Cats.
- We like to design an automatic system to determine whether the recorded animal is a cat or a dog.



Hypothetical Scenario

- Ratio of Cats to Dogs is 1 to 3. (75% are dogs, 25% cats)



Hypothetical Scenario

- Ratio of Cats to Dogs is 1 to 3. (75% are dogs, 25% cats)
- Making a decision without looking at the sensor readings.



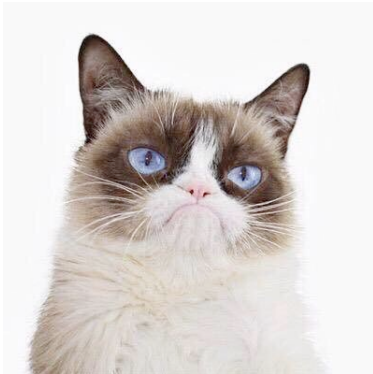
?



Hypothetical Scenario

- Ratio of Cats to Dogs is 1 to 3. (75% are dogs, 25% cats)
- Making a decision without looking at the sensor readings.

$$P(\text{Cat})=0.25$$



?

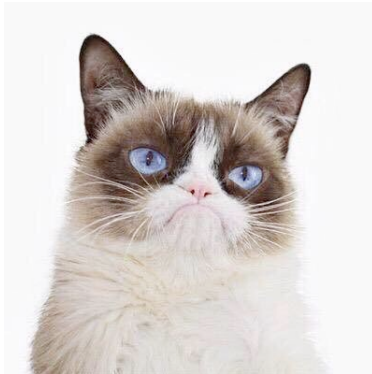
$$P(\text{dog})=0.75$$



Hypothetical Scenario

- Ratio of Cats to Dogs is 1 to 3. (75% are dogs, 25% cats)
- Making a decision without looking at the sensor readings.
- **What is the expected error rate?**

$$P(\text{Cat})=0.25$$



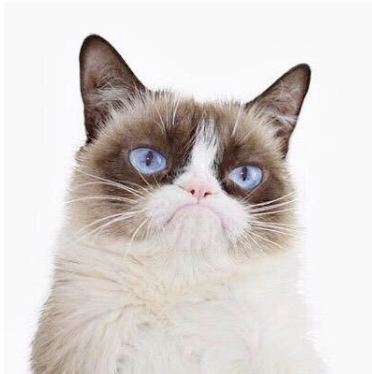
$$P(\text{dog})=0.75$$



Hypothetical Scenario

- Ratio of Cats to Dogs is 1 to 3. (75% are dogs, 25% cats)
- And we have access to some data, the height of the animal being recorded.

