



Machine Learning and Data Mining

Lecture 3.2: Generative Models and Naïve Bayes





RECAP

From Last Lecture



Where are we?

We've covered a few techniques for fitting parameter:

Parameter Estimation:

Frameworks

- Maximum Likelihood Estimation
- Maximum A Posteriori

Optimization Tools

- Gradient Descent



We've covered a few predictive models so far in class:

Predictive Models:

Classifiers (given x , produce discrete y)

- k-Nearest Neighbors Classifier
- Logistic Regression (binary classification only)
- Perceptron (binary classification only)
- Naïve Bayes Classifier (Today)

Regressors (given x , produce continuous y)

- Linear Regression
- k-Nearest Neighbors Regressor

Today's Learning Objectives

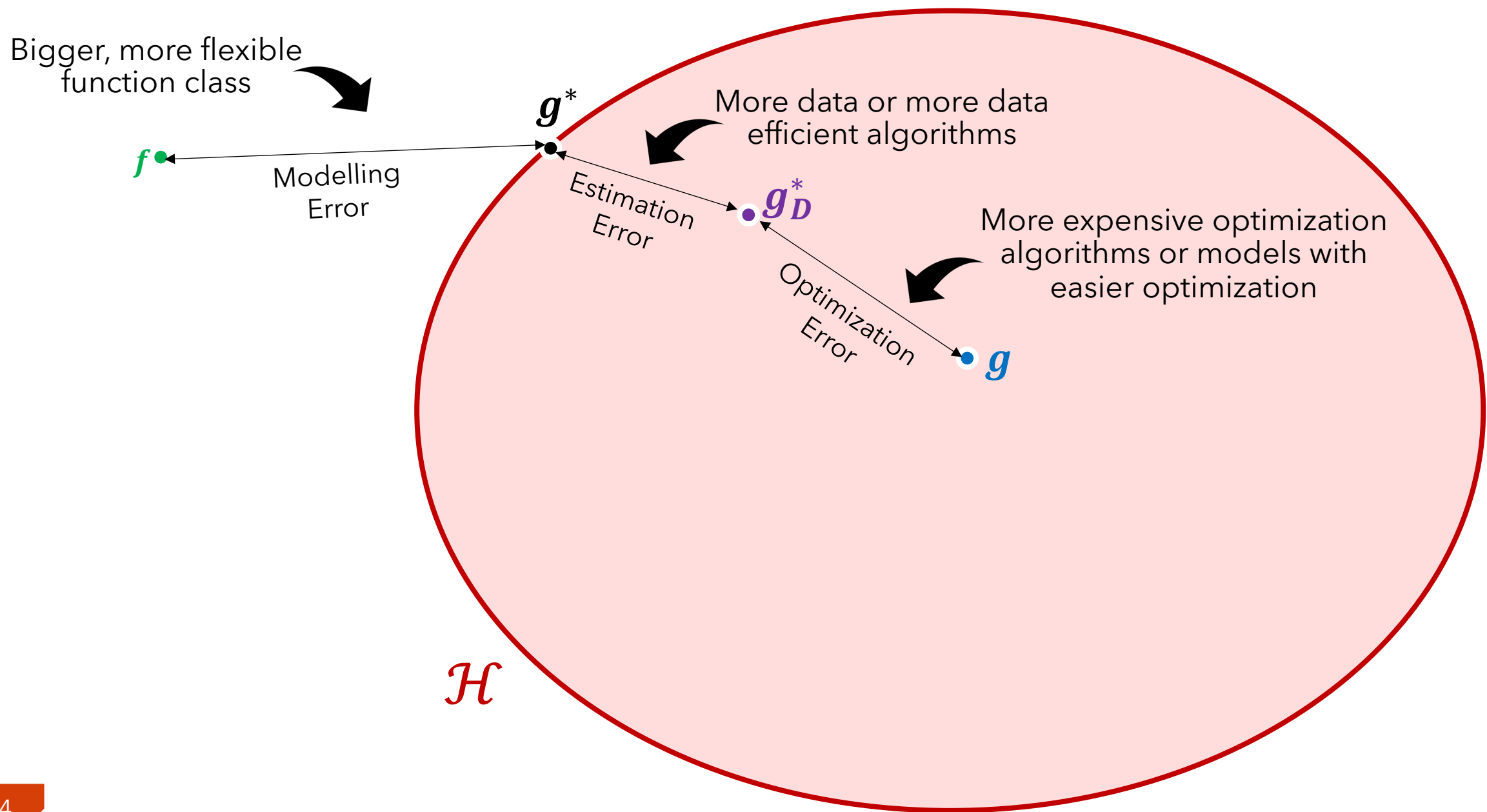
Be able to answer:

- What is a Bayes classifier?
- How do generative and discriminative models differ?
- What is a generative “story” for data?
 - How do we estimate a joint distribution for discrete variables?
- What is the Naïve-Bayes classifier?
 - What assumption makes it “naïve”?
 - What is conditional independence in probability?
 - How do you fit conditional distributions
 - What is Bernoulli Naïve Bayes?
 - What is Categorical Naïve Bayes?
 - What is Gaussian Naïve Bayes?





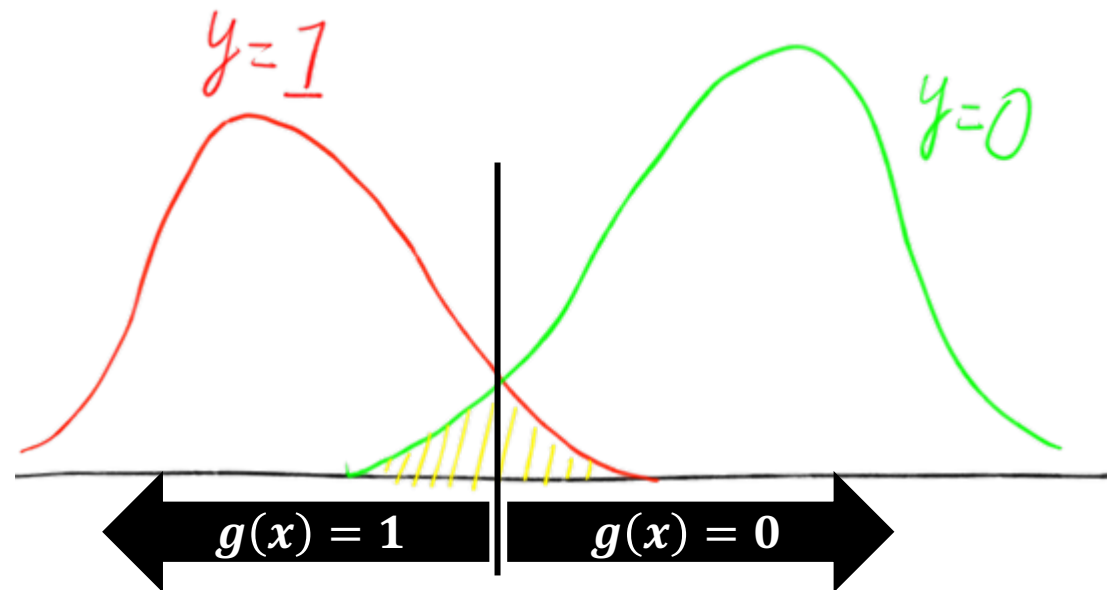
Recall: Sources of Error





Bayes Error

Irreducible error inherit in the function being approximated - nothing we can fix.





