

Naïve Bayes

Noa Lubin & Lior Sidi





Agenda

- Motivation
- Theory
- Practical Example
- Naive Bayes Language Model
- Code
- Summary

Motivation



Guess what's in the box dogs or cats?

Problem: Guess what's in the box dogs or cats?

we **cannot open the box.**

we can only **weight the box.**

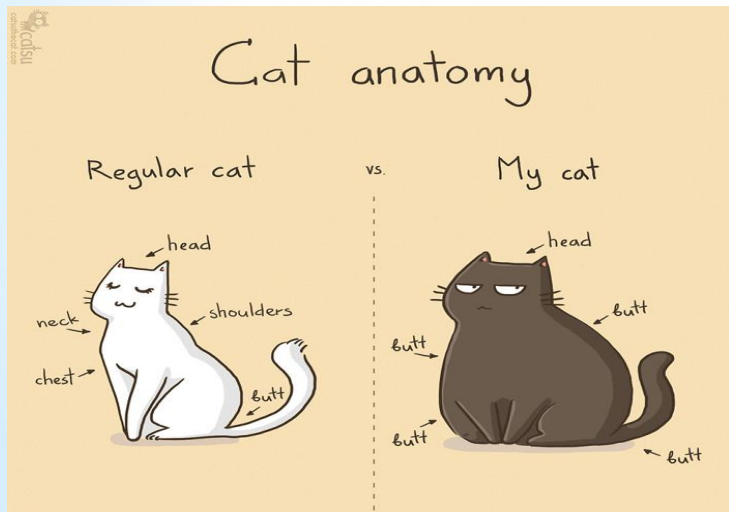


We Know That

Dogs - 471M (56%)

Cats - 373M (44%)

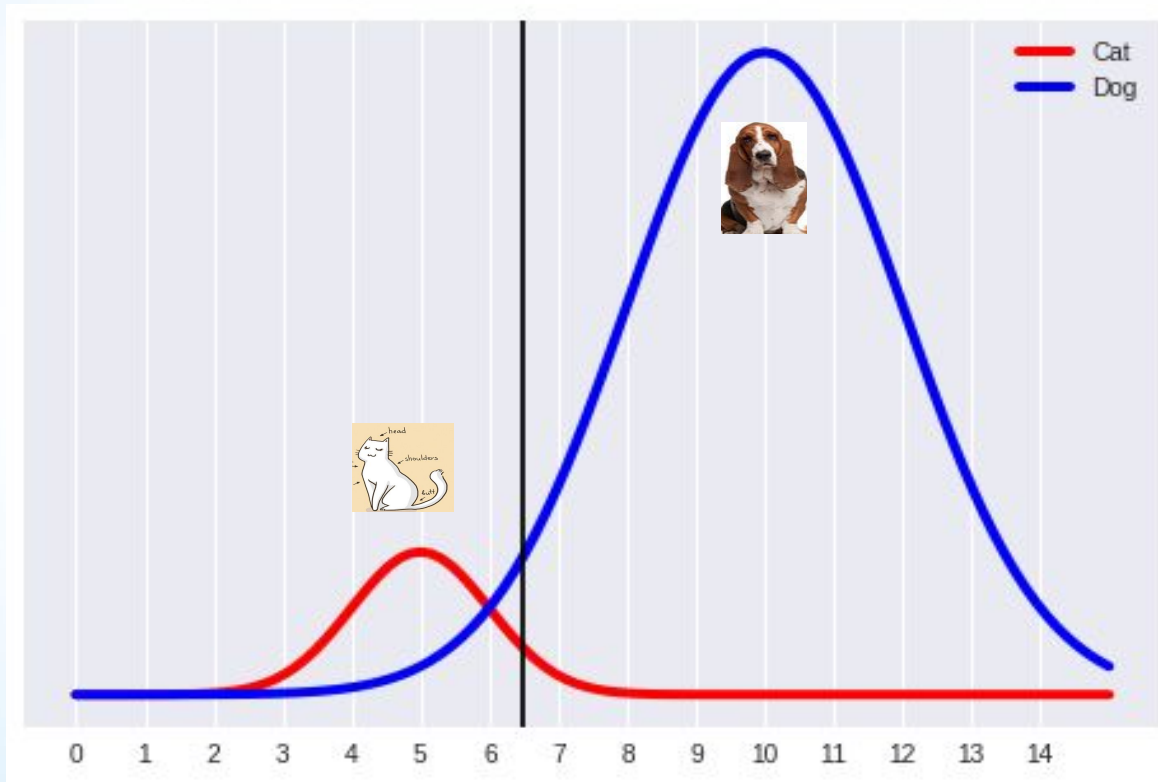
$$p(\text{weight}|\text{Cat}) \sim N(5,1)$$



$$p(\text{weight}|\text{Dog}) \sim N(10,4)$$



A box weighs 6.5kg, is it a Dog or a Cat?