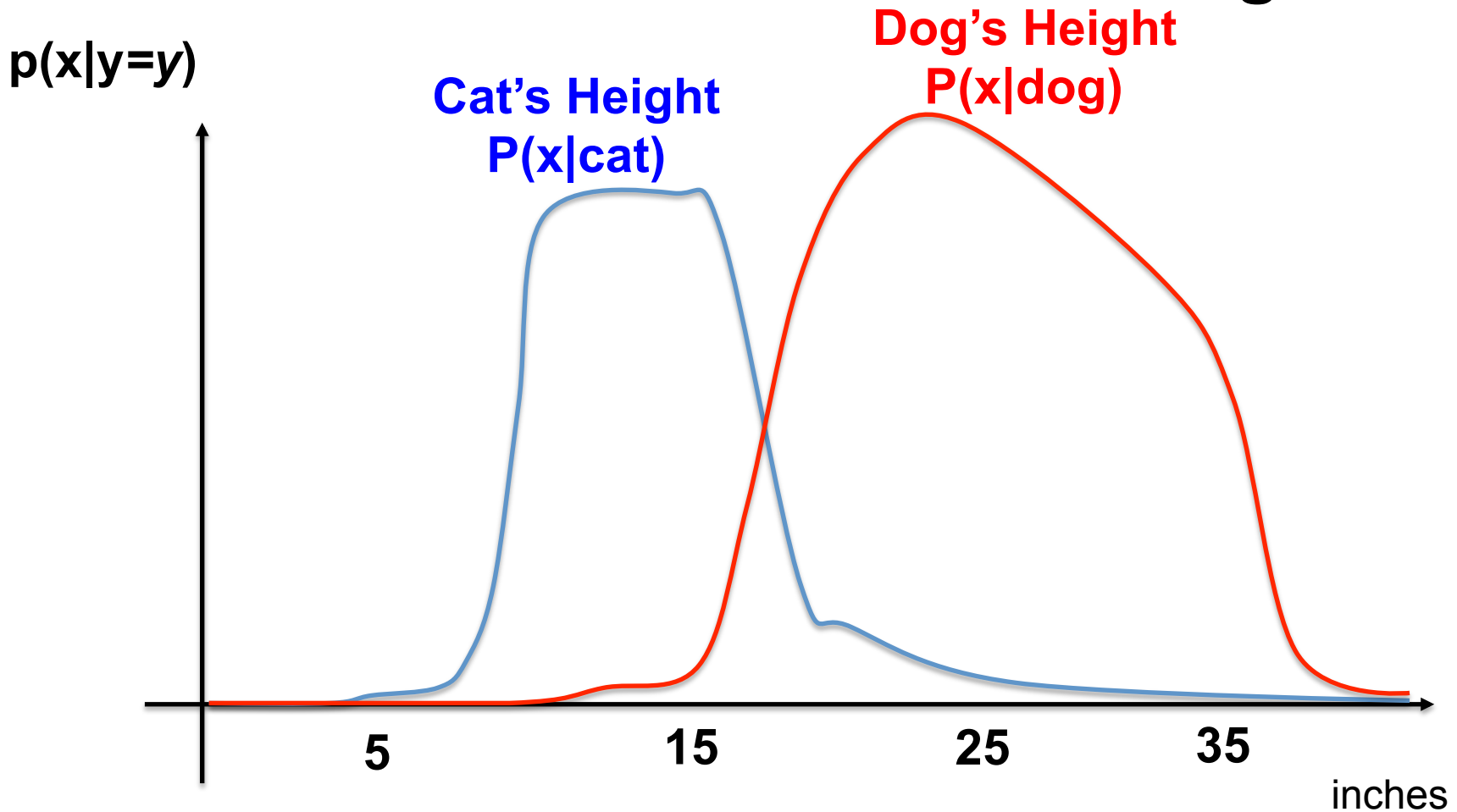
$$P(\text{Cat})=0.25$$



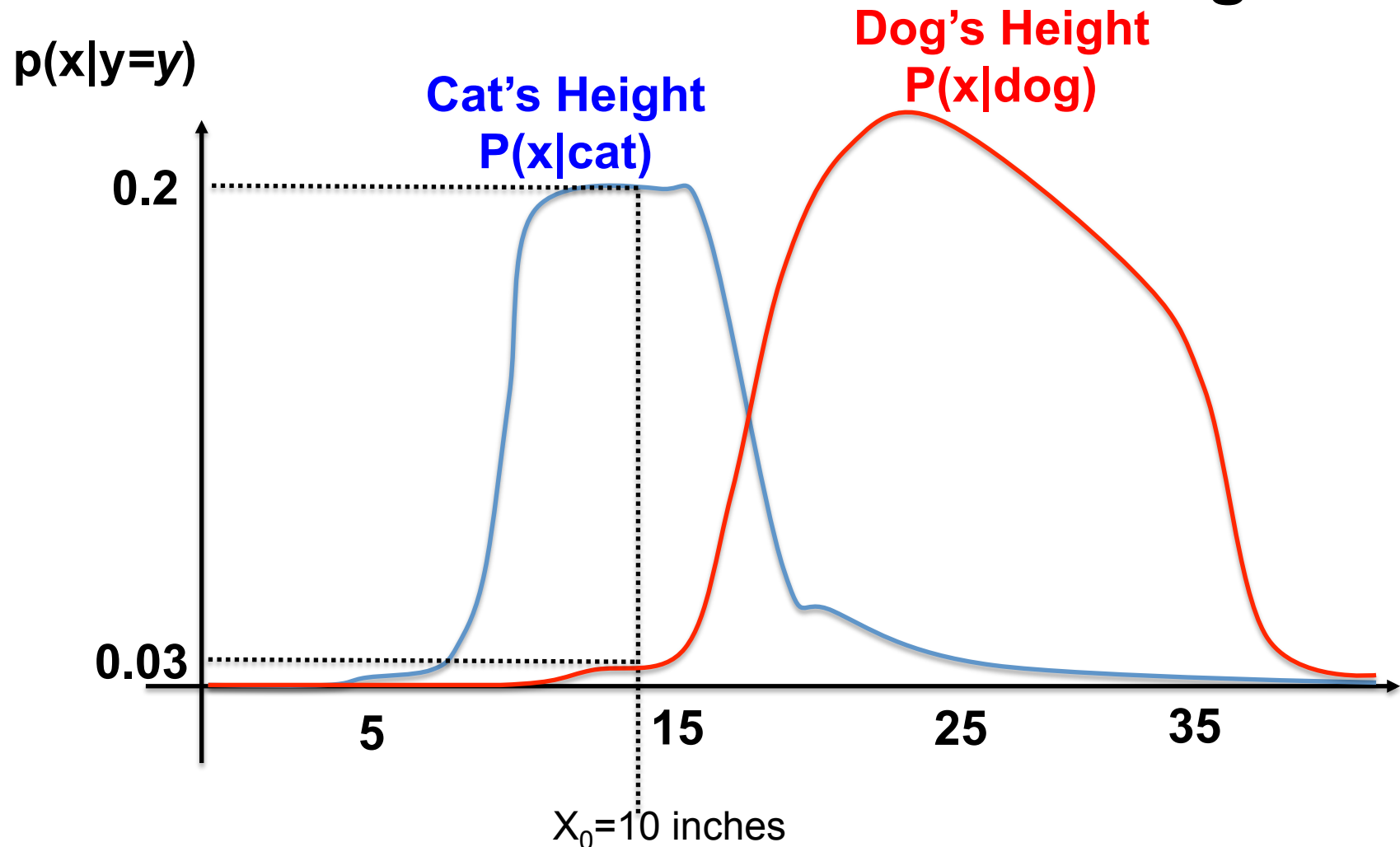
$$P(\text{dog})=0.75$$



Class-Conditional Probability Distribution Function for Height



Class-Conditional Probability Distribution Function for Height



Bayes' Rule

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

$$\textit{posterior} = \frac{\textit{likelihood} \times \textit{prior}}{\textit{evidence}}$$

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \quad \mathbf{\text{Bayes' Rule}} \quad p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

$$p(y = \text{cat} | x = 10) = \frac{p(x = 10 | y = \text{cat})p(\text{cat})}{p(x = 10)}$$



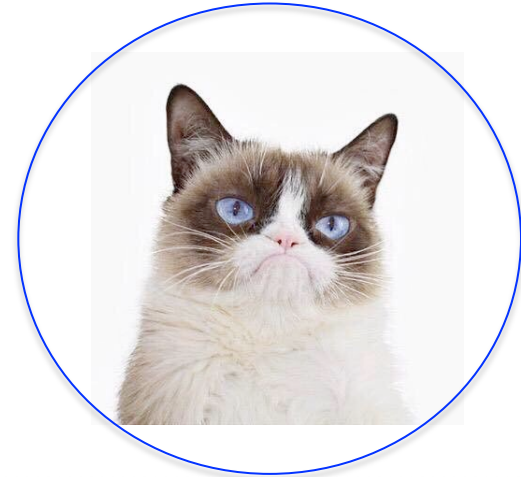
$$p(y = \text{dog} | x = 10) = \frac{p(x = 10 | y = \text{dog})p(\text{dog})}{p(x = 10)}$$



$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \quad \mathbf{\text{Bayes' Rule}} \quad p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

$$p(y = \text{cat} | x = 10) = \frac{p(x = 10 | y = \text{cat})p(\text{cat})}{p(x = 10)}$$

$$p(y = \text{cat} | x = 10) = \frac{0.2 \times 0.25}{p(x = 10)} = \frac{0.05}{p(x = 10)}$$



$$p(y = \text{dog} | x = 10) = \frac{p(x = 10 | y = \text{dog})p(\text{dog})}{p(x = 10)}$$

$$p(y = \text{dog} | x = 10) = \frac{0.03 \times 0.75}{p(x = 10)} = \frac{0.0225}{p(x = 10)}$$