Bayes Error

Irreducible error inherit in the function being approximated - nothing we can fix.

Likes Math	Reads ML News	Thinks MLE Is Cool?
0	0	0
1	0	1
1	1	1
1	1	0





What is a Bayes Classifier?

Optimal Bayes Classifier:

Suppose we know the true distribution $P^*(Y|X)$ and for each x we encounter we predict:

$$\hat{y} = \operatorname{argmax}_y P^*(Y = y|X = x)$$

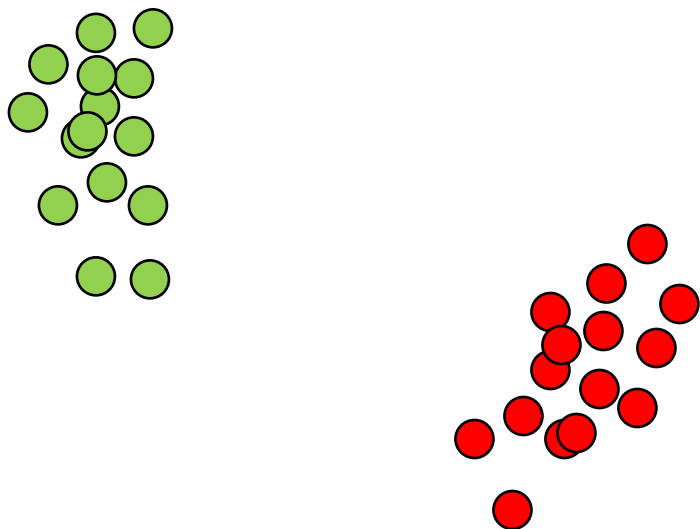
- If we know the true $P^*(Y|X)$, this is optimal.
https://en.wikipedia.org/wiki/Bayes_classifier

Problem: We don't know the true $P(Y|X)$. How to learn it?



Problem: We don't know the true $P(y|x)$. How to learn it?

Consider the following binary classification problem (colors for labels):

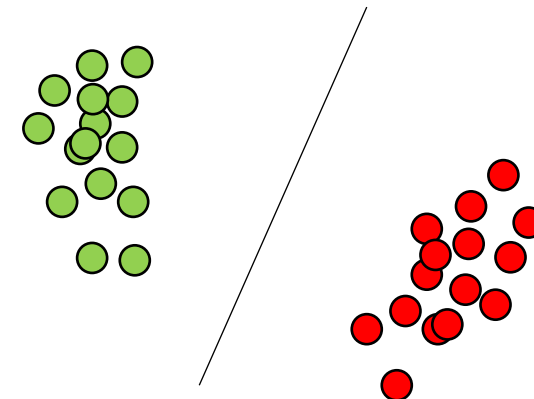


In logistic regression, we directly tried to model $P(y|x)$ -- assuming it was $\sigma(w^T x)$ and learning w with MLE.

An alternative approach would be to model each cluster separately – i.e. modelling $P(x|y)$.

Discriminative Classifiers:

- Learn $P(y|\mathbf{x})$ directly
- Logistic regression is one example
- *Nomenclature note -- people will also refer to algorithms that model no distribution as discriminative (such as kNN).*

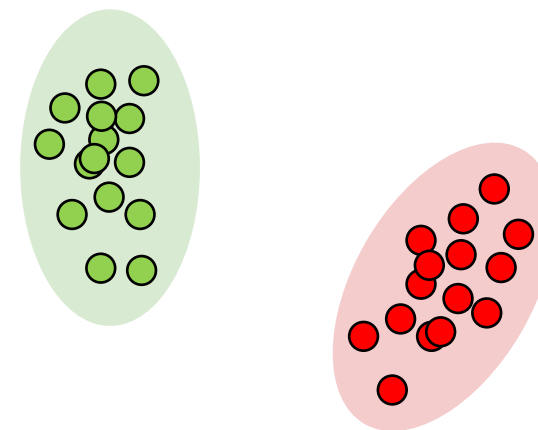


Generative Classifiers:

- Learn $P(\mathbf{x}|y)$ and $P(y)$
- Compute $P(y|\mathbf{x})$ using Bayes Rule

$$P(y|\mathbf{x}) = \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} = \frac{P(\mathbf{x}|y)P(y)}{\sum_y P(\mathbf{x}|y)P(y)}$$

- Naïve Bayes is one example (today)



Both classify according to $\operatorname{argmax} P(y|\mathbf{x})$. Just learn and represent it differently.



Generative Classifiers:

- Learn $P(\mathbf{x}|y)$ and $P(y)$
- Compute $P(y|\mathbf{x})$ using Bayes Rule

$$P(y|\mathbf{x}) = \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} = \frac{P(\mathbf{x}|y)P(y)}{\sum_y P(\mathbf{x}|y)P(y)}$$

When defining a generative classifier, we are rather explicitly making a *fictional* story about how the dataset we observe was generated – the reading calls these *generative stories*.

