# Modeling with Naïve Bayes: Mathematics, Example with School Grades via Hand, & Example with Titanic via Python

SAIL ON • OCTOBER 2016

PAMELA TOMAN • PTOMAN@CS.STANFORD.EDU

# Where are we going tonight?

Educational agenda:

- We're going to model tonight! ...what is that again?
- So what is Kaggle?
- Let's review probabilities...
- How does Naïve Bayes get us a model?
- Let's code this thing!

### Machine learning has three parts

We want our AI to be able to learn (robots? monitoring & labeling danger? personalized help?)

In machine learning, we provide (1) a mathematical equation with lots of unknown variables (a model) and (2) data, and we ask the computer to learn what the best values for the variables are

We implement the mathematical equation as an *algorithm* in some programming language



# What's "a mathematical equation with lots of unknown variables (a model)"?

We choose/we are given (or both) some input variables *x*, *y*, *z* and some target variable *t*.

Our model is an equation for how the input variables x, y, z can be used to produce the target t.

# What's "a mathematical equation with lots of unknown variables (a model)"?

We choose/we are given (or both) some input variables *x*, *y*, *z* and some target variable *t*.

Our model is an equation for how the input variables x, y, z can be used to produce the target t.

Maybe x is hat size, y is weight, and z is foot length. We want to predict t, which is height.

We might think they are "linearly related":

t = ax + by + cz + d

This equation is entirely variables!

- Any individual has their own *x*, *y*, *z*, *t*, so these are variables.
- We don't know *a*, *b*, *c*, or *d*, so these are also variables. (But we can learn estimates for these!)

# What's "a mathematical equation with lots of unknown variables (a model)"?

We choose/we are given (or both) some input variables *x*, *y*, *z* and some target variable *t*.

Our model is an equation for how the input variables x, y, z can be used to produce the target t.

Maybe x is hat size, y is weight, and z is foot length. We want to predict t, which is height.

We might think they are "linearly related":

t = ax + by + cz + d

This equation is entirely variables!

- Any individual has their own *x*, *y*, *z*, *t*, so these are variables.
- We don't know *a*, *b*, *c*, or *d*, so these are also variables. (But we can learn estimates for these!)

We can answer the question with other models too. Tonight we'll talk about Naïve Bayes, which uses probabilities. This time we let *t* be "height > 5.5ft":

$$P(t|x, y, z) = \frac{P(x|t)P(y|t)P(z|t)P(t)}{P(x, y, z)}$$

This equation is still entirely variables!

- Any individual has their own x, y, z, t, so these are variables.
- We don't know P(x|t), P(y|t), P(z|t), P(t), or P(x,y,z), so these are also variables. (But we can learn estimates!)

Training Data: x, y, z (each row contains attributes)

hat size, weight, foot length

Training Data: x, y, z (each row contains attributes)

hat size, weight, foot length

Training Data: t (each row has a target variable)

height







### What is Kaggle?

Kaggle is at <a href="https://www.kaggle.com/">https://www.kaggle.com/</a>

Kaggle is a platform for "data science" – the application of machine learning to problems

- A company (usually) provides data & a target variable
- Participants build models to predict the target variable and they submit their predictions
- Participant solutions are ranked twice: (1) a public leaderboard visible to participants, (2) a private leaderboard visible only to the problem provider (... why twice?)

Participants get bragging rights & sometimes money or jobs

One of their "Getting Started" challenges is predicting who survives on the Titanic

We're going to build a mathematical foundation for Naïve Bayes, and then work on the Titanic

Select 2 cards from a deck of 52 cards with replacement. What's the probability of obtaining 2 kings?

Select 2 cards from a deck of 52 cards *without* replacement. What's the probability of obtaining 2 kings?

Select 2 cards from a deck of 52 cards with replacement. What's the probability of obtaining 2 kings? "Independent" events ~0.0059

Select 2 cards from a deck of 52 cards *without* replacement. What's the probability of obtaining 2 kings?

Select 2 cards from a deck of 52 cards with replacement. What's the probability of obtaining 2 kings? "Independent" events ~0.0059

Select 2 cards from a deck of 52 cards *without* replacement. What's the probability of obtaining 2 kings?

"Dependent" events ~0.0045

Select 2 cards from a deck of 52 cards with replacement. What's the probability of obtaining 2 kings? "Independent" events ~0.0059

If two events A and B are independent events,

then the probability of event A *and* B is given by the following rule:

P(A, B) = P(A) \* P(B)

We read P(X, Y) as "the probability of X and Y (both occurring)".

Select 2 cards from a deck of 52 cards *without* replacement. What's the probability of obtaining 2 kings? "Dependent" events ~0.0045

If two events A and B are dependent events,

then the probability of event A and B is given by the following rule:

P(A, B) = P(A) \* P(B|A) = P(B) \* P(A|B)

Here, P(X|Y) is a "conditional probability": the probability that an event X will occur given that Y has already occurred.

We read P(X|Y) as "the probability of X, given Y".

We roll 2 fair dice, A and B. We write the possible A+B outcomes in a table.

We roll 2 fair dice, A and B. We write the possible A+B outcomes in a table.

What's the probability that A = 2?

We roll 2 fair dice, A and B. We write the possible A+B outcomes in a table.

What's the probability that A = 2?

+		В							
		1	2	3	4	5	6		
	1	2	3	4	5	6	7		
	2	3	4	5	6	7	8		
•	3	4	5	6	7	8	9		
A	4	5	6	7	8	9	10		
	5	6	7	8	9	10	11		
	6	7	8	9	10	11	12		

We roll 2 fair dice, A and B. We write the possible A+B outcomes in a table.

