# INTRO TO DATA SCIENCE LESSON 7: NAIVE BAYES

## LOGISTIC REGRESSION

## QUESTIONS?

**I. NAIVE BAYES**
**II. LAB: IMPLEMENT NAIVE BAYES IN PYTHON**
**III. LAB: USE NAIVE BAYES IN SKLEARN**

# I. BAYESIAN INFERENCE

**Bayes' theorem**. *Here it is again:*

$$P(A|B) = P(B|A) * P(A) / P(B)$$

*Some facts:*
*- This is a simple algebraic relationship using elementary definitions.*
*- It's interesting because it's kind of a "wormhole" between two different "interpretations" of probability.*
*- It's a very powerful computational tool.*

*This term is the **likelihood function**. It represents the joint probability of observing features $\{x_i\}$ given that that record belongs to class C.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*This term is the **prior probability** of c. It represents the probability of a record belonging to class c before the data is taken into account.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*This term is the **normalization constant.** It doesn't depend on C, and is generally ignored until the end of the computation.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*This term is the **posterior probability** of c. It represents the probability of a record belonging to class c after the data is taken into account.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*This term is the **posterior probability** of c. It represents the probability of a record belonging to class c after the data is taken into account.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$
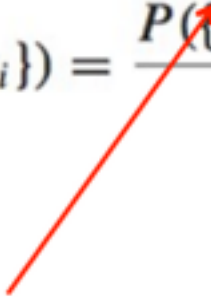
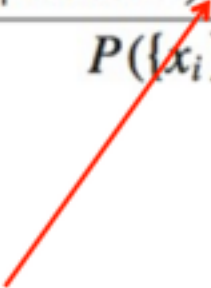*The goal of any Bayesian computation is to find ("learn") the posterior distribution of a particular variable.*

Maximum likelihood estimator (MLE):

What parameters **maximize** the likelihood function?

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

Maximum a posteriori estimate (MAP):

What parameters **maximize** the likelihood function **AND** prior?

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

Problem:

We observe the following coin flips:

HTHH

What is $P(X = \text{Heads})$ ?

Problem:

We observe the following coin flips:

HTHH

What is P( X = Heads) ?   3/4, Why?

Problem:

We observe the following coin flips:

HTHHTHT

What is P( X = Heads) ?

Problem:

We observe the following coin flips:

HTHHTHT

What is P( X = Heads) ?   4/7, Why?

We observe the following coin flips:
HTHHTHT

*Maximum likelihood estimator (MLE):*
*What parameters* **maximize** *the likelihood function?*
Let $P(X = \text{Heads}) = q$, and write Bayes Theorem

$P(q \mid \text{observations}) = P(\text{observations} \mid q) * P(q) / \text{constant}$

*Maximum likelihood estimator (MLE):*
*What parameters* **maximize** *the likelihood function?*
Let $P( X = \text{Heads}) = q$, and write Bayes Theorem

$P(q \mid \text{observations}) = P (\text{observations} \mid q) * P (q) / \text{constant}$

$P(\text{observations} \mid q ) = ?$
$P(q) = ?$

*Maximum likelihood estimator (MLE):*
*What parameters **maximize** the likelihood function?*
Let P( X = Heads) = q, and write Bayes Theorem

P(q | observations) = P (observations | q) * P (q) / constant

P(observations | q ) = Binomial Distribution
P(q) = ????

Binomial Distribution:

$$\Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

P ( HTHHTHT | q ) = P ( X = 4, n = 7 ) =
          = (7 choose 4) * q^4 * (1-q) ^ 3

Binomial Distribution:

$$\Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

P ( HTHHTHT | q ) = P ( X = 4, n = 7) =

$\qquad$ = (7 choose 4) * q^4 * (1-q) ^ 3

After optimizing, the **MLE is 4/7**

A prior distribution is known as **conjugate prior** if its from the same family as the posterior for a certain likelihood function

For the binomial distribution, the conjugate prior is the **Beta distribution**

$$= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1 - x)^{\beta-1}$$

$$= \frac{1}{B(\alpha, \beta)} x^{\alpha-1}(1 - x)^{\beta-1}$$

The **MAP estimate** is the value that maximizes both the likelihood function and prior – the product of the two.

In the coin flip setting is the value that optimizes

$$P ( HTHHTHT \mid q ) * P(q)$$

The **MAP estimate** is the value that maximizes both the likelihood function and prior – the product of the two.