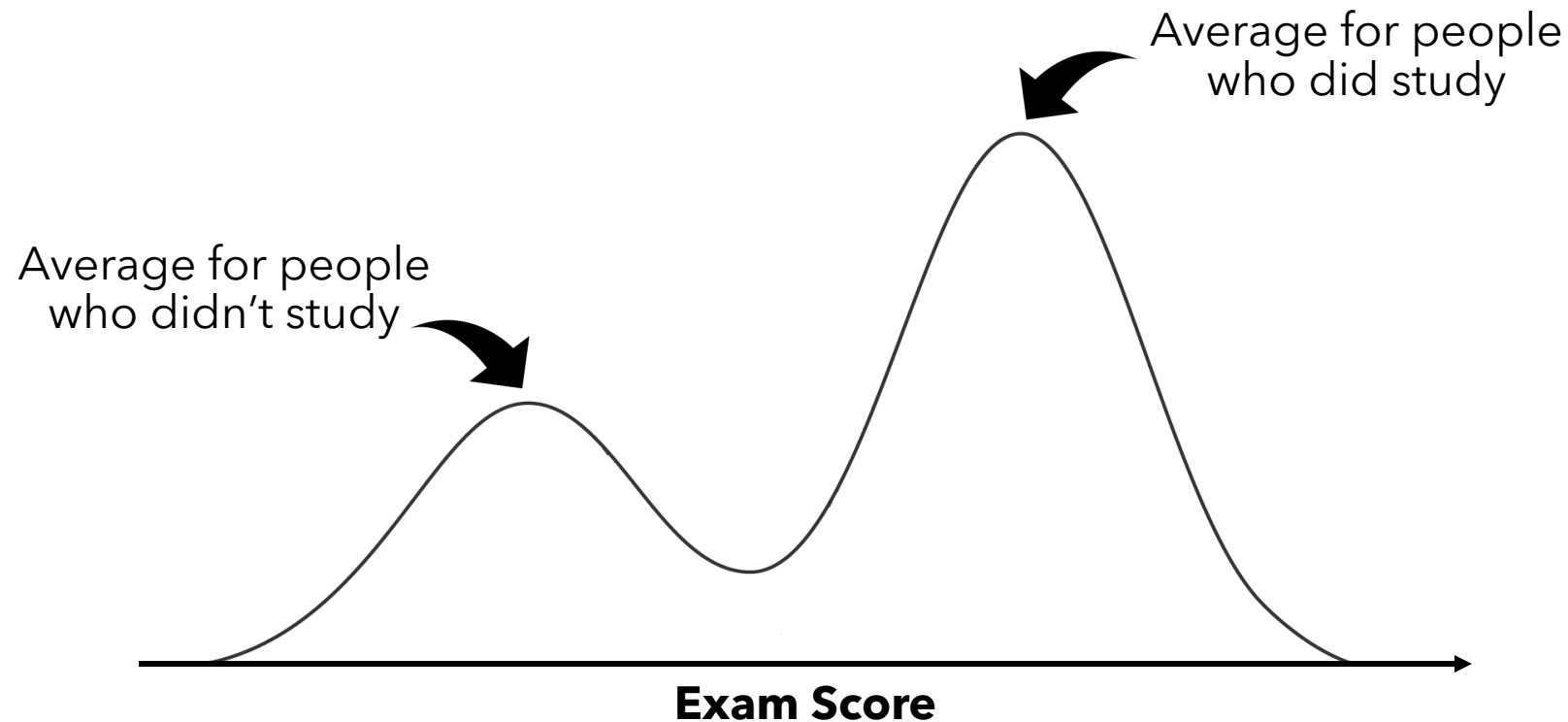
Our modeling decisions combined with the data answer:

$P(y)$ How are classes distributed?

$P(\mathbf{x}|y)$ Given a class, how are features distributed?



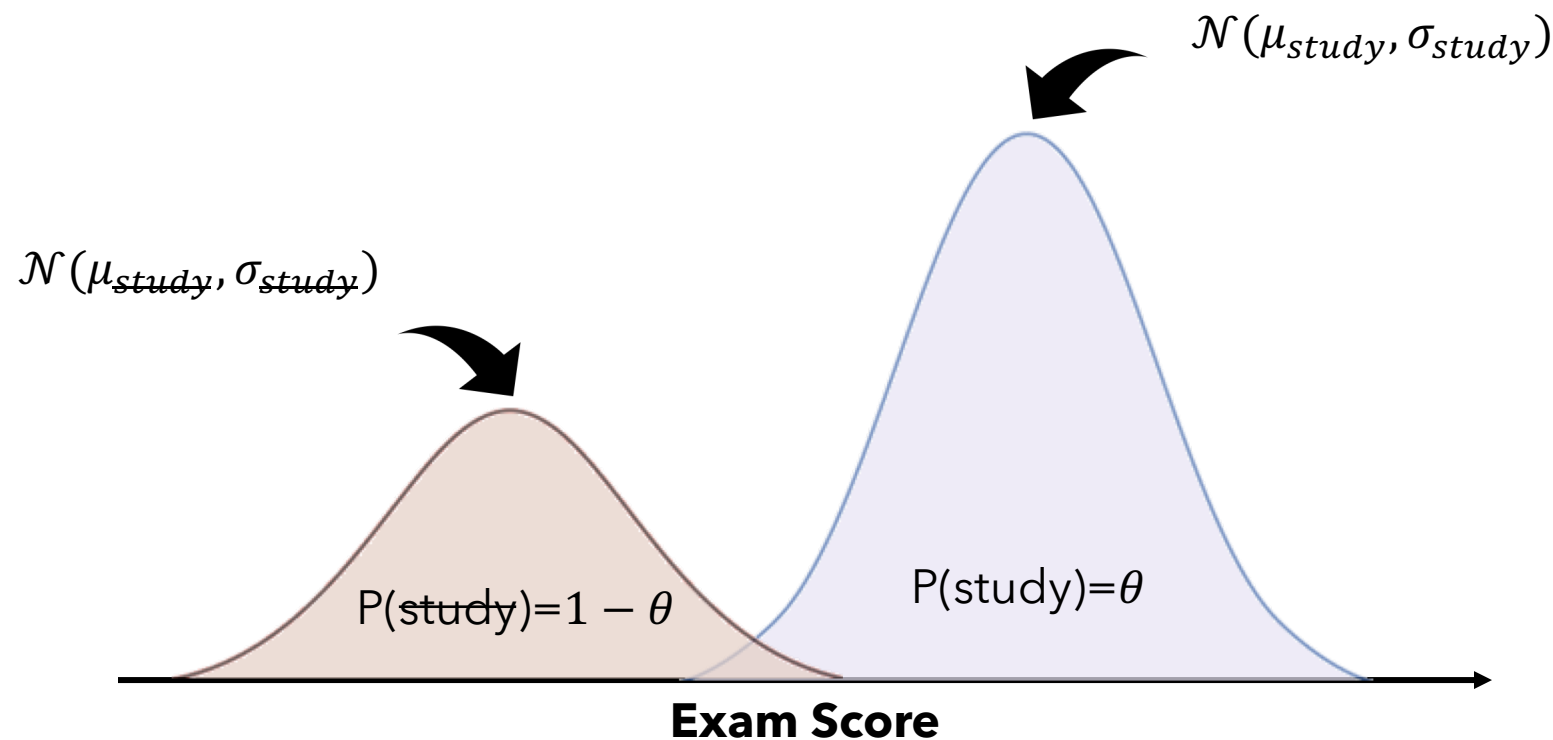
For example, suppose we have a class of students. Each student either studies or doesn't study before taking an exam. Consider the distribution of exam scores below. How might this distribution come to be?