Theory



Reminder:

Independence assumption

Independent events:

E and F are *independent* if $P(EF) = P(E)P(F) \Leftrightarrow P(E|F) = P(E)$

Reminder:

Chain Rule

Conditional probability:

$$P(E|F) = \frac{P(EF)}{P(F)} \quad \text{or} \quad P(EF) = P(E|F)P(F)$$

Chaining rule:

$$P\left(\bigcap_{i=1}^n A_i\right) = P(A_1 A_2 \dots A_n) = P(A_1)P(A_2|A_1)P(A_3|A_1 A_2) \dots P(A_n|A_1 A_2 \dots A_{n-1})$$



Reminder:

Bayes Rule



Bayes' formula:

$$P(F|E) = \frac{P(E|F)P(F)}{P(E)}$$

Estimation

We are interested in:

- Classification - Ex: Decide the topic of a sentence.
- Regression - Ex: Weather
- Generation = Ex: Generate text

We'll look at the generative behavior that causes something.

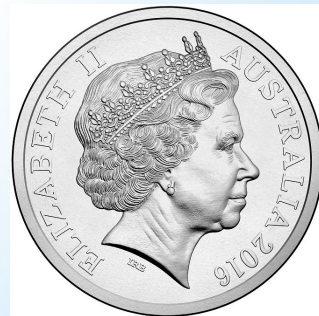
What we formally call: maximum likelihood estimation

Intuition Behind MLE

Let say we have a coin with bernoulli distribution for heads = ?

We can toss the coin 100 times. Let's say that it landed on heads 68 times.

What do you think the probability for heads is?



Intuition Behind MLE

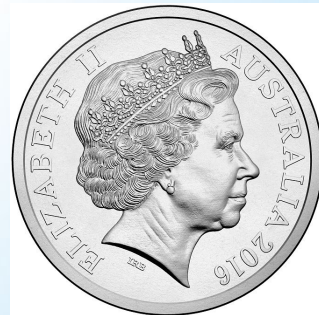
Let say we have a coin with bernoulli distribution for heads = ?

We can toss the coin 100 times. Let's say that it landed on heads 68 times.

What do you think the probability for heads is?

0.68!

You just used the MLE principal



MLE - Maximum Likelihood Estimation

We can also prove that MLE converges to the true parameters when the observation number goes to infinity

Think about the coin example.

If you toss it $n=3$ times and get 2 heads: $p_{MLE} = 0.66$

If you make $n=1000$ and get 705 heads: $p_{MLE} = 0.705$

Let's Go Back to Bayesian Law

Probability of *Dog*
given that $x=6.5$

Probability of
 $x=6.5$ given that y
is Dog

Probability
of y being
Dog

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

Probability of $x=6.5$



Let's Go Back to Bayesian Law

Probability of y
being "True" given
that x is "True"

Probability of x
being "True" given
that y is "True"

Probability
of y being
"True"

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

Probability of x
being "True"



Let's Go Back to Bayesian Law

Posterior

Likelihood

Prior

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

Evidence



MLE - Maximum Likelihood Estimation

If I know the outcome, what is the probability of a certain observation?

Ex: for Dogs and Cats: $P(\text{Dogs}) = \# \text{ dogs} / (\# \text{ dogs} + \# \text{ cats})$

$P(\text{Weight}|\text{Dog})$ = We make a gaussian assumption and Mu and Sigma are calculated based on all observed dogs in the world

Max A-Posterior (MAP) classifier

$$\hat{y}_{MAP} = \operatorname{argmax}_y P(y|x)$$



$$= \operatorname{argmax}_y \frac{P(x|y)P(y)}{P(x)}$$

$$= \operatorname{argmax}_y P(x|y)P(y)$$

Positive const
Not depend on y

likelihood

Prior

Formulate Dog vs Cat Problem with MAP

Y - category {Dog, Cat}

x - weight

Which is larger:

$P(y=\text{Dog}|x=6.5)$ or $P(y=\text{Cat}|x=6.5)$?

Max A-Posterior (MAP) classifier

$$\begin{aligned}\hat{y}_{MAP} &= \operatorname{argmax}_y P(y|x) \\ &\stackrel{\text{Bayes}}{=} \operatorname{argmax}_y \frac{P(x|y)P(y)}{P(x)} \\ &= \operatorname{argmax}_y P(x|y)P(y)\end{aligned}$$

Probability of
weight x given
category y {Cat,
Dog}

Probability of
category y {Cat,
Dog}

What is the most
probable category
 y {Cat, Dog} given
weight x ?

Formulate Dog vs Cat Problem with MAP

We'll transform the problem with Bayesian Rule to observed data:

Which is larger:

$p(x=6.5|y=\text{Dog}) * p(y=\text{Dog})$ or $p(x=6.5|y=\text{Cat}) * p(y=\text{Cat})$?