What's the probability that A = 2?



P(A=2) = 6/36 = 0.167

We roll 2 fair dice, A and B. We write the possible A+B outcomes in a table.

В + Α 

What's the probability that A = 2?

P(A=2) = 6/36 = 0.167

It's revealed that  $A+B \leq 5$ .

We roll 2 fair dice, A and B. We write the possible A+B outcomes in a table.



It's revealed that  $A+B \leq 5$ .

+		В							
		1	2	3	4	5	6		
Α	1	2	3	4	5	6	7		
	2	3	4	5	6	7	8		
	3	4	5	6	7	8	9		
	4	5	6	7	8	9	10		
	5	6	7	8	9	10	11		
	6	7	8	9	10	11	12		

P(A=2) = 6/36 = 0.167

We roll 2 fair dice, A and B. We write the possible A+B outcomes in a table.



It's revealed that  $A+B \leq 5$ .

в + Α 

Now what's the probability that A = 2?

P(A=2) = 6/36 = 0.167

We roll 2 fair dice, A and B. We write the possible A+B outcomes in a table.



It's revealed that  $A+B \leq 5$ .



Now what's the probability that A = 2?



 $P(A=2 | A+B \le 5) = 3/10 = 0.300$ 

P(A=2) = 6/36 = 0.167

source

We roll 2 fair dice, A and B. We write the possible A+B outcomes in a table.



It's revealed that  $A+B \leq 5$ .



Now what's the probability that A = 2?



P(A=2) = 6/36 = 0.167



The probability that A=2 went up given the additional information!

source

Here's the conditional probability rule again:

P(A, B) = P(A) \* P(B|A) = P(B) \* P(A|B)

Here's the conditional probability rule again:	Relative size	Case B	Case B	Total		
	Condition A	w	X	w+x		
$P(A, D) = P(A) \cdot P(D A) = P(D) \cdot P(A D)$	Condition Ā	У	z	y+z		
	Total	w+y	x+z	w+x+y+z		
Let's look at why it works in a picture:						
(Spend a bit of time reflecting on why this works.)						
	$P(A B) \times P(B) = \frac{w}{w+y} \times \frac{w+y}{w+x+y+z} = \frac{w}{w+x+y+z}$					
	$P(B A) \times P(A) = \frac{w}{w+x} \times \frac{w+x}{w+x+y+z} = \frac{w}{w+x+y+z}$				source	

Here's the conditional probability rule again:	Relative size	Case B	Case B	Total		
$D(\Lambda   \mathbf{P}) = D(\Lambda) * D(\mathbf{P}   \Lambda) = D(\mathbf{P}) * D(\Lambda   \mathbf{P})$	Condition A	w	X	w+x		
$P(A, D) = P(A)^{-1}P(D A) = P(D)^{-1}P(A D)$	Condition Ā	У	z	y+z		
	Total	w+y	x+z	w+x+y+z		
Let's look at why it works in a picture:						
(Spend a bit of time reflecting on why this works.)						
	$P(A B) \times P(B)$	$(3) = \frac{w}{w+1}$	$\overline{y} \times \frac{w+y}{w+x+y}$	$\frac{y}{y+z} = \frac{w}{w+x+y+z}$		
Let's consider some more examples:						
P(K1  and  K2) = P(K1) * P(K2 K1) = P(K2) * P(K1 K2)	$P(B A) \times P(A)$	$A) = \frac{w}{w+z}$	$\frac{w+x}{x}$ $\frac{w+x}{w+x+y}$	$\frac{x}{y+z} = \frac{w}{w+x+y+z}$	source	



```
Here's the conditional probability rule again:
                                                                          Relative size Case B Case B
                                                                                                         Total
                                                                          Condition A
                                                                                         w
                                                                                                         w+x
                                                                                                 x
        P(A, B) = P(A) * P(B|A) = P(B) * P(A|B)
                                                                          Condition Ā
                                                                                                         y+z
                                                                                         ν
                                                                                                 z
                                                                              Total
                                                                                       w+y
                                                                                                x+z
                                                                                                      w+x+y+z
Let's look at why it works in a picture:
(Spend a bit of time reflecting on why this works.)
                                                                       P(A|B) \times P(B) = \frac{w}{w+w} \times \frac{w+y}{w+x+v+z} = \frac{w}{w+x+v+z}
Let's consider some more examples:
                                                                       P(B|A) \times P(A) = \frac{w}{w+x} \times \frac{w+x}{w+x+y+z} = \frac{w}{w+x+y+z}
P(K1 \text{ and } K2) = P(K1) * P(K2|K1) = P(K2) * P(K1|K2)
                                                                                                                       source
P(A \text{ and passed}) = P(A) * P(passed|A) = P(passed) * P(A|passed)
P(A+B \le 5 \text{ and } A=2) = P(A+B \le 5) * P(A=2 | A+B \le 5) = P(A=2) * P(A+B \le 5 | A=2)
```

Here's our conditional probability rule again:

P(A, B) = P(A) \* P(B|A) = P(B) \* P(A|B)

## Bayes' Law (= Bayes' Theorem, Bayes' Rule)

Here's our conditional probability rule again:

```
P(A, B) = P(A) * P(B|A) = P(B) * P(A|B)
```

Let's rewrite it with P(B|A) on the left using algebra:

P(A) \* P(B|A) = P(B) \* P(A|B)

Here's our conditional probability rule again:

P(A, B) = P(A) \* P(B|A) = P(B) \* P(A|B)

Let's rewrite it with P(B|A) on the left using algebra:

P(A) \* P(B|A) = P(B) \* P(A|B)P(B|A) = P(B) \* P(A|B)P(A)

Here's our conditional probability rule again:

```
P(A, B) = P(A) * P(B|A) = P(B) * P(A|B)
```

Let's rewrite it with P(B|A) on the left using algebra:

```
P(A) * P(B|A) = P(B) * P(A|B)
P(B|A) = \frac{P(B) * P(A|B)}{P(A)}
P(B|A) = \frac{P(A|B) * P(B)}{P(A)}
```

This equation is called "Bayes' Law"!
Remember Bayes' law:

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

Remember Bayes' law:

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

Let's apply it.