A simple generative story for this exam scores:

1. Flip a coin ($P(\text{heads})=\theta$) to decide if someone studies (heads=study, tails=no)
2. If **study**, sample the exam score from a Gaussian distribution $\mathcal{N}(\mu_{\text{study}}, \sigma_{\text{study}})$
3. If **no study**, sample the exam score from a Gaussian distribution $\mathcal{N}(\mu_{\text{no study}}, \sigma_{\text{no study}})$





How to fit these distributions?

- $P(\text{study}) = \text{Bernoulli}(\theta_{\text{study}})$
 - $\theta_{\text{study}} = \# \text{ studying students} / \text{total students}$
- $P(\text{score} \mid \text{study}) = \mathcal{N}(\text{score}; \mu_{\text{study}}, \sigma_{\text{study}})$
 - $\mu_{\text{study}}, \sigma_{\text{study}} = \text{average and variance of scores of studying students}$
- $P(\text{score} \mid \text{not study}) = \mathcal{N}(\text{score}; \mu_{\text{not study}}, \sigma_{\text{not study}})$
 - $\mu_{\text{not study}}, \sigma_{\text{not study}} = \text{average and variance of scores of non-studying}$



Given an observed exam score, how would this story help us predict whether the student studied? Bayes rule to the rescue again.

$$\begin{aligned}P(\textit{study}|\textit{score}) &\propto P(\textit{score}|\textit{study})P(\textit{study}) \\ &= \mathcal{N}(\textit{score}; \mu_{\textit{study}}, \sigma_{\textit{study}})\theta\end{aligned}$$

$$\begin{aligned}P(\textit{no study}|\textit{score}) &\propto P(\textit{score}|\textit{no study})P(\textit{no study}) \\ &= \mathcal{N}(\textit{score}; \mu_{\textit{study}}, \sigma_{\textit{study}})(1 - \theta)\end{aligned}$$

If $P(\textit{study}|\textit{score}) > P(\textit{no study}|\textit{score})$, then predict study. Otherwise predict no study.



Study?	Score
Yes	87
Yes	89
Yes	87
Yes	88
Yes	91
Yes	97
Yes	91
Yes	93
No	80
No	77
No	64
No	75
No	91

$$P(study) = \frac{8}{14} \qquad P(no\ study) = \frac{6}{14}$$

$$P(score|study) = \mathcal{N}(score; \mu_{study} = 90.375, \sigma_{study} = 3.19)$$

$$P(score|no\ study) = \mathcal{N}(score; \mu_{no\ study} = 77.4, \sigma_{no\ study} = 8.68)$$

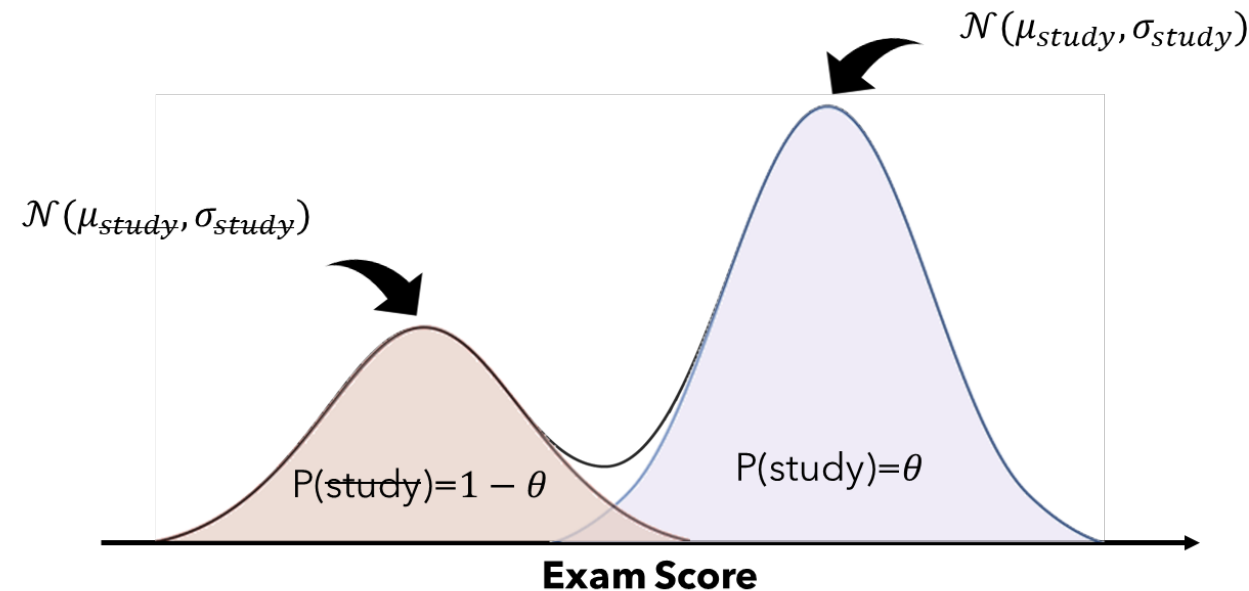
See a new score of 82, did the student study?

$$\begin{aligned} P(study|82) &\propto P(82|study)P(study) \\ &= \mathcal{N}(82; \mu_{study}, \sigma_{study}) \frac{8}{14} = 0.00227 \end{aligned}$$

$$\begin{aligned} P(no\ study|82) &\propto P(82|no\ study)P(no\ study) \\ &= \mathcal{N}(82; \mu_{no\ study}, \sigma_{no\ study}) \frac{6}{14} = 0.0171 \end{aligned}$$



This example had binary labels ($y=\text{study} / \text{no study}$) and only a single continuous feature ($x=\text{exam score}$). But we can do similar things with many more features or classes.





Consider the task of classifying email as either Normal or Spam based on the words in the email.

Output: $y \in \{Normal, Spam\}$ or equivalently $y \in \{0,1\}$

Input: List of M words occurring in an email. How to represent this?