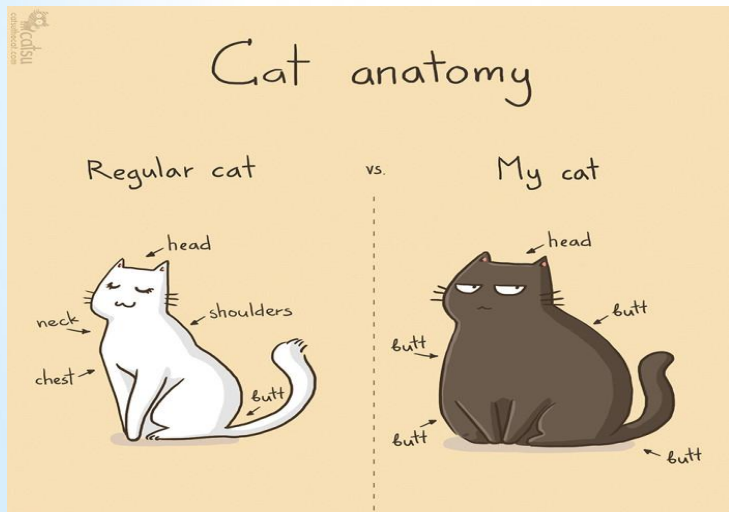
Reminder: We Know That

Dogs - 471M (56%)

Cats - 373M (44%)

$p(\text{weight}|\text{Cat}) \sim N(5,1)$

$p(\text{weight}|\text{Dog}) \sim N(10,4)$



Gaussian Naive Bayes with MLE

$$P(x_i | y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}\right)$$

$$\Pr(y) = \frac{N_y}{N}$$

$$p(x=6.5|y=\text{dog}) = 0.68$$

$$p(y=\text{dog}) = 0.56$$

$$p(x=6.5|y=\text{cat}) = 0.13$$

$$p(y=\text{cat}) = 0.44$$

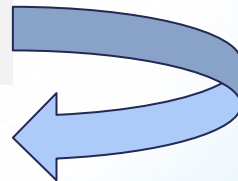
$0.68 * 0.56 > 0.13 * 0.44 \rightarrow$ **It's a dog!**



MLE is a specific case of MAP

In the special case when prior follows a uniform distribution, MAP can be written as:

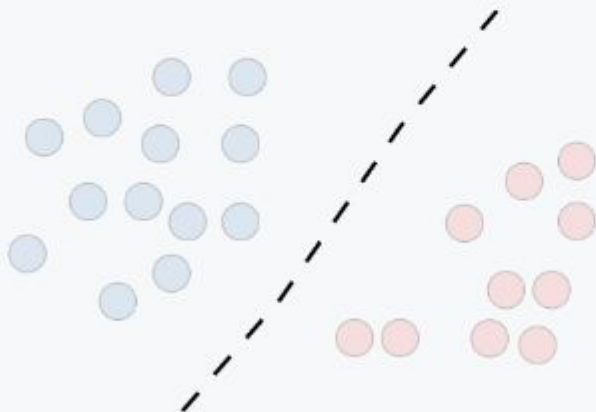
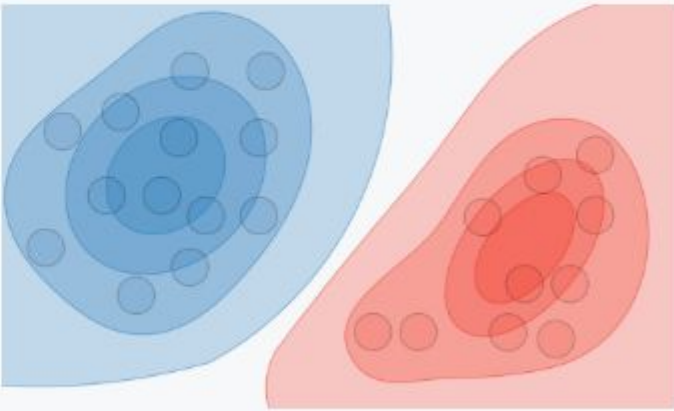
$$\begin{aligned}\hat{y}_{MAP} &= \operatorname{argmax}_y P(y|x) \\ &= \operatorname{argmax}_y \frac{P(x|y)P(y)}{P(x)} \\ &= \operatorname{argmax}_y P(x|y)P(y) \\ &= \operatorname{argmax}_y P(x|y).\end{aligned}$$



Types of Naive Bayes Classifiers

- ***Gaussian Naive Bayes*** - Used when we are dealing with continuous data and uses Gaussian distribution.
- ***Bernoulli Naive Bayes*** - Used for discrete data, where features are only in binary form.
- ***Multinomial Naive Bayes*** - Widely used classifier for document classification which keeps the count of frequent words present in the documents.
- ***Complement Naive Bayes*** - Used for imbalanced data

Generative vs Discriminative Models

	Discriminative model	Generative model
Goal	Directly estimate $P(y x)$	Estimate $P(x y)$ to then deduce $P(y x)$
What's learned	Decision boundary	Probability distributions of the data
Illustration		
Examples	Regressions, SVMs	GDA, Naive Bayes

Practical Example



Boy or Girl?