$p(\text{evidence})$
Calculated Based on Law of Total Probability

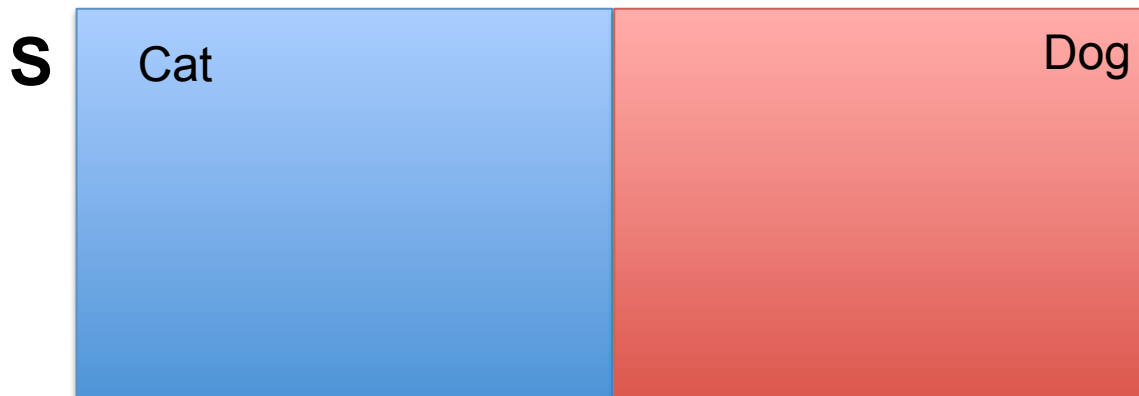
$$p(x = 10) = ?$$

S



$p(\text{evidence})$ Calculated Based on Law of Total Probability

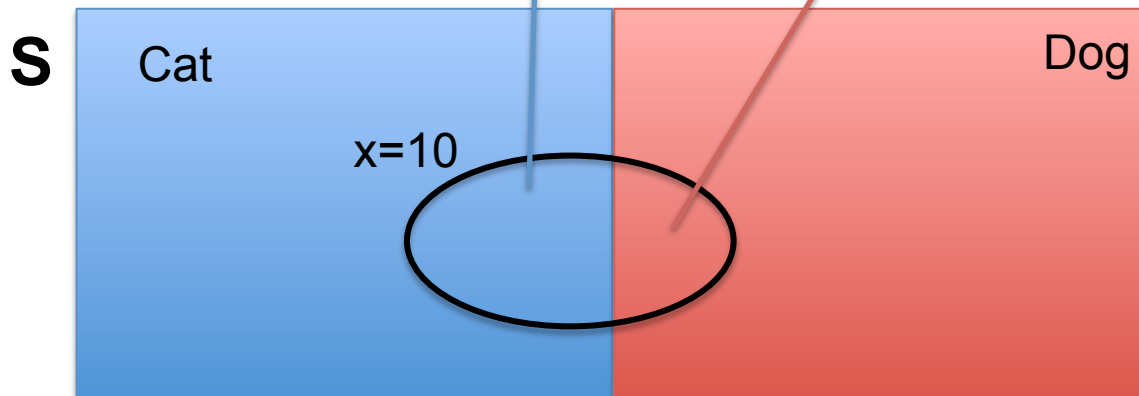
$$p(x = 10) = ?$$



$$S = \text{Cat} \cup \text{Dog}$$

$p(\text{evidence})$ Calculated Based on Law of Total Probability

$$p(x = 10) = p(x = 10 \cap y = \text{Cat}) + p(x = 10 \cap y = \text{Dog})$$



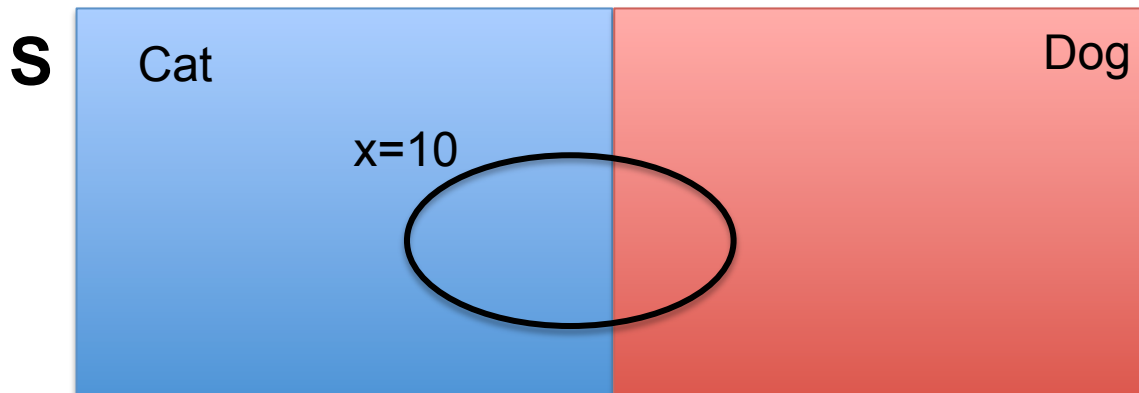
$p(\text{evidence})$ Calculated Based on Law of Total Probability

$$p(x = 10) = p(x = 10 \cap y = \text{Cat}) + p(x = 10 \cap y = \text{Dog})$$

$$p(x = 10) = p(x = 10 | y = \text{cat})p(y = \text{cat}) + p(x = 10 | y = \text{dog})p(y = \text{dog})$$

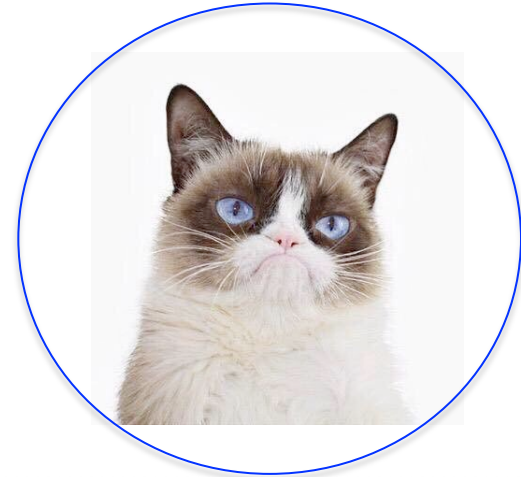
$$p(x = 10) = 0.2 \times 0.25 + 0.03 \times 0.75$$

$$= 0.0725$$



$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \quad \mathbf{\text{Bayes' Rule}} \quad p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

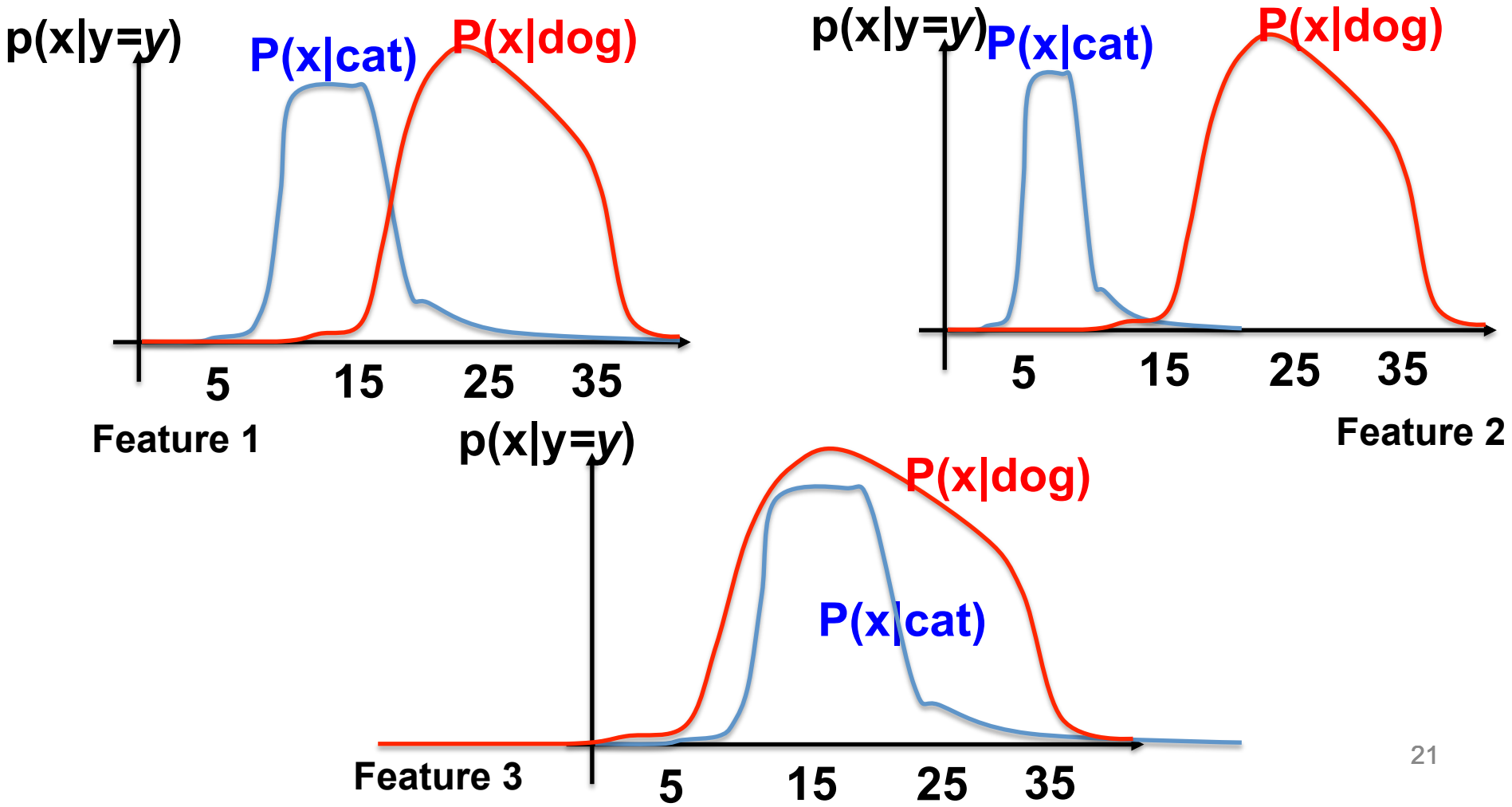
$$p(y = \text{cat} | x = 10) = \frac{0.2 \times 0.25}{p(x = 10)} = \frac{0.05}{p(x = 10)} = 0.69$$