- Popular approach is “bag-of-words”

Option 1: Binary Bag of Words

With a dictionary of size d , an email is represented as a binary vector $x = [x_1, x_2, \dots, x_d]$ where $x_i \in \{0,1\}$

- $x_i = 1$ if the email contains the i^{th} word in the dictionary at least once
- $x_i = 0$ otherwise

Option 2: Multinomial Bag of Words

With a dictionary of size d , an email is represented as a binary vector $x = [x_1, x_2, \dots, x_d]$ where $x_i \in \mathbb{N}_0$

- x_i : non-negative integer occurrence count of the i^{th} word in the dictionary



Option 1: Binary Bag of Words

With a dictionary of size d , an email is represented as a binary vector $x = [x_1, x_2, \dots, x_d]$ where $x_i \in \{0, 1\}$

- $x_i = 1$ if the email contains the i^{th} word in the dictionary at least once
- $x_i = 0$ otherwise

... source for a partner
abroad who can
accommodate HUGE
RESOURCES.

Each document represented by a vector showing which words occurred:

$$\mathbf{x} = \begin{matrix} & \text{partner} & & \text{huge} & & \text{resources} & \\ [0, & 1, & 0, & 0, & 0, & 1, & \dots, & 0, & 0, & 0, & 1, & 0] \\ & \text{assignment} & & \text{research} & & \text{meeting} & \end{matrix}$$

Vector has length equal to the number of words in our vocabulary (d).



A simple generative story for this spam emails:

1. Flip a weighted coin to determine Normal or Spam
2. If **Normal**, sample an \mathbf{x} vector from $P(\mathbf{x}|\text{normal})$
3. If **Spam**, sample an \mathbf{x} vector from $P(\mathbf{x}|\text{spam})$

Would need to learn:

Prior distribution $P(y)$ as the prior probability of being Normal, Spam, or Advertisement

- $P(\text{normal}) = \# \text{ normal} / \text{total}$. $P(\text{spam}) = \# \text{ spam} / \text{total}$. $P(\text{ad}) = \# \text{ ad} / \text{total}$.

Conditional distributions $P(\mathbf{x} | y = \text{normal})$, $P(\mathbf{x} | y = \text{spam})$

- $P(\mathbf{x} | y = \text{normal}) = (\# \text{ normal emails where the exact set of words in } \mathbf{x} \text{ occur}) / (\# \text{ normal emails})$
- $P(\mathbf{x} | y = \text{spam}) = (\# \text{ spam emails where the exact set of words in } \mathbf{x} \text{ occur}) / (\# \text{ spam emails})$

Problem: How do we fit these distributions?



Estimating Joint Distributions over Discrete Variables





Joint Distributions over Discrete Variables

Recipe for writing a joint distribution of M discrete variables:

1. Make a table listing all value combinations of the variables. (M Boolean variables $\rightarrow 2^M$ rows)
2. For each value combination, define how probable it is
3. Due to axioms of probability, these probabilities must sum to 1.

<i>A</i>	<i>B</i>	<i>C</i>	Prob
0	0	0	0.1
0	0	1	0.2
0	1	0	0.3
0	1	1	0.05
1	0	0	0.025
1	0	1	0.025
1	1	0	0.1
1	1	1	0.2

$P(A=0, B=0, C=0)$

$P(A=1, B=0, C=1)$

Example: Joint distribution over Boolean random variables A, B, C



Learning Joint Distributions over Discrete Variables

Recipe for estimating a joint distribution of M discrete variables from data:

1. Make a table listing all value combinations of the variables. (M Boolean variables $\rightarrow 2^M$ rows)
2. For each value combination, make an MLE estimate of the probability of that entry:

$$P(\text{row}) = \frac{\text{examples matching row exactly}}{\text{total number of examples}}$$

(This is the MLE for a categorical distribution)

<i>A</i>	<i>B</i>	<i>C</i>	Prob
0	0	0	?
0	0	1	?
0	1	0	?
0	1	1	?
1	0	0	?
1	0	1	?
1	1	0	?
1	1	1	?

count(1,0,1 examples)
total examples



Example: Joint distribution over Boolean random variables A, B, C



Example of Estimating Joint Distributions over Discrete Variables

Our Data:

Likes Math	Reads ML News	At Least Junior
0	0	1
1	0	1
1	1	1
0	0	1
1	1	1
0	1	1
0	1	1
1	1	1
1	1	0
1	0	0

Out Joint Distribution:

Likes Math	Reads ML News	At Least Junior	Prob
0	0	0	
0	0	1	
0	1	0	
0	1	1	
1	0	0	
1	0	1	
1	1	0	
1	1	1	



Example of Estimating Joint Distributions over Discrete Variables

Our Data:

Likes Math	Reads ML News	At Least Junior
0	0	1
1	0	1
1	1	1
0	0	1
1	1	1
0	1	1
0	1	1
1	1	1
1	1	0
1	0	0

Out Joint Distribution:

Likes Math	Reads ML News	At Least Junior	Prob
0	0	0	0
0	0	1	0.2
0	1	0	0
0	1	1	0.2
1	0	0	0.1
1	0	1	0.1
1	1	0	0.1
1	1	1	0.3



Estimating Joint Distributions over Discrete Variables





Generative Stories - Example 2: Email Classification

Prior distribution $P(y)$ as the prior probability of being Normal, Spam, or Advertisement

- $P(\text{normal}) = \# \text{ normal} / \text{total}$. $P(\text{spam}) = \# \text{ spam} / \text{total}$. $P(\text{ad}) = \# \text{ ad} / \text{total}$.

Conditional distributions $P(\mathbf{x} | y = \text{normal})$, $P(\mathbf{x} | y = \text{spam})$

- $P(\mathbf{x} | y = \text{normal}) = (\# \text{ normal emails where the exact set of words in } \mathbf{x} \text{ occur}) / (\# \text{ normal emails})$
- $P(\mathbf{x} | y = \text{spam}) = (\# \text{ spam emails where the exact set of words in } \mathbf{x} \text{ occur}) / (\# \text{ spam emails})$

$P(\mathbf{x} | y = \text{normal})$

x_1	x_2	...	x_{d-1}	x_d	Prob
0	0	...	0	0	?
0	0	...	0	1	?
0	0	...	1	0	?
0	0	...	1	1	?
0	0	...	0	0	?
:	:	:	:	:	?
1	1	...	1	0	?
1	1	...	1	1	?

$P(\mathbf{x} | y = \text{spam})$

x_1	x_2	...	x_{d-1}	x_d	Prob
0	0	...	0	0	?
0	0	...	0	1	?
0	0	...	1	0	?
0	0	...	1	1	?
0	0	...	0	0	?
:	:	:	:	:	?
1	1	...	1	0	?
1	1	...	1	1	?

These are huge! $2 * 2^d$ entries.