Pregnant woman at week 20. Task: Boy or Girl ?

Measurements:

1. Weight gain (Kg)
2. Avg amount of chocolate eaten weekly (grams)



Collect training data

Weight Gain	Chocolate Craving	Gender
10	4	Boy
7	5	Boy
13	1	Boy
11	1	Boy
10	6	Boy
9	0	Boy
6	5	Boy
9	3	Boy
15	5	Boy
1	1	Boy
4	0	Girl
5	1	Girl
10	3	Girl
7	0	Girl
3	1	Girl
7	1.5	Girl
5	1	Girl
3	0	Girl
5	0	Girl
3	0.5	Girl

Data Processing: Discretization

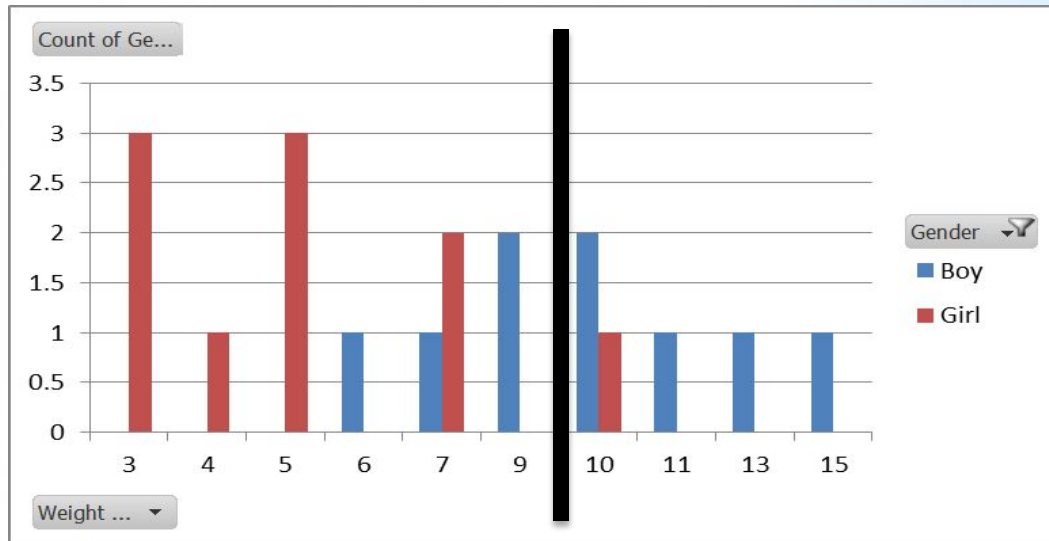
Pregnant woman at week 20.

Task: Boy or Girl ? : Gender (G)



Weight gain: Weight Gain (W)

1. Low if less than 9Kg
2. High if more than 9kg



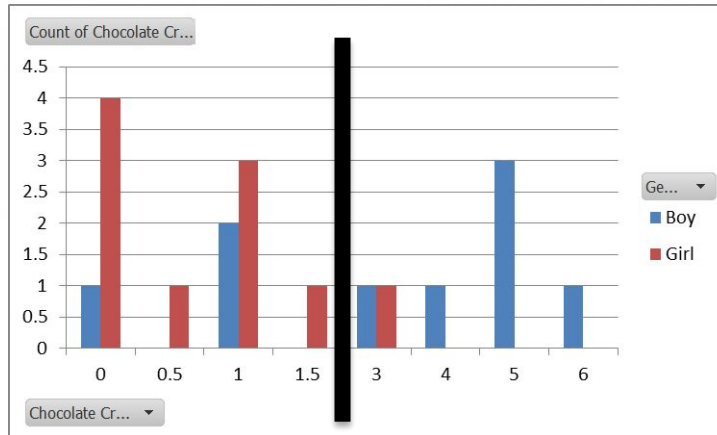
Data Processing: Discretization

Pregnant woman at week 20

Task: Boy or Girl ? : Gender (G)

Avg amount of chocolate eaten weekly
(bars): Chocolate Craving (CC)

1. Low if less than 2 bars
2. High if more than 2 bars



We Have Observed Labeled Data

[illegible]

Let's Predict The Outcome

Pregnant Woman:

Weight Gain: 10 kg \rightarrow 1

Chocolate Craving: 1 bar \rightarrow 0

Argmax gender $p(\text{gender} | w=1, cc=0)$?