$$p(y = \text{dog} | x = 10) = \frac{0.03 \times 0.75}{p(x = 10)} = \frac{0.0225}{p(x = 10)} = 0.31$$



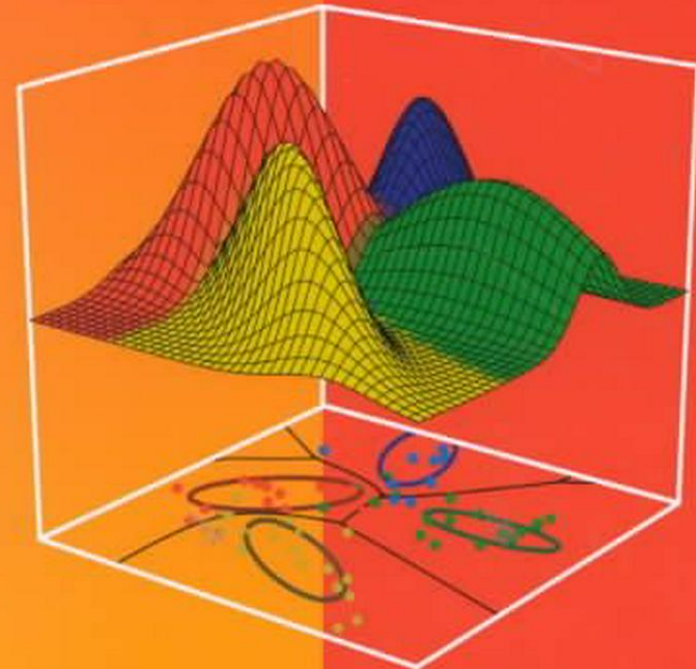
Which Feature x Would You Choose?