In general, if I measure **d** things in **\mathbf{x}** with each having **m** options, $P(\mathbf{x}|y)$ will have **$c * m^d$** parameters to learn for a **c** class problem. **Yikes.**

The real problem is that we must consider all combinations of features jointly.

The Naïve Bayes Assumption:

Each feature is **conditionally independent** given the class label.

$$P(\mathbf{x}|y) = \prod_{i=1}^d P(x_i|y)$$



Reminders about (conditional) independence





Independence: If two random variables X and Y are independent, then:

$$P(X, Y) = P(X)P(Y), \quad P(X|Y) = P(X), \quad P(Y|X) = P(Y)$$

Conditional Independence: If two random variables X and Y are conditionally independent given Z , then:

$$P(X, Y|Z) = P(X|Z)P(Y|Z), \quad P(X|Y, Z) = P(X|Z), \quad P(Y|X, Z) = P(Y|Z)$$



Reminders about (conditional) independence

Example: Define three random variables:

- **H:** A person's height
- **V:** How many words they know
- **A:** A person's age

Are height and vocabulary size independent?

No. Children tend to have lower vocabularies than adults (on average). So without knowing the person's age: $P(V, H) \neq P(V)P(H)$

However, if I tell you their age... it's likely that these are independent. $P(V, H|A) = P(V|A)P(H|A)$

Events that are dependent in general, can be made independent given some other observation.



def children:
Short people with
small vocabularies.



Reminders about (conditional) independence





The Naïve Bayes Assumption:

Each feature is **conditionally independent** given the class label.

$$P(\mathbf{x}|y) = \prod_{i=1}^d P(x_i|y)$$



How does this help?

$$P(y|\mathbf{x}) \propto P(\mathbf{x}|y)P(y)$$

Bayes Theorem

$$= P(y) \prod_{i=1}^d P(x_i|y)$$

Naïve Bayes
Assumption

Can make predictions this way:

$$\operatorname{argmax}_{c=1,2,3,\dots,k} P(y = c) \prod_{i=1}^d P(x_i|y = c)$$



A “naïve” generative story for this spam emails:

1. Flip a weighted coin to determine Normal or Spam.
2. For each word i , flip a weighted coin to determine if it is included in the email. Let the weight of the coin depend on whether the email is normal or spam.

Would need to learn:

Prior distribution $P(y)$ as the prior probability of being Normal, Spam, or Advertisement

- $P(\text{normal}) = \# \text{ normal} / \text{total}$. $P(\text{spam}) = \# \text{ spam} / \text{total}$. $P(\text{ad}) = \# \text{ ad} / \text{total}$.

Conditional distributions $P(x | y = \text{normal})$, $P(x | y = \text{spam})$

- $P(x_i | y = \text{normal}) = (\# \text{ normal emails with word } i) / (\# \text{ normal emails})$
- $P(x_i | y = \text{spam}) = (\# \text{ spam emails with word } i) / (\# \text{ spam emails})$

Problem: How do we fit these distributions?



Updating our Generative Story for Spam Emails to be "Naïve"

Would need to learn:

Prior distribution $P(y)$ as the prior probability of being Normal, Spam, or Advertisement

- $P(\text{normal}) = \# \text{ normal} / \text{total}$. $P(\text{spam}) = \# \text{ spam} / \text{total}$. $P(\text{ad}) = \# \text{ ad} / \text{total}$.

Conditional distributions $P(x_i | y = \text{normal})$, $P(x_i | y = \text{spam})$

- $P(x_i | y = \text{normal}) = (\# \text{ normal emails with word } i) / (\# \text{ normal emails})$
- $P(x_i | y = \text{spam}) = (\# \text{ spam emails with word } i) / (\# \text{ spam emails})$

$P(x_i | y = \text{spam})$

x_i	Prob
0	?
1	?

x d

$P(x_i | y = \text{normal})$

x_i	Prob
0	?
1	?

x d

Now we only need to learn $2 \times 2d$ values. ($2 \times d$ if you take advantage of axioms)



Parameter Cost of Learning $P(\mathbf{x}|\mathbf{y})$:

In general, if I measure \mathbf{d} things in \mathbf{x} with each having \mathbf{m} options, $P(\mathbf{x}|\mathbf{y})$ will have $\mathbf{c} * \mathbf{m}^{\mathbf{d}}$ free parameters to learn for a \mathbf{c} class problem. **Yikes.**

With the Naïve Bayes assumption, we only need $\mathbf{c} * \mathbf{md}$ for the $P(x_i|\mathbf{y})$ distributions and $\mathbf{c} - 1$ for the class prior $P(\mathbf{y})$ for this setting.



The Naïve Bayes Model Steps:

1. **Learn the conditional** $P(x_i|y = c)$ **for each feature** x_i **and class** c (training)
2. **Estimate** $P(y = c)$ **as a fraction of records with** $y = c$ **for each class** c (training)
3. **For a new example** $x = [x_1, \dots, x_m]^T$, **predict:** (testing)

$$\operatorname{argmax}_{c=1,2,3,\dots,k} P(y = c) \prod_{i=1}^d P(x_i|y = c)$$



An example: Suppose I measure three things about each of you:

- **UpperLowerClass** {Lower, Upper}
- **LikesMath** {Yes, No}
- **ReadsML** {Yes, No}

And want to predict whether you think the class is easy/moderate/hard:

- **CS434** {Easy, Moderate, Hard}

What do examples look like?

([Lower, Yes, No], Moderate), ([Upper, Yes, Yes], Easy),
([Lower, No, Yes], Difficult), ([Upper, Yes, Yes], Moderate)



Now need to learn the following distributions:

$$P(\text{UpperLower}|\text{CS434})$$

$P(\text{lower} \mid \text{easy})$	$P(\text{lower} \mid \text{moderate})$	$P(\text{lower} \mid \text{hard})$
$P(\text{upper} \mid \text{easy})$	$P(\text{upper} \mid \text{moderate})$	$P(\text{upper} \mid \text{hard})$

4 free parameters

$$P(\text{CS434})$$

$P(\text{CS434}=\text{easy})$
$P(\text{CS434}=\text{moderate})$
$P(\text{CS434}=\text{hard})$

2 free parameters

$$P(\text{LikesMath}|\text{CS434})$$

$P(0 \mid \text{easy})$	$P(0 \mid \text{moderate})$	$P(0 \mid \text{hard})$
$P(1 \mid \text{easy})$	$P(1 \mid \text{moderate})$	$P(1 \mid \text{hard})$

3 free parameters

$$P(\text{ReadsML}|\text{CS434})$$

$P(0 \mid \text{easy})$	$P(0 \mid \text{moderate})$	$P(0 \mid \text{hard})$
$P(1 \mid \text{easy})$	$P(1 \mid \text{moderate})$	$P(1 \mid \text{hard})$