Everyone who gets an A will pass. The average passing rate is 85%, and 20% of the class gets As. Of the people who pass, what percentage gets As?

Remember Bayes' law:

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

Let's apply it.

Everyone who gets an A will pass. The average passing rate is 85%, and 20% of the class gets As. Of the people who pass, what percentage gets As?

```
P(A|passed) = P(passed|A) * P(A)
P(passed)
```

Remember Bayes' law:

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

Let's apply it.

Everyone who gets an A will pass. The average passing rate is 85%, and 20% of the class gets As. Of the people who pass, what percentage gets As?

P(A|passed) = <u>P(passed|A) \* P(A)</u> P(passed)

Remember Bayes' law:

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

Let's apply it.

Everyone who gets an A will pass. The average passing rate is 85%, and 20% of the class gets As. Of the people who pass, what percentage gets As?

```
P(A|passed) = \frac{P(passed|A) * P(A)}{P(passed)}P(A|passed) = \frac{1.00 * 0.20}{0.85}P(A|passed) \approx 0.23
```

Remember Bayes' law:

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

Let's apply it.

Everyone who gets an A will pass. The average passing rate is 85%, and 20% of the class gets As. Of the people who pass, what percentage gets As?



Okay, back to modeling....

Let's say we want to infer whether you got an A.

Let's say we want to infer whether you got an A.

You are very shy and say nothing, but some classmates are more forthcoming.

Let's say we want to infer whether you got an A.

You are very shy and say nothing, but some classmates are more forthcoming.

We do know a few things about you:

- Did you pass the class?
- Did you look happy when grades came out last week?
- What's your GPA (or our guess at it)?
- What's the average grade of your three best friends in the class?

Let's say we want to infer whether you got an A.

You are very shy and say nothing, but some classmates are more forthcoming.

We do know a few things about you:

- Did you pass the class?
- Did you look happy when grades came out last week?
- What's your GPA (or our guess at it)?
- What's the average grade of your three best friends in the class?

We create a table for our information:

Let's say we want to infer whether you got an A.

You are very shy and say nothing, but some classmates are more forthcoming.

We do know a few things about you:

- Did you pass the class?
- Did you look happy when grades came out last week?
- What's your GPA (or our guess at it)?
- What's the average grade of your three best friends in the class?

#### We create a table for our information:

Person	Passed?	Happy?	GPA?	Friends'?	Grade
1	Yes	No	3.8	А	А
2	No	No	2.1	D	F
99 (you!)	Yes	Yes	3.9	В	????

We want to know:

We want to know:

P(A | Yes, Yes, 3.9, B)

We want to know:

P(A | Yes, Yes, 3.9, B)

Bayes' rule means we could plug values into this equation & get the answer:

We want to know:

P(A | Yes, Yes, 3.9, B)

Bayes' rule means we could plug values into this equation & get the answer:

P(A | Yes, Yes, 3.9, B) = <u>P(Yes, Yes, 3.9, B | A) \* P(A)</u>

P(Yes, Yes, 3.9, B)

We want to know:

P(A | Yes, Yes, 3.9, B)

Bayes' rule means we could plug values into this equation & get the answer:

P(A | Yes, Yes, 3.9, B) = <u>P(Yes, Yes, 3.9, B | A) \* P(A)</u>

P(Yes, Yes, 3.9, B)

But this is *hard*!! We don't have nearly enough data.

We want to know:

P(A | Yes, Yes, 3.9, B)

Bayes' rule means we could plug values into this equation & get the answer:

But this is *hard*!! We don't have nearly enough data.

 $\rightarrow$  In particular:

Looking only at the people in the class who got As, how many people match your answers?

Probably 0. But that means we have small data, not that you didn't get an A...

We can resolve the small data problem through making a "naïve" assumption:

We can resolve the small data problem through making a "naïve" assumption:

Within each target class (A or not-A), the probability of each variable is independent of others:

- Probability of passing is independent from happiness, GPA, and friends, among those who (didn't) get an A
- Happiness is independent from passing, GPA, and friends, among those who (didn't) get an A

•

So instead of:

We can resolve the small data problem through making a "naïve" assumption:

Within each target class (A or not-A), the probability of each variable is independent of others:

- Probability of passing is independent from happiness, GPA, and friends, among those who (didn't) get an A
- Happiness is independent from passing, GPA, and friends, among those who (didn't) get an A

•

So instead of:

P(A | Yes, Yes, A, B) = <u>P(Yes, Yes, 3.9, B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

We will make a "Naïve Bayes" assumption and instead use:

We can resolve the small data problem through making a "naïve" assumption:

Within each target class (A or not-A), the probability of each variable is independent of others:

- Probability of passing is independent from happiness, GPA, and friends, among those who (didn't) get an A
- Happiness is independent from passing, GPA, and friends, among those who (didn't) get an A

•

So instead of:

We will make a "Naïve Bayes" assumption and instead use:

We can resolve the small data problem through making a "naïve" assumption:

Within each target class (A or not-A), the probability of each variable is independent of others:

- Probability of passing is independent from happiness, GPA, and friends, among those who (didn't) get an A
- Happiness is independent from passing, GPA, and friends, among those who (didn't) get an A

•

So instead of:

# If we just want *whether* you got an A, we can simplify even more....

Let's keep simplifying. So currently our formula will tell us the *probability* you got an A.