Mathematical Formulation

$$\text{Argmax}_G p(G | W, CC)$$



$$= \text{Argmax}_G p(W, CC | G) * P(G) / P(W, CC)$$

$$= \text{Argmax}_G p(W, CC | G) * P(G)$$

* Assuming
independence

$$= \text{Argmax}_G p(W|G) * P(CC | G) * P(G)$$

Likelihood

	W = 0	W = 1
G = Boy	0.3	0.7
G = Girl	0.9	0.1

	CC = 0	CC = 1
G = Boy	0.4	0.6
G = Girl	0.9	0.1

G = Boy	0.5
G = Girl	0.5

[illegible]

Likelihood

	W = 0	W = 1
G = Boy	0.3	0.7
G = Girl	0.9	0.1

	CC = 0	CC = 1
G = Boy	0.4	0.6
G = Girl	0.9	0.1

G = Boy	0.5
G = Girl	0.5

Let's Go Back to Our Woman: $w=1$, $cc=0$

$\text{Argmax}_G P(W|G) * P(CC | G) * P(G)$

Boy \rightarrow

$$P(W=1|G=\text{boy}) * P(CC = 0 | G=\text{boy}) * P(G=\text{boy}) = 0.7 * 0.4 * 0.5 = 0.14$$

Girl \rightarrow

$$P(W=1|G=\text{girl}) * P(CC = 0 | G=\text{girl}) * P(G=\text{girl}) = 0.1 * 0.9 * 0.5 = 0.045$$

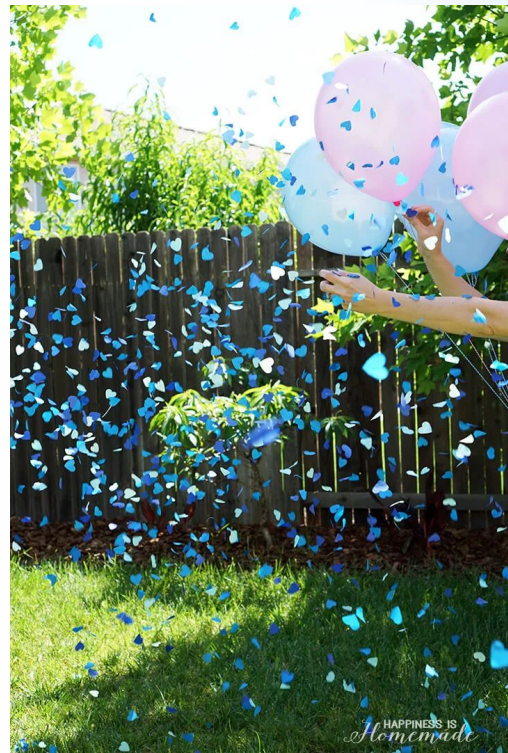
$0.14 > 0.045 \rightarrow$ **It's a boy!**

Likelihood

	W = 0	W = 1
G = Boy	0.3	0.7
G = Girl	0.9	0.1

	CC = 0	CC = 1
G = Boy	0.4	0.6
G = Girl	0.9	0.1

G = Boy	0.5
G = Girl	0.5



Likelihood

	W = 0	W = 1
G = Boy	0.3	0.7
G = Girl	0.9	0.1

	CC = 0	CC = 1
G = Boy	0.4	0.6
G = Girl	0.9	0.1

Let's Go Back to Our Woman: $w=1$, $cc=0$

This is not $P(\text{boy}|W,CC)$!!! What do

Boy → **we need to add to calculate it?**

$$P(W=1|G=\text{boy}) * P(CC = 0 | G=\text{boy}) * P(G=\text{boy}) = 0.7 * 0.4 * 0.5 = 0.14$$

Girl →

$$P(W=1|G=\text{girl}) * P(CC = 0 | G=\text{girl}) * P(G=\text{girl}) = 0.1 * 0.9 * 0.5 = 0.045$$

$0.14 > 0.045 \rightarrow$ **It's a boy!**

Since Probabilities are equal

Are W and CC Independent? Or Could We be Naive?

- Conditional independence:

$$P(W, CC | G) = P(W|G) * P(CC|G)$$

W and CC are independent given G

Each feature has its own conditional probability table:

Question: $P(W=0|G=Girl) * P(CC=0|G=Girl) = ?$ $P(W=0, CC=0|G=Girl)$

Are W and CC Independent?

	W = 0	W = 1
G = Boy	0.3	0.7
G = Girl	0.9	0.1

	CC = 0	CC = 1
G = Boy	0.4	0.6
G = Girl	0.9	0.1

	W=0, CC=0	W=0, CC=1	W=1, CC=0	W=1, CC=1
G= Boy	0.1	0.2	0.3	0.4
G= Girl	0.9	0	0	0.1

$$P(W=0|G=Girl) * P(CC=0|G=Girl) \neq P(W=0, CC=0|G=Girl)$$

[illegible]

Are W and CC Independent?