3 free parameters



Let's try fitting one of these to the class:

$P(\text{UpperLower}|\text{CS434})$

$P(\text{lower} \mid \text{easy})$	$P(\text{lower} \mid \text{moderate})$	$P(\text{lower} \mid \text{hard})$
$P(\text{upper} \mid \text{easy})$	$P(\text{upper} \mid \text{moderate})$	$P(\text{upper} \mid \text{hard})$

$P(\text{CS434})$

$P(\text{CS434}=\text{easy})$
$P(\text{CS434}=\text{moderate})$
$P(\text{CS434}=\text{hard})$

$P(\text{LikesMath}|\text{CS434})$

$P(0 \mid \text{easy})$	$P(0 \mid \text{moderate})$	$P(0 \mid \text{hard})$
$P(1 \mid \text{easy})$	$P(1 \mid \text{moderate})$	$P(1 \mid \text{hard})$

$P(\text{ReadsML}|\text{CS434})$

$P(0 \mid \text{easy})$	$P(0 \mid \text{moderate})$	$P(0 \mid \text{hard})$
$P(1 \mid \text{easy})$	$P(1 \mid \text{moderate})$	$P(1 \mid \text{hard})$

Series of polls to fill tables: A=Easy B=Moderate C=Hard

The zero-probability problem:

$$\operatorname{argmax}_{c=0,1} P(y = c) \prod_{i=1}^d P(w_i | y = c)$$

If any one of these terms is 0 for an instance, whole thing is 0.

Why might that happen?

What if a new email contains a word we never saw in the training emails?

$$P(w|spam) = 0 \text{ and } P(w|not\ spam) = 0$$

Laplace Smoothing for Binary Variables

For binary variable x_i , add a small prior to $p(x_i|y = c)$:

- Bernoulli

$$p(x_i|y = c) = \frac{(\# \text{ of times } x_i \text{ is true and } y = c) + 1}{(\# \text{ times } y = c) + 2}$$

Taking our $P(\text{LowerUpper} | \text{CS434}=\text{easy})$ example:

$$p(\text{lower} | \text{CS434} = \text{easy}) = \frac{(\# \text{ of fresh\&soph who think cs434 is easy}) + 1}{(\# \text{ fresh\&soph}) + 2}$$

This is just adding a prior to the estimated conditional distributions.
Specifically, a $\text{Beta}(1,1)$ prior.

Laplace Smoothing for Categorical Variables

For categorical variable x_i , add a small prior to $p(x_i|y = c)$:


- Categorical

$$p(x_i = a|y = c) = \frac{(\# \text{ of times } x_i \text{ is } a \text{ and } y = c) + 1}{(\# \text{ times } y = c) + \text{\textcolor{brown}{\#classes}}}$$

This is just adding a prior to the estimated conditional distributions.
Specifically, a Dirichlet prior.

Taking the product of a bunch of probabilities is prone to underflow errors:

$$\operatorname{argmax}_{c=0,1} P(y = c) \prod_{i=1}^d P(x_i | y = c)$$

0.0001*0.001*0.02*0.01*... 

Let's do this in log space:

$$\operatorname{argmax}_{c=0,1} \log P(y = c) + \sum_{i=1}^d \log P(x_i | y = c)$$

- **Generative Model:** Estimate $P(y|\mathbf{x})$ by learning $P(\mathbf{x}|y)$ and $P(y)$. However $P(\mathbf{x}|y)$ can have way too many parameters to be fit effectively.
- **Naïve Bayes Assumption:** Assume features are conditionally independent given the class labels: $P(\mathbf{x}|y) = \prod P(x_i|y)$
- Training a Naïve Bayes classifier comes down to fitting distributions $P(x_i|y)$ and $P(y)$ either with MLE or MAP (MAP is more robust to data sparsity)
- Naïve Bayes is cheap and survives tens of thousands of attributes easily. Also does okay even when conditional independent doesn't hold.
- Any density estimator can be plugged in to estimate $p(x_i|y)$ for Naïve Bayes
 - Real valued attributes can be discretized or directly modeled using simple continuous distributions such as Gaussian (Normal) distribution