P(A | Yes, Yes, A, B) = <u>P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

If we just want to know *whether* you got an A, we can simplify even more.

# If we just want *whether* you got an A, we can simplify even more....

Let's keep simplifying. So currently our formula will tell us the *probability* you got an A.

#### P(A | Yes, Yes, A, B) = <u>P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

If we just want to know *whether* you got an A, we can simplify even more.

Whether you got an A is just:

#### P(A | whatever) >? P(not A | whatever)

 P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A)
 ?
 P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

 P(Yes, Yes, 3.9, B)
 P(Yes, Yes, 3.9, B)
 P(Yes, Yes, 3.9, B)

Since the denominator is the same, we can drop the denominator.... The key here is that whether you got an A or you didn't, we still have the same observations.

Starting with the definition of conditional probability:

Starting with the definition of conditional probability:

P(A, Yes, Yes, 3.9, B) = P(A) \* P(Yes, Yes, 3.9, B|A) = P(Yes, Yes, 3.9, B) \* P(A|Yes, Yes, 3.9, B)

Starting with the definition of conditional probability:

P(A, Yes, Yes, 3.9, B) = P(A) \* P(Yes, Yes, 3.9, B|A) = P(Yes, Yes, 3.9, B) \* P(A|Yes, Yes, 3.9, B)

Rearranging the equations into a format known as "Bayes' law":

Starting with the definition of conditional probability:

P(A, Yes, Yes, 3.9, B) = P(A) \* P(Yes, Yes, 3.9, B|A) = P(Yes, Yes, 3.9, B) \* P(A|Yes, Yes, 3.9, B)

Rearranging the equations into a format known as "Bayes' law":

P(A | Yes, Yes, 3.9, B) = <u>P(Yes, Yes, 3.9, B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Starting with the definition of conditional probability:

P(A, Yes, Yes, 3.9, B) = P(A) \* P(Yes, Yes, 3.9, B|A) = P(Yes, Yes, 3.9, B) \* P(A|Yes, Yes, 3.9, B)

Rearranging the equations into a format known as "Bayes' law":

P(A | Yes, Yes, 3.9, B) = <u>P(Yes, Yes, 3.9, B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Deciding we can deal with small data through the "Naïve Bayes" assumption of independence:

Starting with the definition of conditional probability:

P(A, Yes, Yes, 3.9, B) = P(A) \* P(Yes, Yes, 3.9, B|A) = P(Yes, Yes, 3.9, B) \* P(A|Yes, Yes, 3.9, B)

Rearranging the equations into a format known as "Bayes' law":

P(A | Yes, Yes, 3.9, B) = <u>P(Yes, Yes, 3.9, B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Deciding we can deal with small data through the "Naïve Bayes" assumption of independence:

P(A | Yes, Yes, 3.9, B) = <u>P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Starting with the definition of conditional probability:

P(A, Yes, Yes, 3.9, B) = P(A) \* P(Yes, Yes, 3.9, B|A) = P(Yes, Yes, 3.9, B) \* P(A|Yes, Yes, 3.9, B)

Rearranging the equations into a format known as "Bayes' law":

P(A | Yes, Yes, 3.9, B) = <u>P(Yes, Yes, 3.9, B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Deciding we can deal with small data through the "Naïve Bayes" assumption of independence:

P(A | Yes, Yes, 3.9, B) = <u>P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Identifying we just want *whether*, not a probability:

Starting with the definition of conditional probability:

P(A, Yes, Yes, 3.9, B) = P(A) \* P(Yes, Yes, 3.9, B|A) = P(Yes, Yes, 3.9, B) \* P(A|Yes, Yes, 3.9, B)

Rearranging the equations into a format known as "Bayes' law":

P(A | Yes, Yes, 3.9, B) = <u>P(Yes, Yes, 3.9, B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Deciding we can deal with small data through the "Naïve Bayes" assumption of independence:

P(A | Yes, Yes, 3.9, B) = <u>P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Identifying we just want *whether*, not a probability:

 P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A)
 P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

 P(Yes, Yes, 3.9, B)
 P(Yes | Yes | not A) \* P(S, Yes, 3.9, B)
# So finally we have a tractable (usable) formula for whether you got an A...

Starting with the definition of conditional probability:

P(A, Yes, Yes, 3.9, B) = P(A) \* P(Yes, Yes, 3.9, B|A) = P(Yes, Yes, 3.9, B) \* P(A|Yes, Yes, 3.9, B)

Rearranging the equations into a format known as "Bayes' law":

P(A | Yes, Yes, 3.9, B) = <u>P(Yes, Yes, 3.9, B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Deciding we can deal with small data through the "Naïve Bayes" assumption of independence:

P(A | Yes, Yes, 3.9, B) = <u>P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Identifying we just want *whether*, not a probability:

 P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A)
 ?
 P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

 P(Yes, Yes, 3.9, B)
 P(Yes, Yes, 3.9, B)
 P(Yes, Yes, 3.9, B)

Recognizing the denominator doesn't matter:

# So finally we have a tractable (usable) formula for whether you got an A...

Starting with the definition of conditional probability:

P(A, Yes, Yes, 3.9, B) = P(A) \* P(Yes, Yes, 3.9, B|A) = P(Yes, Yes, 3.9, B) \* P(A|Yes, Yes, 3.9, B)

Rearranging the equations into a format known as "Bayes' law":

P(A | Yes, Yes, 3.9, B) = <u>P(Yes, Yes, 3.9, B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Deciding we can deal with small data through the "Naïve Bayes" assumption of independence:

P(A | Yes, Yes, 3.9, B) = <u>P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A)</u> P(Yes, Yes, 3.9, B)

Identifying we just want *whether*, not a probability:

 P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A)
 ?
 P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

 P(Yes, Yes, 3.9, B)
 P(Yes | Not A) \* P(Yes | not A) \* P(B | not A) \* P(not A)

Recognizing the denominator doesn't matter:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

So here is our formula:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

Using the data we collected from your forthcoming classmates, we can estimate *all* these probs.

Passed?	Нарру?	GPA?	Friends'?	Grade			Dassed	- Voc	Passed - No
Yes	No	3.8	А	А			Passeu	- 165	rasseu – No
No	No	2.1	D	F		Α			
Yes	Yes	3.1	В	В		Not A			
No	No	3.3	А	F					
Yes	Yes	3.2	В	В		CDA	<b>CDA</b>	004	<b>CD4</b>
						GPA	GPA	GPA	GPA
Yes	Yes	3.9	В	????		> 3.5	3.0-3.5	2.5-3.0	J 1.5-2.5
					Α				
	Passed?YesNoYesNoYesYes	Passed?Happy?YesNoNoNoYesYesNoNoYesYesYesYesYesYesYesYes	Passed?         Happy?         GPA?           Yes         No         3.8           No         No         2.1           Yes         Yes         3.1           No         No         3.3           Yes         Yes         3.2                Yes         Yes         3.9	Passed?Happy?GPA?Friends'?YesNo3.8ANoNo2.1DYesYes3.1BNoNo3.3AYesYes3.2BYesYes3.9B	Passed?Happy?GPA?Friends'?GradeYesNo3.8AANoNo2.1DFYesYes3.1BBNoNo3.3AFYesYes3.2BBYesYes3.9B????	Passed?Happy?GPA?Friends'?GradeYesNo3.8AANoNo2.1DFYesYes3.1BBNoNo3.3AFYesYes3.2BBYesYes3.9B???	Passed?Happy?GPA?Friends'?GradeYesNo3.8AANoNo2.1DFYesYes3.1BBNoNo3.3AFYesYes3.2BBYesYes3.9B???AAA <td>Passed?Happy?GPA?Friends'?GradeYesNo3.8AANoNo2.1DFYesYes3.1BBNoNo3.3AFYesYes3.2BBYesYes3.9B???AModelAAYesYes3.9BYesYes3.9AYes<t< td=""><td>Passed?         Happy?         GPA?         Friends'?         Grade           Yes         No         3.8         A         A           No         No         2.1         D         F           Yes         Yes         3.1         B         B           No         No         3.3         A         F           Yes         Yes         3.2         B         B               GPA         GPA         GPA           Yes         Yes         3.9         B         ????         A             A         Yes         3.9         B         ????         A         </td></t<></td>	Passed?Happy?GPA?Friends'?GradeYesNo3.8AANoNo2.1DFYesYes3.1BBNoNo3.3AFYesYes3.2BBYesYes3.9B???AModelAAYesYes3.9BYesYes3.9AYes <t< td=""><td>Passed?         Happy?         GPA?         Friends'?         Grade           Yes         No         3.8         A         A           No         No         2.1         D         F           Yes         Yes         3.1         B         B           No         No         3.3         A         F           Yes         Yes         3.2         B         B               GPA         GPA         GPA           Yes         Yes         3.9         B         ????         A             A         Yes         3.9         B         ????         A         </td></t<>	Passed?         Happy?         GPA?         Friends'?         Grade           Yes         No         3.8         A         A           No         No         2.1         D         F           Yes         Yes         3.1         B         B           No         No         3.3         A         F           Yes         Yes         3.2         B         B               GPA         GPA         GPA           Yes         Yes         3.9         B         ????         A             A         Yes         3.9         B         ????         A

So here is our formula:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

Using the data we collected from your forthcoming classmates, we can estimate *all* these probs.

?	Happy?	GPA?	Friends'?	Grade
	No	3.8	А	А
	No	2.1	D	F
	Yes	3.1	В	В
	No	3.3	А	F
	Yes	3.2	В	В
	Yes	3.9	В	????

So here is our formula:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

Using the data we collected from your forthcoming classmates, we can estimate *all* these probs.

		Passed = Yes
	Α	A 1
	Not A	Not A
	CDA	
	GPA	GPA GPA GPA
	> 3.5	> 3.5 3.0-3.5 2.5-3
Α	Α	A
	Not A GPA > 3.5	Not A         GPA         GPA         GPA         Secondary         GPA         Secondary         GPA         Secondary         Secon

So here is our formula:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

Using the data we collected from your forthcoming classmates, we can estimate *all* these probs.

Person	Passed?	Нарру?	GPA?	Friends'?	Grade
1	Yes	No	3.8	А	А
2	No	No	2.1	D	F
\$	Yes	Yes	3.1	В	В
ļ	No	No	3.3	А	F
	Yes	Yes	3.2	В	В
99 (you!)	Yes	Yes	3.9	В	????

So here is our formula:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

Using the data we collected from your forthcoming classmates, we can estimate *all* these probs.