Same Concepts, But in 2D

Richard O. Duda
Peter E. Hart
David G. Stork

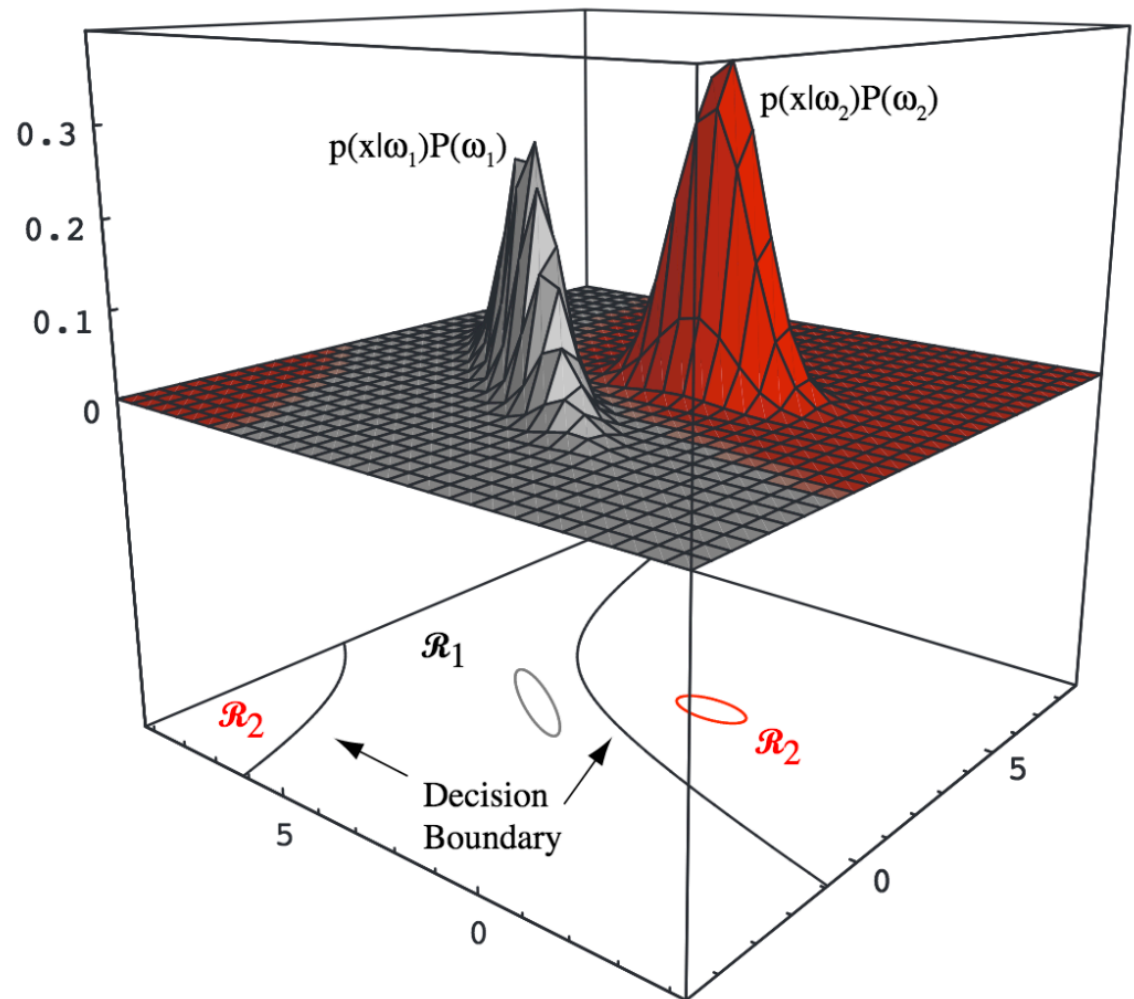
Pattern Classification



Pattern Classification,
Duda, Hart, and Stork

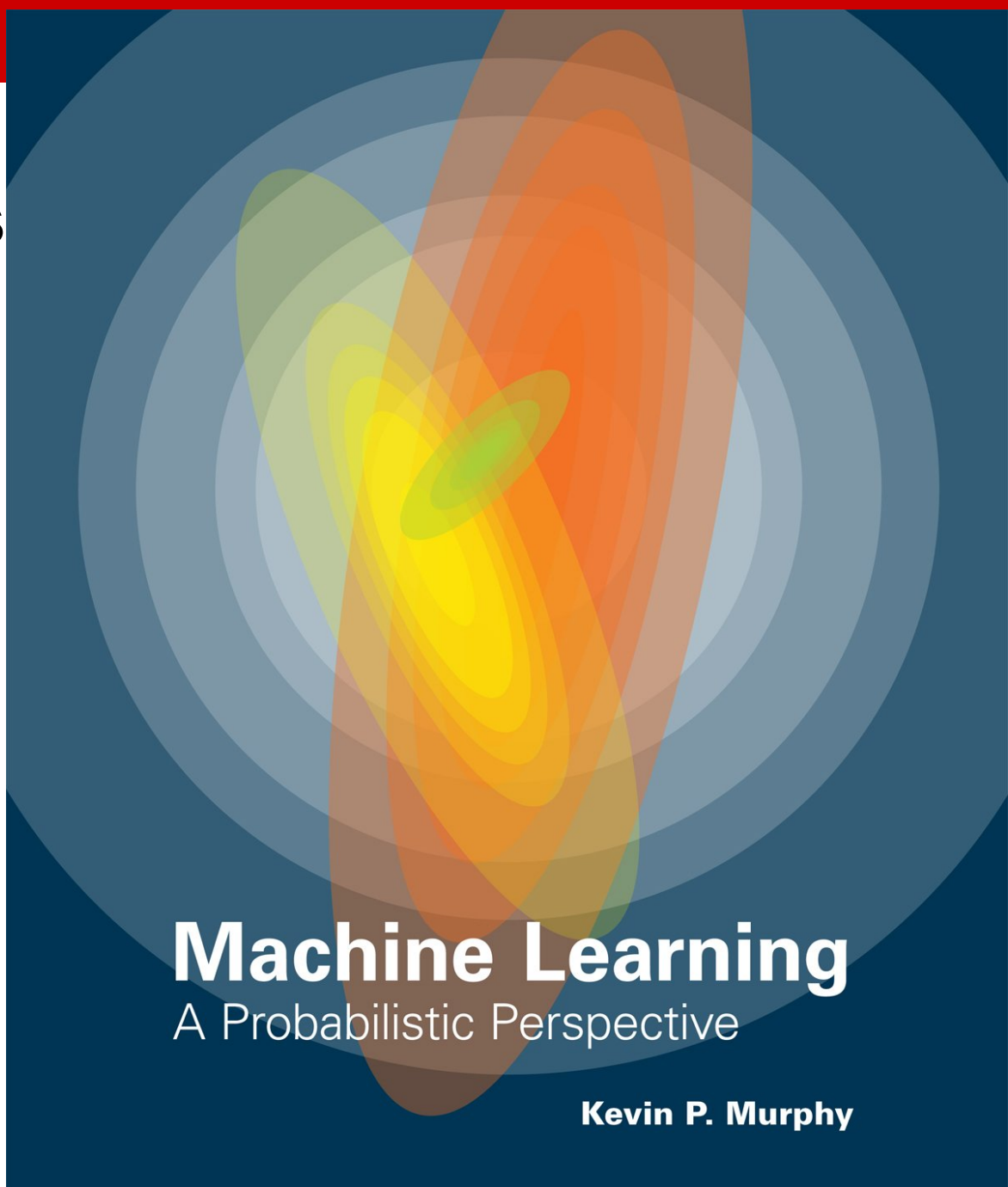
Second Edition

Same Concepts But in 2D



Same Concepts But in 2D

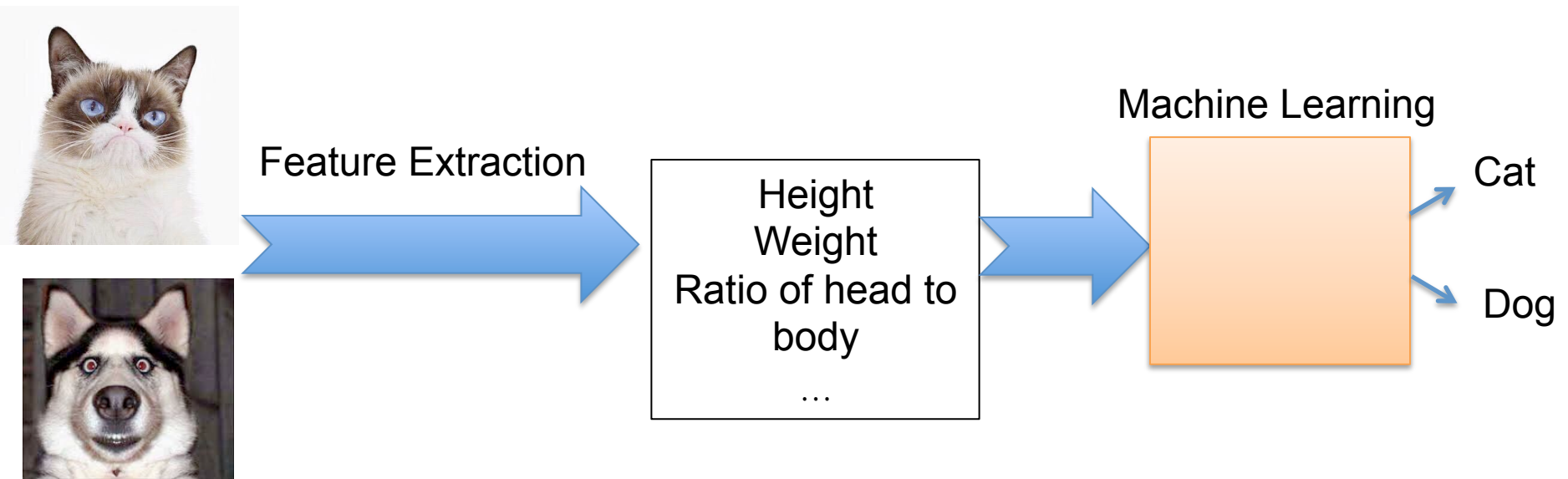
Machine Learning
A Probabilistic Perspective
Kevin Murphy



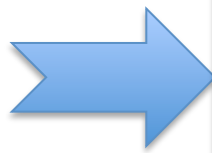
Machine Learning
A Probabilistic Perspective