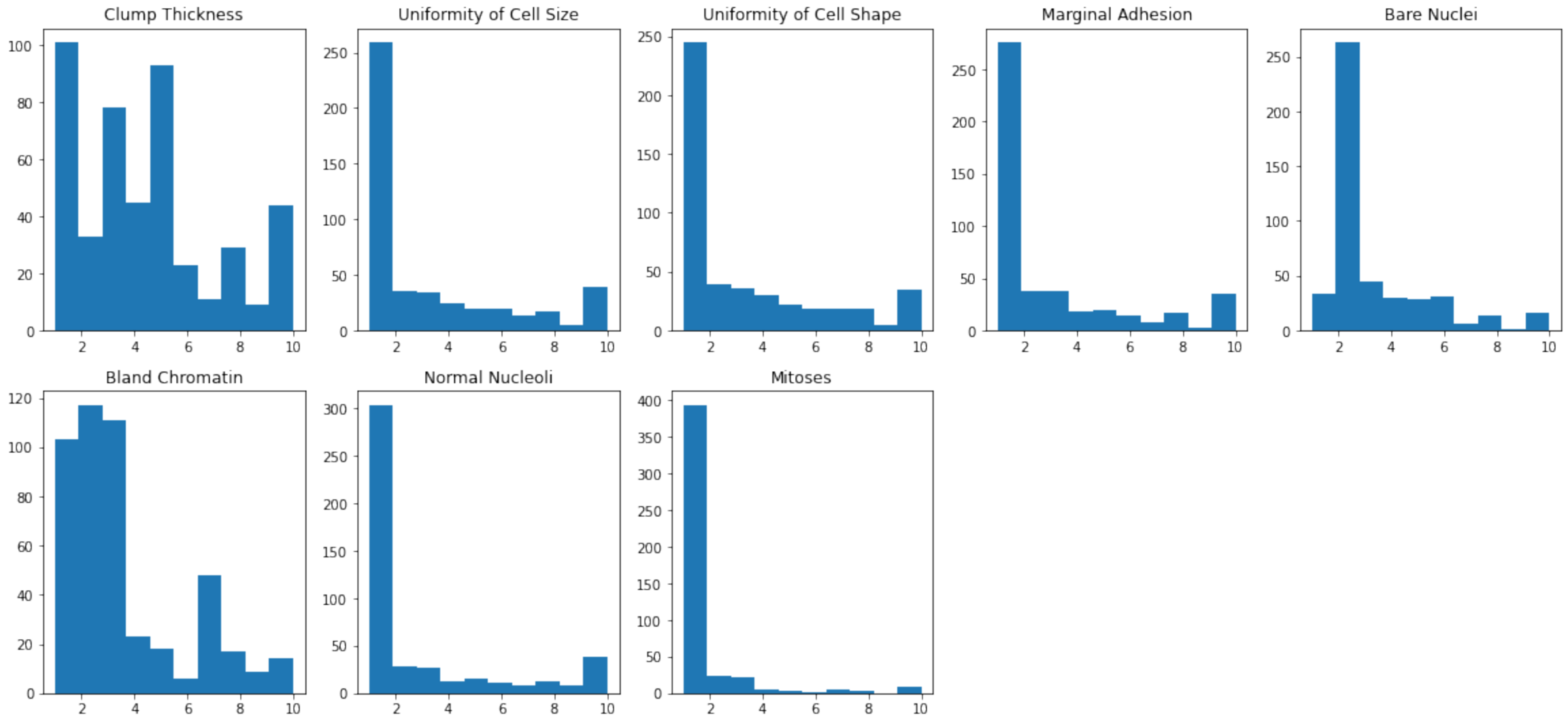


In-Class Coded Example 1:

Categorical Naïve Bayes on Tumor Classification



Categorical Naïve Bayes on Tumor Classification





Let $y \in \{0,1\}$ map to benign / malignant respectively (assume Bernoulli)

- Two parameters $P(y=\text{Benign})$ and $P(y=\text{Malignant})$

For each feature i , assume:

Let $x_i|y \in \{1,2,3,4,5,6,7,8,9,10\}$ be Categorical

- 10×2 parameter $P(x_i = 0 | y = 0), P(x_i = 0 | y = 1), P(x_i = 1 | y = 0), \dots$

Need to learn the distributions:

- $P(y)$ and $P(x_i|y_i) \forall i$

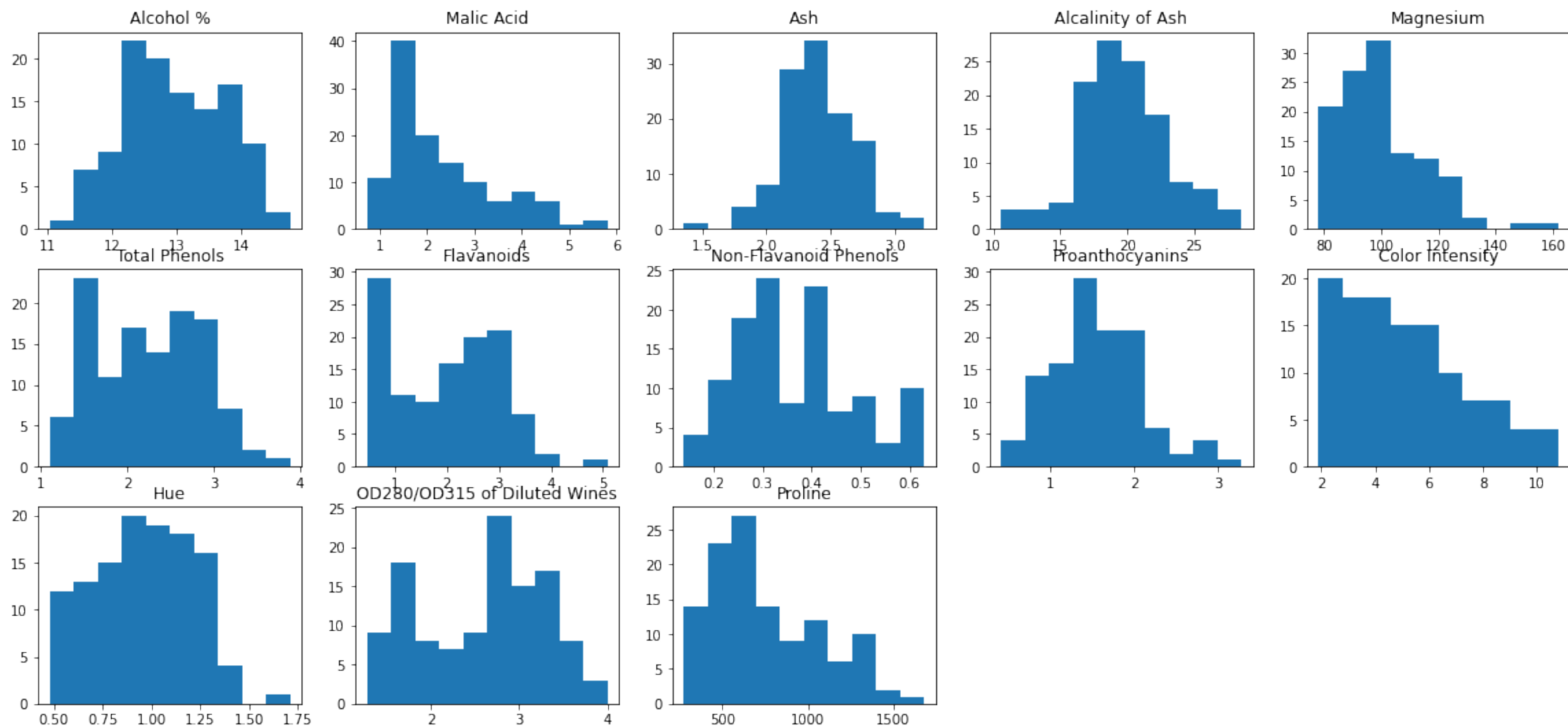


In-Class Coded Example 2:

Gaussian Naïve Bayes on Winegrape Classification



Gaussian Naïve Bayes on Winegrape Classification





Gaussian Naïve Bayes on Winegrape Classification

Let $y \in \{0,1,2\}$ map to three different types of grapes (assume Categorical)

- Two parameters $P(y=G1)$, $P(y=G2)$, and $P(y=G3)$

For each feature i , assume:

Let $x_i|y \in \{-\infty, \infty\}$ be Gaussian

- $3*2$ parameters -- $\mu_{G1}, \sigma_{G1}, \mu_{G2}, \sigma_{G2}, \mu_{G3}, \sigma_{G3}$

Need to learn the distributions:

- $P(y)$ and $P(x_i|y_i) \forall i$



You can check out the Colab notebook here:

https://colab.research.google.com/drive/1_7306psnP0FgrgUgAGQUm8XEL_3QZ6FO?usp=sharing



Next Time: We'll talk about maximum margin classifiers – specifically SVMs!