	W = 0	W = 1
G = Boy	0.3	0.7
G = Girl	0.9	0.1

	CC = 0	CC = 1
G = Boy	0.4	0.6
G = Girl	0.9	0.1

	W=0, CC=0	W=0, CC=1	W=1, CC=0	W=1, CC=1
G= Boy	0.1	0.2	0.3	0.4
G= Girl	0.9	0	0	0.1

Weight Gain	Chocolate Craving	Gender
1	1	Boy
0	1	Boy
1	0	Boy
1	0	Boy
1	1	Boy
1	0	Boy
0	1	Boy
1	1	Boy
1	1	Boy
0	0	Boy
0	0	Girl
0	0	Girl
1	1	Girl
0	0	Girl
0	0	Girl
0	0	Girl
0	0	Girl
0	0	Girl
0	0	Girl
0	0	Girl
0	0	Girl
0	0	Girl

$0.9 * 0.9 \neq 0.9 \rightarrow$ **Weight and Chocolate Consumption are Dependant**

Let's Repeat Without Naive Assumption

	W=0, CC=0	W=0, CC=1	W=1, CC=0	W=1, CC=1
G= Boy	0.1	0.2	0.3	0.4
G= Girl	0.9	0	0	0.1

Let's Go Back to Our Woman: $w=1, cc=0$

$\text{Argmax}_G P(W, CC | G) * P(G)$

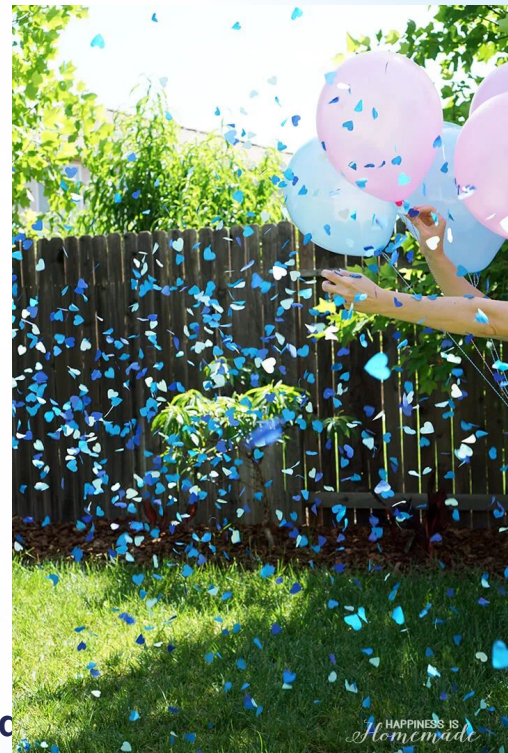
Boy \rightarrow

$$P(W=1, CC=0 | G=\text{boy}) * P(G=\text{boy}) = 0.3 * 0.5 = 0.15$$

Girl \rightarrow

$$P(W=1, CC=0 | G=\text{girl}) * P(G=\text{girl}) = 0 * 0.5 = 0$$

$0.15 > 0 \rightarrow$ **Still a boy (but could have changed)**



Naive Bayes Language Model

Based on: CIS 391 – Introduction to Artificial Intelligence



Example: Classifying Spam Mails

- Desired output '**Label**' an article to one of three categories (multiclass)
 - Spam
 - Family
 - Work



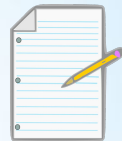
Data

- Collection of 1.5M emails (documents) where
 - 500K labeled as spam
 - 500K labeled as family
 - 500K labeled as work

Hi,
How are you?
You must buy the new nike
shoes
Here is a link:
<http://www.nike.com/bla>

Multinomial Naive Bayes for Text Classification

- $P(w_i | c)$ is the conditional probability of word w_i occurring in document of class c
- $P(c)$ is the prior probability of a document occurring in class c
- n_d is the number of such tokens in d
- $P(c|d)$ is the probability of class given document



The model: $\text{Argmax}_c P(c|d) = \text{Argmax}_c p(d|c) * p(c) \mid p(d) = \text{Argmax}_c p(c) * p(d|c)$

$$p(d|c) = p(w_1, w_2, w_3, w_4, \dots, w_n | c)$$

$$= p(w_1|c) * p(w_2|w_1, c) * p(w_3|w_1, w_2, c) * \dots * p(w_n|w_1, \dots, w_{n-1}, c)$$

Assuming independence = $p(w_1|c) * p(w_2|c) * p(w_3|c) * \dots * p(w_n|c)$

$$c_{NB} = \underset{c_j \in C}{\text{argmax}} P(c_j) \prod_{i \in \text{positions}} P(w_i | c_j)$$

Multinomial Naive Bayes for Text Classification

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(w_i | c_j)$$

Instead of multiplication of probabilities, use **sum of logs** to avoid underflow

How will we calculate the parameters?

- From training corpus, extract *Vocabulary*
- Calculate required $P(c_j)$ and $P(w_k | c_j)$ terms
 - For each c_j in C do :

$$\hat{P}(c_j) = \frac{\text{doccount}(C = c_j)}{N_{doc}}$$

In our case?

How will we calculate the parameters?

- From training corpus, extract *Vocabulary*
- Calculate required $P(c_j)$ and $P(w_k | c_j)$ terms
 - For each c_j in C do :

$$\hat{P}(c_j) = \frac{\text{doccount}(C = c_j)}{N_{doc}}$$

In our case: $P(\text{spam}) = P(\text{family}) = p(\text{work}) = 1/3$

How will we calculate the parameters?

- For each word w_k in *Vocabulary*

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

What happens if a word didn't appear in any document?