Kevin P. Murphy

Classic Machine Learning

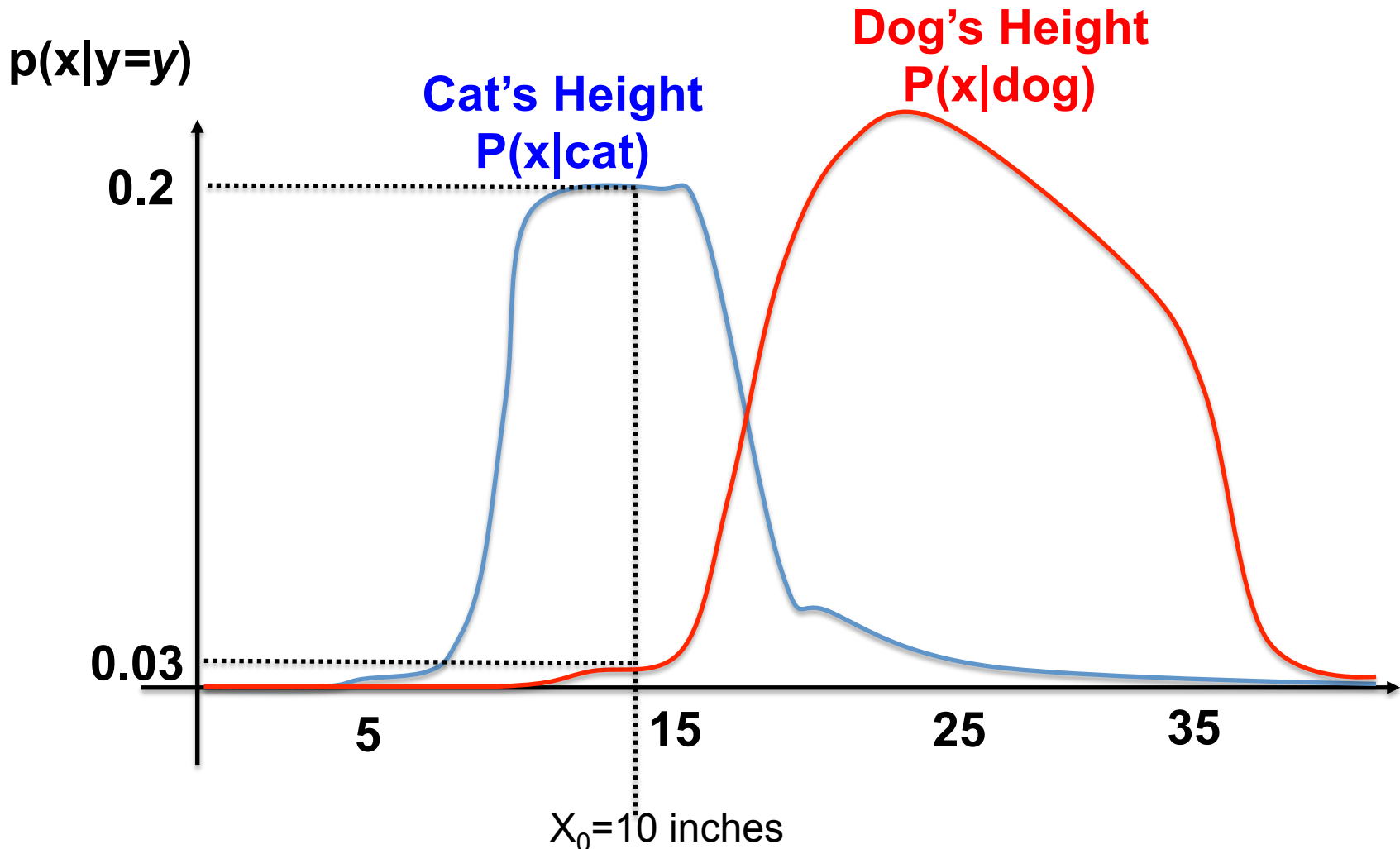


Deep Learning



Cat
Dog

How to Get Class-Conditional Probability Distribution Functions?