Person	Passed?	Нарру?	GPA?	Friends'?	Grade
1	Yes	No	3.8	А	А
2	No	No	2.1	D	F
3	Yes	Yes	3.1	В	В
4	No	No	3.3	А	F
5	Yes	Yes	3.2	В	В
99 (you!)	Yes	Yes	3.9	В	????

So here is our formula:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

Using the data we collected from your forthcoming classmates, we can estimate *all* these probs.

Person	Passed?	Нарру?	GPA?	Friends'?	Grade
1	Yes	No	3.8	А	А
2	No	No	2.1	D	F
3	Yes	Yes	3.1	В	В
4	No	No	3.3	А	F
5	Yes	Yes	3.2	В	В
99 (you!)	Yes	Yes	3.9	В	????

So here is our formula:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

Using the data we collected from your forthcoming classmates, we can estimate *all* these probs.

erson	Passed?	Нарру?	GPA?	Friends'?	Grade			Passed	= Ves	Passed
1	Yes	No	3.8	А	А			Tassea	- 105	Tusseu -
2	No	No	2.1	D	F	Α		1		0
3	Yes	Yes	3.1	В	В	N	ot A	2		2
4	No	No	3.3	А	F					-
5	Yes	Yes	3.2	В	В		CDA	CDA	CDA	CDA
							GPA	GPA	GPA	GPA
99 (you!)	Yes	Yes	3.9	В	????		> 3.5	3.0-3.5	2.5-3.	.0 1.5-2.
						Α	1	0	0	0

So here is our formula:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

Using the data we collected from your forthcoming classmates, we can estimate *all* these probs.

			THEIR I	Grade			Passed :		Passed = No
Yes	No	3.8	А	А			T USSEU -	103	
No	No	2.1	D	F	Α		1		0
Yes	Yes	3.1	В	В	N	ot A	2		2
No	No	3.3	А	F					_
Yes	Yes	3.2	В	В		CDA	CDA	CDA	CDA
						GPA	GPA 2 0 2 5		GPA
Yes	Yes	3.9	В	????		> 3.5	3.0-3.5	2.5-3.0	1.5-2.5
					Α	1	0	0	0
	No Yes No Yes  Yes	NoNoNoNoYesYesNoNoYesYesYesYes	NoNo3.0NoNo2.1YesYes3.1NoNo3.3YesYes3.2YesYes3.9	NoNoS.6ANoNo2.1DYesYes3.1BNoNo3.3AYesYes3.2BYesYes3.9B	NoNoS.0AANoNo2.1DFYesYes3.1BBNoNo3.3AFYesYes3.2BBYesYes3.9B????	NoNoS.0AAANoNo2.1DFAYesYes3.1BBNNoNo3.3AFYesYes3.2BBYesYes3.9B???AA <td< td=""><td>NoNo3.0AAANoNo2.1DFANot AYesYes3.1BBNot ANoNo3.3AFAYesYes3.2BBGPA &gt; 3.5YesYes3.9B???AMMathematical StructureImage: StructureAAYesYes3.9B???AMMathematical StructureImage: StructureAA</td><td>NoNo3.0AAANoNo2.1DFAAYes3.1BBNot A2NoNo3.3AFDAYesYes3.2BBBAGPA &gt; 3.5GPA 3.0-3.5YesYes3.9B???A1</td><td>NoNoS.5AAANoNo2.1DFYes3.1BB<math>Not A</math>2NoNo3.3AFYesYes3.2BBYesYes3.9B???A100</td></td<>	NoNo3.0AAANoNo2.1DFANot AYesYes3.1BBNot ANoNo3.3AFAYesYes3.2BBGPA > 3.5YesYes3.9B???AMMathematical StructureImage: StructureAAYesYes3.9B???AMMathematical StructureImage: StructureAA	NoNo3.0AAANoNo2.1DFAAYes3.1BBNot A2NoNo3.3AFDAYesYes3.2BBBAGPA > 3.5GPA 3.0-3.5YesYes3.9B???A1	NoNoS.5AAANoNo2.1DFYes3.1BB $Not A$ 2NoNo3.3AFYesYes3.2BBYesYes3.9B???A100

Not A

So here is our formula:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

Using the data we collected from your forthcoming classmates, we can estimate *all* these probs.

on	Passed?	Нарру?	GPA?	Friends'?	Grade			Passed	= Yes	Passed = No	
	Yes	No	3.8	А	А			T d55Cd ·			
2	No	No	2.1	D	F	Α	<b>\</b>	1		0	
3	Yes	Yes	3.1	В	В	N	lot A	2		2	
4	No	No	3.3	А	F					-	
5	Yes	Yes	3.2	В	В		CDA	CDA	CDA	CDA	
										GPA	
99 (you!)	Yes	Yes	3.9	В	????		> 3.5	3.0-3.5	2.5-3.0	1.5-2.5	
						Α	1	0	0	0	

Not A

3

So here is our formula:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

Using the data we collected from your forthcoming classmates, we can estimate *all* these probs.

erson	Passed?	Нарру?	GPA?	Friends'?	Grade			Passed	= Yes	Passed =
1	Yes	No	3.8	А	А	_		Tusseu	- 103	1 45504 -
2	No	No	2.1	D	F	Α		1		0
3	Yes	Yes	3.1	В	В	N	ot A	2		2
4	No	No	3.3	A	F					_
5	Yes	Yes	3.2	В	В		CDA	CDA	CDA	CDA
							GPA	GPA	GPA	GPA
99 (you!)	Yes	Yes	3.9	В	????		> 3.5	3.0-3.5	2.5-3.0	J 1.5-2.5
						Α	1	0	0	0

Not A

3

1

So here is our formula:

P(Yes | A) \* P(Yes | A) \* P(3.9 | A) \* P(B | A) \* P(A) >? P(Yes | not A) \* P(Yes | not A) \* P(3.9 | not A) \* P(B | not A) \* P(not A)

Using the data we collected from your forthcoming classmates, we can estimate *all* these probs.

n	Passed?	Нарру?	GPA?	Friends'?	Grade			Passed	= Yes	Passed = No	
	Yes	No	3.8	А	А			T GSSCG	105		
	No	No	2.1	D	F	Α		1		0	
	Yes	Yes	3.1	В	В	N	ot A	2		2	
	No	No	3.3	А	F					_	
	Yes	Yes	3.2	В	В		CDA	CDA	CDA	CDA	
									GPA 2 5 2	GPA	
9 (you!)	Yes	Yes	3.9	В	????		> 3.5	3.0-3.5	2.5-3.	0 1.5-2.5	
						Α	1	0	0	0	

Not A

0

0

3

1

#### Your turn!

That's a lot of potentially new math.

Check your understanding at the lowest level by filling out the worksheet by hand.

If that's straightforward, try the deeper thought questions....

#### Naïve Bayes Modeling Example: Estimating Your Hidden Grade From More Forthcoming Classmates' Statements

You're shy and don't want to say whether you got an A. But many of your classmates are pretty forthcoming. Can we figure your grade out from what they say about themselves plus a handful of facts about you – with mathematical rigor?

Person	Passed the class?	Looked happy seeing	GPA estimate?	Average grade of friends?	True grade
	chu351	grades?		or menus	
1	Yes	No	3.8	A	A
2	No	No	2.1	D	F
3	Yes	Yes	3.1	В	В
4	Yes	No	3.3	A	D
5	Yes	Yes	3.6	В	В
6	Yes	No	4.0	A	A
7	Yes	Yes	3.7	A	A
8	Yes	Yes	3.2	A	В
9	Yes	Yes	3.1	В	В
10	No	No	1.8	D	F
11	Yes	No	3.3	В	A
12	Yes	No	3.6	A	В
13	Yes	Yes	2.7	В	С
14	No	No	3.1	С	F
15	Yes	Yes	1.9	С	D
99 (YOU)	Yes	Yes	3.9	В	2222

To get a mathematically justifiable prediction of your grade, we tally up the counts:

Passed	Yes	No	Ha	арру	Yes	No
A	4	0	A			
Not A	8	3	Ne	ot A		

GPA	>3.5	3.0-3.5	2.5-3.0	1.5-2.5	<1.5
A					
Not A					

Friends'	A	В	с	D	F
А					
Not A					

Then we calculate two Naïve Bayes scores for person #99 (you!), whose grade we want to predict:

	Passed * Happy * GP	'A * Friends' * P(A)
Setting an A:	P(Yes   A) * P(Yes   A) * P(>3.5	5   A) * P(B   A) * P(A)
	4/4 ≓ 1.0 * *	• • =

Not getting an A:

Given the observed data and our Naïve Bayes model, which target class (A or not A) is most likely?

Pamela Toman (ptoman@cs.stanford.edu) SAIL ON, October 2016



## Deeper thought questions for you...

Why did we "discretize" (put into buckets) the GPA variable? How do we avoid counting a border case like 3.5 in more than one bucket?

If we need to bucket variables (like GPA), what are some good ways of choosing the bucket size?

Does the Naïve Bayes method get better with a bigger sample size (more data)? Why?

What do we need to estimate the actual probability that you got an A?

Can we use this framework to figure out your most likely grade (A-F)? How?

What happens if we want to make a prediction about someone who has a 1.4 GPA but we don't have anyone in the dataset with that characteristic? How could we fix that problem?

Suppose the teacher tells us P(A) is really 45% even though our data estimate says 27% (our sample was skewed). Which value should we use in our equation for P(A)? Why?

#### Naïve Bayes Modeling Example: Estimating Your Hidden Grade From More Forthcoming Classmates' Statements

#### Your turn! (answers)

Our prediction, based on the 15 rows of data, is that you **don't** get an A. \*sad trombones\*

The score for getting an A: 0.0127

The score for not getting an A: 0.0186

You're shy and don't want to say whether you got an A. But many of your classmates are pretty forthcoming. Can we figure your grade out from what they say about themselves plus a handful of facts about you – with mathematical rigor?

Person	Passed the	Looked happy	GPA	Average grade	True grade
	class?	seeing	estimate?	of friends?	
		grades?			
1	Yes	No	3.8	A	A
2	No	No	2.1	D	F
3	Yes	Yes	3.1	В	В
4	Yes	No	3.3	A	D
5	Yes	Yes	3.6	В	В
6	Yes	No	4.0	A	A
7	Yes	Yes	3.7	A	A
8	Yes	Yes	3.2	A	В
9	Yes	Yes	3.1	В	В
10	No	No	1.8	D	F
11	Yes	No	3.3	В	A
12	Yes	No	3.6	A	В
13	Yes	Yes	2.7	В	С
14	No	No	3.1	С	F
15	Yes	Yes	1.9	с	D
99 (YOU)	Yes	Yes	3.9	В	????

To get a mathematically justifiable prediction of your grade, we tally up the counts:

Passed	Yes	No	Hap	ру	Yes	No
A	4	0	A			
Not A	8	3	Not.	A		

GPA	>3.5	3.0-3.5	2.5-3.0	1.5-2.5	<1.5
A					
Not A					

Friends'	Α	В	С	D	F
A					
Not A					

Then we calculate two Naïve Bayes scores for person #99 (you!), whose grade we want to predict:

	Passed * Haj	opy * GPA	* Friends' * P(A)	
Getting an A:	P(Yes   A) * P(Yes	s   A) * P(>3.5   A	) * P(B   A) * P(A)	
	4/4 ≠ 1.0 *	•	• •	=

Not getting an A:

Given the observed data and our Naïve Bayes model, which target class (A or not A) is most likely?