Discounting - Handling Missing Data

Example: Lidstone discounting

$$p_{Lid}(x) = \frac{C(x) + \lambda}{|S| + \lambda|X|}$$

Vocabulary size

For $\lambda=1$ it's called laplace discounting

Bonus: We can prove that $p_{Lid}(x) > p_{MLE}(x)$

How will we calculate the parameters?

- For each word w_k in *Vocabulary*

$$\begin{aligned}\hat{P}(w_i | c) &= \frac{\text{count}(w_i, c) + 1}{\sum_{w \in V} (\text{count}(w, c) + 1)} \\ &= \frac{\text{count}(w_i, c) + 1}{\left(\sum_{w \in V} \text{count}(w, c) \right) + |V|}\end{aligned}$$

Naive Bayes as a Generative Model

Notice that we can use the model we build to classify if an email is spam or not in order to build a generative model! This model is called “Language Model”

How?

$$P(c_j)$$

We pick a class by tossing a bernoulli coin with

Let's say we got $c=\text{spam}$

Generative Story

Now we generate the text accordingly.

We pick the first word by sampling $p(w|c=\text{spam})$

Code



Let's Think About This Together

1. What are the main hyper-parameters?
2. Can it work for Multi-class data?
3. Does it handle Categorical data?
4. Does it handle missing data?
5. Is it sensitive to outliers?
6. What if some features are correlated?
7. Is it Interpretable?
8. Can it be parallelized?
9. Speed of training
10. Speed of prediction

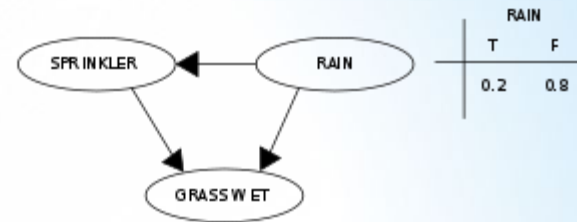
Let's Think About This Together

1. What are the main hyper-parameters? Distribution type, smoothing parameter
2. Can it work for Multi-class data? Yes
3. Does it handle Categorical data? Yes
4. Does it handle missing data? Yes
5. Is it sensitive to outliers? No
6. What if some features are correlated? Naive assumptions don't work
7. Is it Interpretable? Yes
8. Can it be parallelized? Yes
9. Speed of training - Fast (Linear Time)
10. Speed of prediction - Fast (Linear Time)