Machine Learning with Many Features

$$p(y | X) = \frac{p(X | y)p(y)}{p(X)}$$

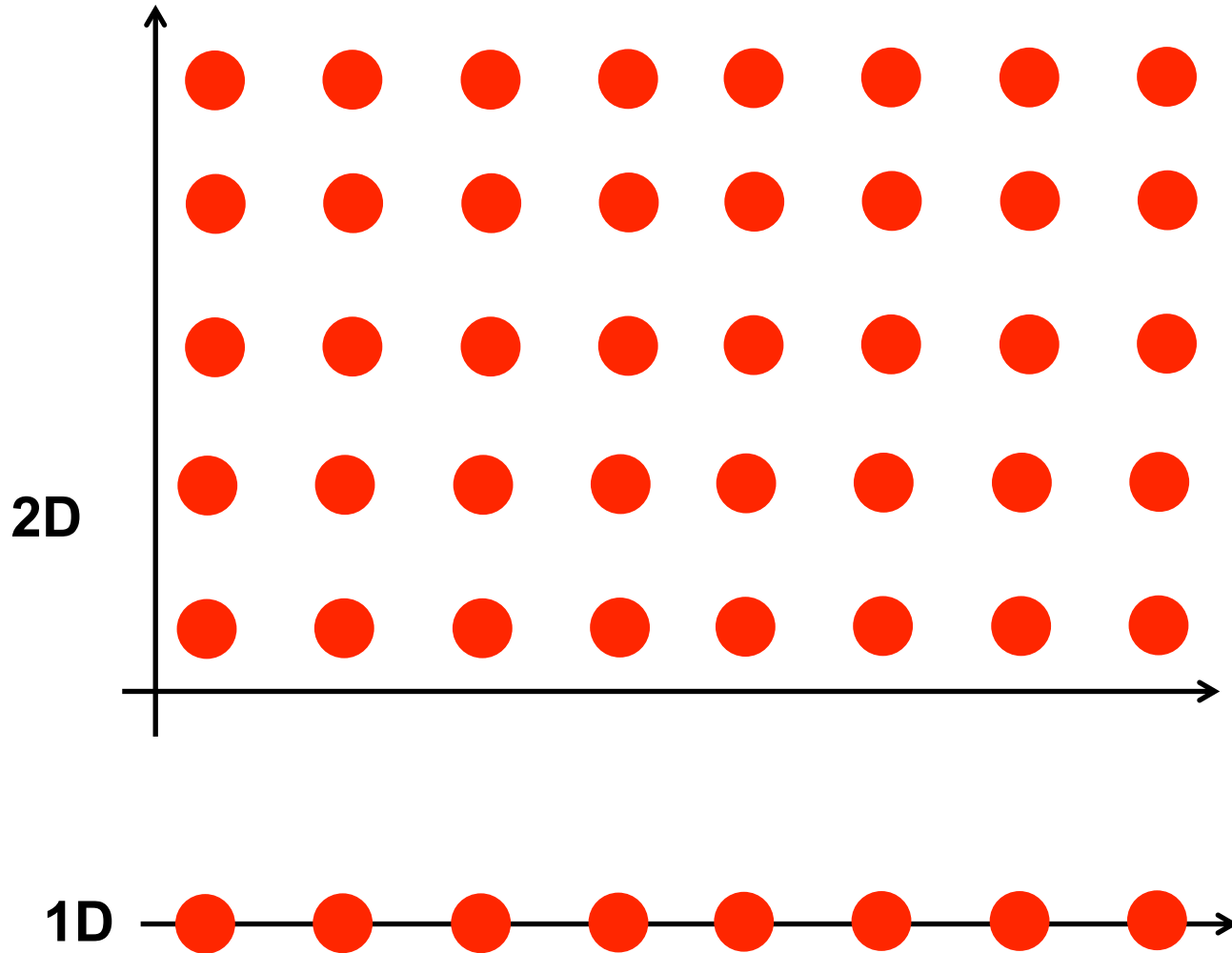
$$X = (x_1, x_2, \dots, x_k)$$

The Curse of Dimensionality

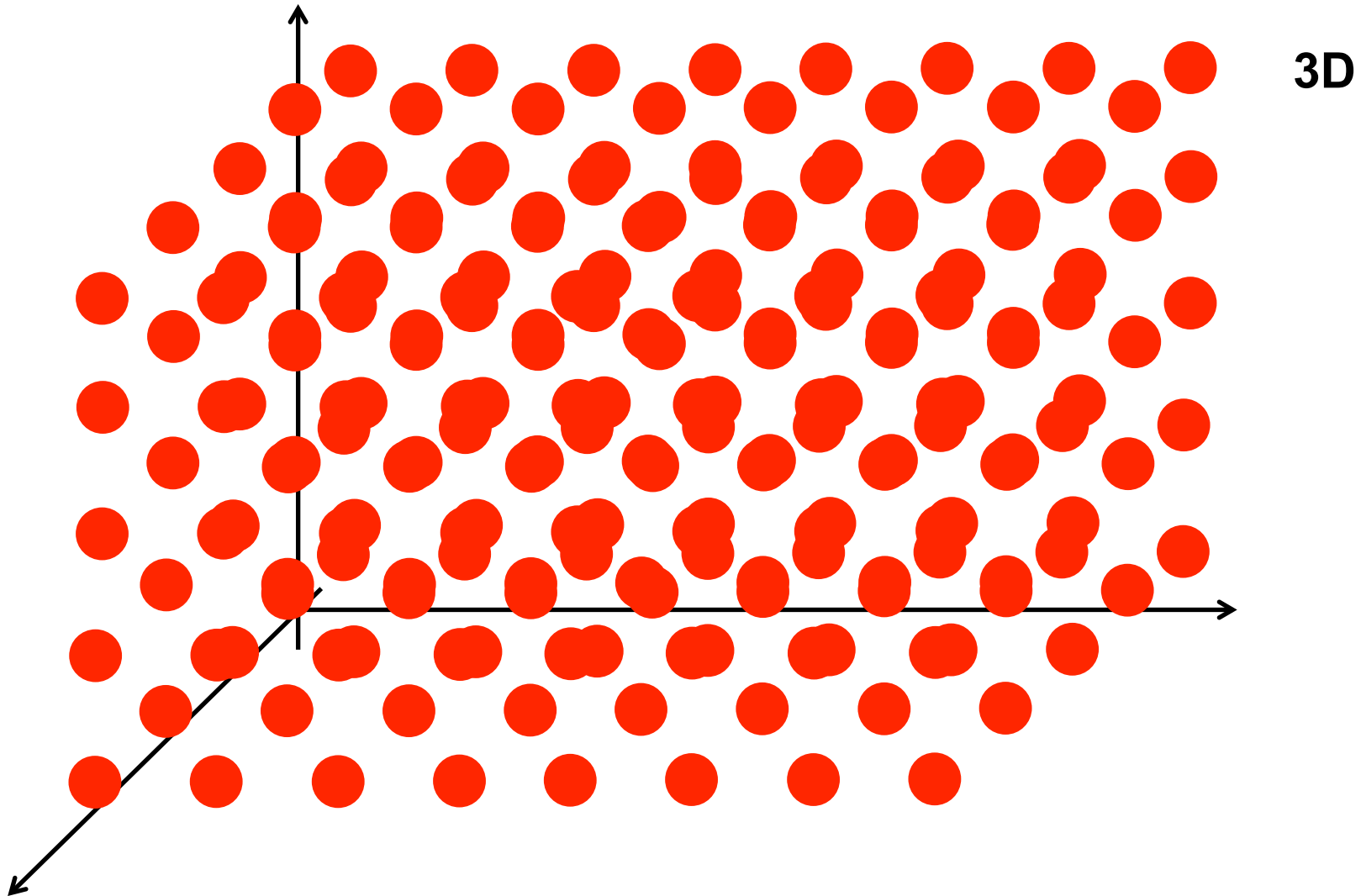
The Curse of Dimensionality



The Curse of Dimensionality



The Curse of Dimensionality



The Curse of Dimensionality

Naïve Bayes Rule

Bayes Rule

$$p(y | x_1, x_2, \dots, x_k) = \frac{p(x_1, x_2, \dots, x_k | y) p(y)}{p(X)}$$

$$\approx \frac{p(x_1 | y) p(x_2 | y) \dots p(x_k | y) p(y)}{p(X)}$$

Naïve Bayes Rule

We have assumed that features x_1, x_2, \dots , and x_k are conditionally independent given y

$$p(x_1, x_2, \dots, x_k | y) \approx p(x_1 | y) p(x_2 | y) \dots p(x_k | y)$$

Document Classification Using Naïve Bayes Rule

(Homework III)

Document Classification

- Imagine we are given a document. We would like to classify it. For example:
 - An email is spam or ham?
 - A review is positive or negative? (sentiment analysis)
 - The subject of the document is Math, Physics, or Chemistry?
 - Authorship identification
 - ...

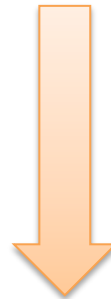
Document Classification

A review is positive or negative? (sentiment analysis)

Hahaha I knew it.

██████████ 13 September 2018

Straight trash. Why bother making a predator movie if you are just going to ignore all the stories? Heck they even ignored all we have learned from the movies. It seems anyone can defeat a predator nowadays. Back in the day, it took Arnie and a whole squad of macho man to try to take on just one. Now any ole joe smoe can go toe to toe with the predator. Stupid.



Classifier

Positive

?

Negative

Document Classification

A review is positive or negative? (sentiment analysis)

- Rule based approach:
 - If the review contains:
 - “What an awful movie” OR
 - “I need my money back!” OR
 - “I wish I had got sick so I couldn’t end up going to watch this movie!”
- Then it is Negative!

Document Classification

A review is positive or negative? (sentiment analysis)

- Rule based approach:
 - If the review contains:
 - “What a fun movie” OR
 - “I am going to watch it again!” OR
 - “This movie is the best thing that has happened to human race!”
- Then it is Positive!

Document Classification

A review is positive or negative? (sentiment analysis)

- Machine learning approach:
 - A training set of m labeled documents $(R_1, c_1), (R_2, c_2), \dots, (R_m, c_m)$ $c \downarrow 1, c \downarrow 2, \dots, c \downarrow m \in \{Positive, Negative\}$
 - Train a classifier that automatically assigns an unlabeled review to its correct class.
 - Many different machine learning techniques for this problem; here we use Naïve Bayes' Rule.

Document Representation

- Different representations.
- Here we use “*the bag of words*” representation. Order of the words doesn’t matter, just the words and how many times they occur in the text.

The Bag of Word Representation

- Review to be classified: “this was a good movie. This was the best movie of the series.”