Pamela Toman (ptoman@cs.stanford.edu) SAIL ON, October 2016

All models connect input variables to a target output variable through math/algorithms

- The Naïve Bayes model goes through each target class individually, and then it estimates the probability that the data we observed came from that target class. It predicts the most likely target class.
- Other models follow other assumptions & use other approaches

All models connect input variables to a target output variable through math/algorithms

- The Naïve Bayes model goes through each target class individually, and then it estimates the probability that the data we observed came from that target class. It predicts the most likely target class.
- Other models follow other assumptions & use other approaches

The best & most commonly used models are *mathematically justifiable* 

- We like rigor
- We like having a guarantee that the solution is correct (or optimal) under some assumptions
- No one takes you seriously if you don't have a defensible reason to do what you did

All models connect input variables to a target output variable through math/algorithms

- The Naïve Bayes model goes through each target class individually, and then it estimates the probability that the data we observed came from that target class. It predicts the most likely target class.
- Other models follow other assumptions & use other approaches

The best & most commonly used models are *mathematically justifiable* 

- We like rigor
- We like having a guarantee that the solution is correct (or optimal) under some assumptions
- No one takes you seriously if you don't have a defensible reason to do what you did

#### What makes this *machine learning*?

We defined a **general** approach for using conditional probabilities to get the most likely target class. We only used & programmed generic logic. The same approach works anywhere – with predicting cancer, with predicting whether an email is something you'll care about, etc. ...!

#### Let's turn to the Titanic

The Titanic data is much more than 15 rows (yay!)

But hand counting was painful on only 15 rows.

Now that we thoroughly understand the math & algorithm... Let's use a computer!

# Opening up the Jupyter Notebook

Quickstart from scratch:

- 1. Install the <u>Anaconda</u> distribution of Python (other distributions won't necessarily come with all that we need)
- 2. Download the .ipynb file to /some/path/for/my/ipynb
  - Jupyter notebooks (\*.ipynb) are for "literate programming" intermingling code and detailed text discussions (even pictures!)
- 3. From Kaggle's website, download the Titanic train.csv and test.csv files to /some/path/for/my/ipynb
- 4. Start Jupyter Notebook in a parent folder to where the .ipynb file is
  - Navigate to /some/path/for/my/ipnb in the command line & run `jupyter notebook` in the command line

-OR-

- Open up Jupyter Notebook & move the .ipynb & .csv files so they show up in the file list
- 5. Click on the *SAILON\_Titanic.ipynb* file within your web browser

Hama		-	o x
- Home - T			
ilocalhost:8889/tree	C Q Search	☆ 自 ♥ ↓	<b>⋒</b> =
📁 jupyter			
Files Running Clusters Conda			
Select items to perform actions on them.		Upload New - 2	
- * *			
SAILON_Titanic.ipynb			
201610_NaiveBayes.pptx			
Naïve Bayes Example.docx			
SAILON_Titanic.zip			
test.csv			
Titanic_NaiveBayes.pptx			
train.csv			



Bradley

tutorial on Jupyter Notebooks.

ook 🖹

When coding, we almost always use libraries for the math (which we need to understand to not make silly mistakes & get stuck). Starting to code raises new challenges: thinking on & usefully preparing lots of data.

For even more info, there is a Youtube <u>tutorial on Jupyter Notebooks</u>.

			_	- 🗆	$\times$
SAILON_Titanic	× +				
i localhost:8889/note	books/SAILON_Titanic.i C Q Search	☆自	♥	↓ ⋒	≡
💭 jupyter	SAILON_Titanic (autosaved)			ę	
File Edit View	Insert Cell Kernel Widgets Help	Pytho	n [conda i	root] O	
B + % 4 B	↑     ↓     ▶     ■     C     Markdown     □     □     CellToolbar     ▲	ü D			
					~

#### Kaggle: Titanic

We're finally here: Let's predict who lives and who dies on the Titanic given the passenger manifest records.

#### Getting and examining the data

The first step is always reading in data and looking at it to get a sense for what it contains.

Here's one way to do that in Python:

In [10]: **from** IPython.display **import** display, HTML # this lets us get pretty tabular **import** pandas **as** pd # this is a useful library for working with data

> all\_data = pd.read\_csv("train.csv", index\_col=0) # assumes train.csv is in display(all\_data[0:5])

	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	С
Passengerld										Γ
1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	N
			Cumings, Mrs. John Bradley							

#### To get an even better handle on this...

1) Try to answer the thought questions posed earlier

2) Tinker with & fully understand the Jupyter Notebook we worked on tonight

- 3) Build an amazing Naïve Bayes model for the Titanic
- Create a Kaggle account (<u>https://www.kaggle.com</u>)
- Review the <u>Titanic data and tutorials</u> (look for DIY Tutorials and Kaggle Kernels look especially for ones that name their method)
- Build a Naïve Bayes model that predicts who will survive (starting from our notebook or elsewhere)
- Explore excluding, including, and deriving new columns to try to improve performance
- Keep your best models & predictions we'll share insights later on!

4) Write your own implementation of Naïve Bayes in Python (from scratch!) and/or write out the derivation of the Naïve Bayes model (without peaking!) – this is the way to verify you have a deep & thorough understanding