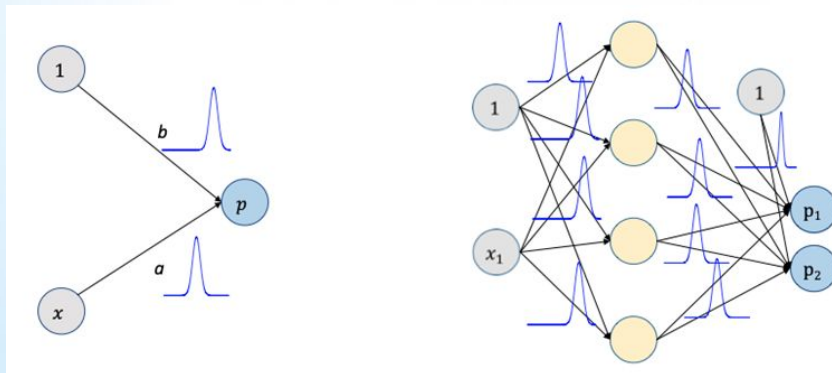
Non-Naive Bayesian Methods

- Bayesian Networks
- Deep Bayesian Networks

RAIN	SPRINKLER	
	T	F
F	0.4	0.6
T	0.01	0.99



SPRINKLER	RAIN	GRASS WET	
		T	F
F	F	0.0	1.0
F	T	0.8	0.2
T	F	0.9	0.1
T	T	0.99	0.01

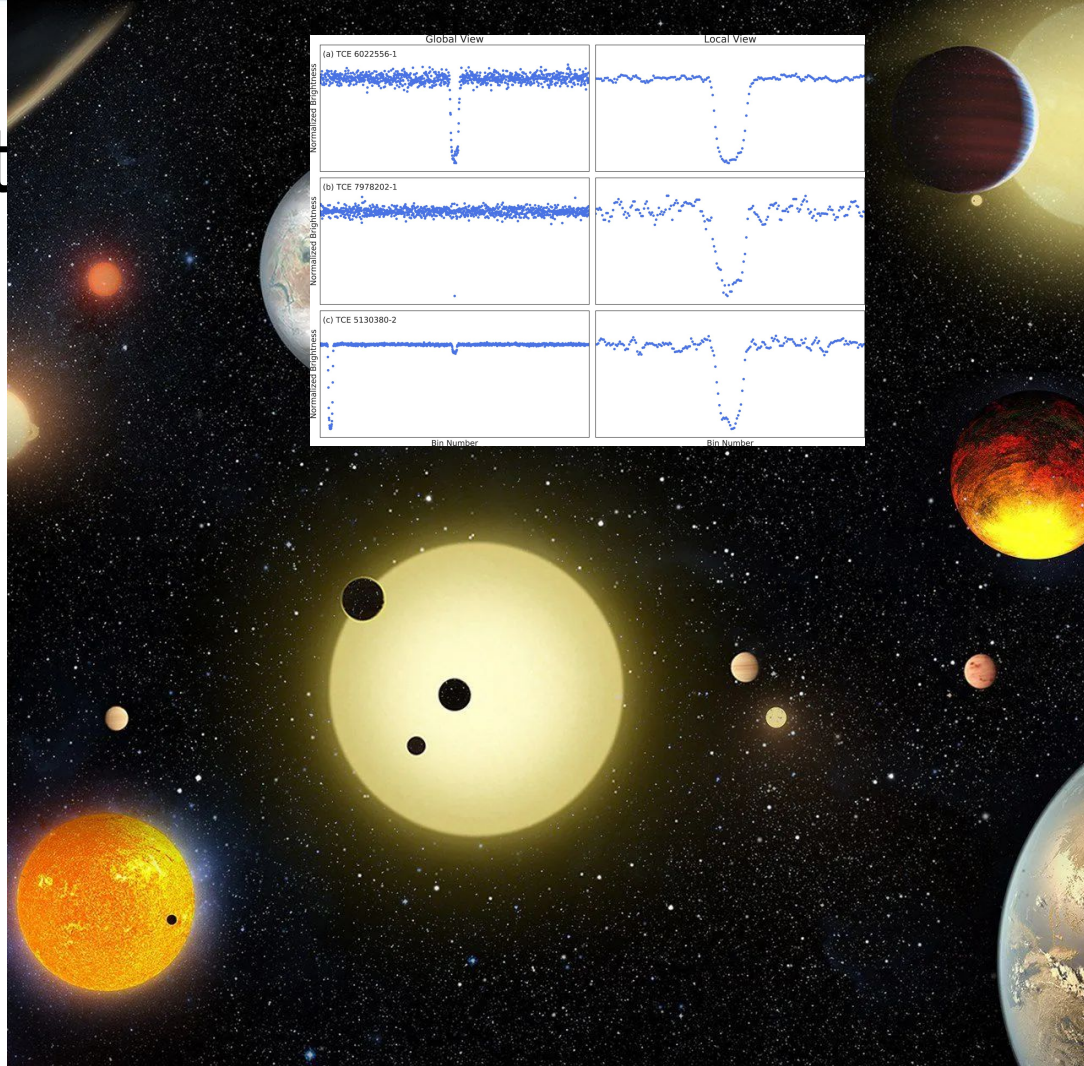
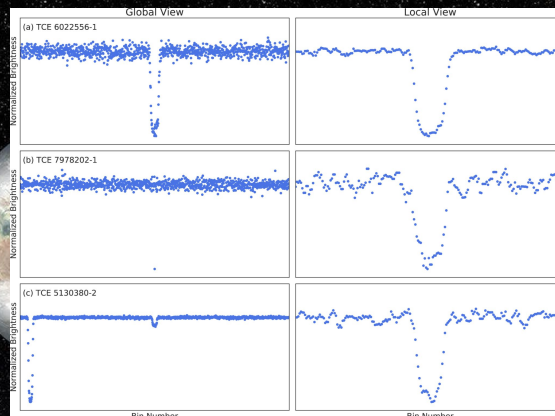


VESPA Exoplanet

The generative category includes models such as **vespa** (Morton & Johnson 2011; Morton 2012; Morton et al. 2016) that require the class priors $p(y = 1)$ and $p(y = 0)$ and likelihoods $P(X|y = 1)$ and $P(X|y = 0)$ to estimate the posterior $P(y = 1|X)$ according to the Bayes' theorem:

$$p(y = 1|X) = \frac{p(X|y = 1)p(y = 1)}{\sum_y p(X|y)p(y)} \quad (1)$$

where $y = 1$ represents an exoplanet, $y = 0$ represents a false positive, and X is a representation of the transit signal. Such generative approaches require the detailed knowledge of the likelihood, $P(X|y)$, and prior, $P(y)$, for each class (exoplanet vs false positive), and class scenario (e.g., BEB). While in general both the likelihood and priors can be learned in a data-driven approach (using ML), **vespa** estimates them by sim-



Summary



Naive Bayes

- Is a generative probabilistic model (the only one we'll see in SL class)
- Can be used for classification, regression and generation
- Uses MLE and MAP principals

Pros and Cons

Pros	Cons
<ol style="list-style-type: none">1. Very fast2. Good for big data with big velocity3. Ignores interactions and therefore needs less data4. Works well with multiclass problems5. Works with missing data6. Can be used as a generative model (e.g. generate text)	<ol style="list-style-type: none">1. Small dataset leads to instability (probabilities can be 0 or 1 with high variance)2. Doesn't perform well for imbalanced datasets3. Continuous features require binning or assumption of a distribution