The Bag of Words Representation of the Review

this	2
was	2
a	1
good	1
movie	2
the	1
best	1
of	1
the	1
series	1

Document Classification

c is the class, d is the document we like to classify

$$p(c | d) = \frac{p(d | c)p(c)}{p(d)}$$

Document Classification

c is the class, d is the document we like to classify

$$\begin{aligned}c_{MAP} &= \operatorname{argmax}_{c \in \{Positive, Negative\}} p(c | d) \\ &= \operatorname{argmax}_{c \in \{Positive, Negative\}} \frac{p(d | c)p(c)}{p(d)} \\ &= \operatorname{argmax}_{c \in \{Positive, Negative\}} p(d | c)p(c)\end{aligned}$$

Document Classification

c is the class, d is the document we like to classify

$$\begin{aligned}c_{MAP} &= \operatorname{argmax}_{c \in \{Positive, Negative\}} p(d | c) p(c) \\ &= \operatorname{argmax}_{c \in \{Positive, Negative\}} p(x_1, x_2, \dots, x_n | c) p(c)\end{aligned}$$

Document Classification

c is the class, d is the document we like to classify

$$\begin{aligned}
 c_{MAP} &= \operatorname{argmax}_{c \in \{Positive, Negative\}} p(d | c) p(c) \\
 &= \operatorname{argmax}_{c \in \{Positive, Negative\}} p(x_1, x_2, \dots, x_n | c) p(c) \\
 &\approx \operatorname{argmax}_{c \in \{Positive, Negative\}} p(x_1 | c) p(x_2 | c) \dots p(x_n | c) p(c)
 \end{aligned}$$

Conditional Independence Assumption (Naïve Bayes' Rule)

Document Classification

c is the class, d is the document we like to classify

$$\begin{aligned}
 c_{MAP} &= \operatorname{argmax}_{c \in \{Positive, Negative\}} p(d | c) p(c) \\
 &= \operatorname{argmax}_{c \in \{Positive, Negative\}} p(x_1, x_2, \dots, x_n | c) p(c) \\
 &\approx \operatorname{argmax}_{c \in \{Positive, Negative\}} p(x_1 | c) \times p(x_2 | c) \times \dots \times p(x_n | c) p(c)
 \end{aligned}$$

And we can easily learn every term on the right hand side from the training data.

Document Classification

Naïve Bayes

c is the class, d is the document we like to classify, x_i s are the words.

$$c_{NB} = \operatorname{argmax}_{c \in \{Positive, Negative\}} p(c) \prod_i p(x_i | c)$$

And we can easily learn every term on the right hand side from the training data.

Document Classification, Naïve Bayes: Learning from Training Data

$$c_{NB} = \operatorname{argmax}_{c \in \{Positive, Negative\}} p(c) \prod_i p(x_i | c)$$

$N_{reviews}$: number of reviews
In training set

$$\hat{p}(c_p) = \frac{\text{reviewcount}(Positive)}{N_{reviews}}$$

reviewcount(Positive): number of
positive reviews.

$$\hat{p}(c_N) = \frac{\text{reviewcount}(Negative)}{N_{reviews}}$$

reviewcount(Negative): number of
negative reviews.

Document Classification, Naïve Bayes: Learning from Training Data

$$\hat{p}(x_i | c_p) = \frac{\text{count}(x_i, \text{positive reviews})}{\sum_{x \in V} \text{count}(x, \text{positive reviews})}$$

Fraction of times word x_i shows up among all words in the positive reviews.

$\text{count}(x_i, \text{positive reviews})$: how many times the word x_i has appeared in positive reviews. (you have to repeat this process for every word x_i).

V : The vocabulary. All the words that show up in training reviews.

Document Classification, Naïve Bayes: Learning from Training Data

$$\hat{p}(x_i | c_N) = \frac{\text{count}(x_i, \text{negative reviews})}{\sum_{x \in V} \text{count}(x, \text{negative reviews})}$$

Fraction of times word x_i shows up among all words in the negative reviews.

count(x_i , negative reviews): how many times the word x_i has appeared in negative reviews. (you have to repeat this process for every word x_i).

V : The vocabulary. All the words that show up in reviews.

Document Classification, Naïve Bayes: Learning from Training Data

Two numerical Challenge:

1- What if a word in a test review has never showed up in positive or negative training reviews?

$$\hat{p}(x_i | c_N) = \frac{\text{count}(x_i, \text{negative reviews})}{\sum_{x \in V} \text{count}(x, \text{negative reviews})} = 0$$

$$c_{NB} = \operatorname{argmax}_{c \in \{\text{Positive}, \text{Negative}\}} p(c) \prod_i p(x_i | c)$$

Document Classification, Naïve Bayes: Learning from Training Data

Two numerical Challenge:

1- What if a word in a test review has never showed up in positive or negative training reviews?

$$\hat{p}(x_i | c_N) = \frac{\text{count}(x_i, \text{negative reviews}) + 1}{\sum_{x \in V} (\text{count}(x, \text{negative reviews}) + 1)}$$

$$\hat{p}(x_i | c_p) = \frac{\text{count}(x_i, \text{positive reviews}) + 1}{\sum_{x \in V} (\text{count}(x, \text{positive reviews}) + 1)}$$

Document Classification, Naïve Bayes: Learning from Training Data

Two numerical Challenge:

- 1- What if a word in a test review has never showed up in positive or negative training reviews?
- 2- We are multiplying many small numbers together. How to fight against underflow?

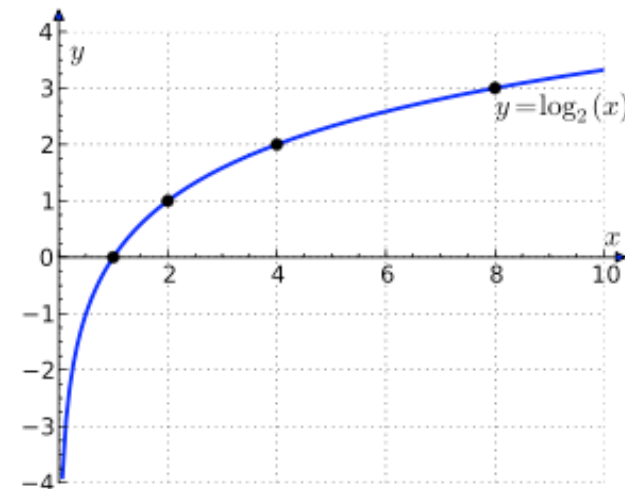
$$c_{NB} = \operatorname{argmax}_{c \in \{Positive, Negative\}} p(c) \prod_i p(x_i | c)$$

Document Classification, Naïve Bayes: Learning from Training Data

Two numerical Challenge:

- 1- What if a word in a test review has never showed up in positive or negative training reviews?
- 2- We are multiplying many small number fight against underflow?

$$c_{NB} = \operatorname{argmax}_{c \in \{Positive, Negative\}} \log(p(c) \prod_i p(x_i | c))$$



Document Classification, Naïve Bayes: Learning from Training Data

Two numerical Challenge:

- 1- What if a word in a test review has never showed up in positive or negative training reviews?
- 2- We are multiplying many small numbers together. How to fight against underflow?

$$\begin{aligned}
 c_{NB} &= \operatorname{argmax}_{c \in \{Positive, Negative\}} \log(p(c) \prod_i p(x_i | c)) \\
 &= \operatorname{argmax}_{c \in \{Positive, Negative\}} [\log p(c) + \sum_i \log(p(x_i | c))]
 \end{aligned}$$

Reading Assignment:

Shimodaira, Hiroshi. "Text classification using naive bayes."
Learning and Data Note 7 (2014): 1-9.

<https://web.stanford.edu/class/cs124/lec/naivebayes.pdf>