In the coin flip setting is the value that optimizes
$P ( HTHHTHT \mid q ) * P(q)$
$= (7 \text{ choose } 4) \, q \char`\^ 4 * (1 - q) \char`\^3 * q\char`\^(a-1) * (1-a) \char`\^(b-1)$

The **MAP estimate** is the value that maximizes both the likelihood function and prior – the product of the two.

In the coin flip setting is the value that optimizes
   $P ( HTHHTHT \mid q ) * P(q)$
   $= (7 \text{ choose } 4) \, q^{\wedge}4 * (1-q)^{\wedge}3 * q^{\wedge}(a-1) * (1-a)^{\wedge}(b-1)$
   $= q^{\wedge}(4+a-1) * (1-q)^{\wedge}(3+b-1)$

The **MAP estimate** is the value that maximizes both the likelihood function and prior – the product of the two.

In the coin flip setting is the value that optimizes

P ( HTHHTHT | q )  * P(q)

= (7 choose 4) q ^ 4 * (1 -q ) ^3 * q^(a−1) * (1−a) ^(b−1)

= q^(4 + a −1) * (1−q)^ ( 3 + b − 1)

After optimizing, the **MAP is (4 + a −1) / ( 7 + a + b − 2 )**

Why do you care?

Why do you care?

Many problems are binary and are estimated using counts...

Why do you care?

Many problems are binary and are estimated using counts...

Ex. 1:
Sample 100 people and ask if they support a politician?

Why do you care?

Many problems are binary and are estimated using counts...

Ex. 1:
Sample 100 people and ask if they support a politician?
23 say Yes – Is the correct prediction 23/100?
What's the prior?

Ex. 2:

Need to choose between multiple categories to present (for ads, products, news).

Ex. 2:

Need to choose between multiple categories to present (for ads, products, news).

You can compute response % for each category

Ex. 2:

Need to choose between multiple categories to present (for ads, products, news).

You can compute response % for each category

But each should have a unique prior – **unique psuedo counts**

*Suppose we have a dataset with features $x_1, ..., x_n$ and a class label $c$. What can we say about classification using Bayes' theorem?*

Suppose we have a dataset with features $x_1, ..., x_n$ and a class label $c$. What can we say about classification using Bayes' theorem?

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

Bayes' theorem can help us to determine the probability of a record belonging to a class, given the data we observe.

The idea of Bayesian inference, then, is to **update** our beliefs about the distribution of $c$ using the data ("evidence") at our disposal.

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

Then we can use the posterior for prediction.

*Q: What piece of the puzzle we've seen so far looks like it could intractably difficult in practice?*

*Remember the likelihood function?*

$$P(\{x_i\}\,|\,C) = P(\{x_1, x_2, \ldots, x_n\})\,|\,C)$$

Remember the likelihood function?

$$P(\{x_i\}\,|\,C) = P(\{x_1, x_2, \ldots, x_n\})\,|\,C)$$

Observing this exactly would require us to have enough data for every possible combination of features to make a reasonable estimate.

Q: What piece of the puzzle we've seen so far looks like it could intractably difficult in practice?

A: Estimating the full likelihood function.

Q: So what can we do about it?

*Q: So what can we do about it?*

*A: Make a simplifying assumption. In particular, we assume that the features $x_i$ are conditionally independent from each other:*

Q: So what can we do about it?

A: Make a simplifying assumption. In particular, we assume that the features $x_i$ are conditionally independent from each other:

$$P(\{x_i\}|C) = P(x_1, x_2, ..., x_n|C) \approx P(x_1|C) * P(x_2|C) * ... * P(x_n|C)$$

Q: So what can we do about it?

A: Make a simplifying assumption. In particular, we assume that the features $x_i$ are conditionally independent from each other:

$$P(\{x_i\}|C) = P(x_1, x_2, ..., x_n|C) \approx P(x_1|C) * P(x_2|C) * ... * P(x_n|C)$$

This "naïve" assumption simplifies the likelihood function to make it tractable.

$$P(\{x_i\}|C) = P(x_1, x_2, ..., x_n|C) \approx P(x_1|C) * P(x_2|C) * ... * P(x_n|C)$$

Q: Given that we can compute this value, what do we do with it?

$$P(\{x_i\}|C) = P(x_1, x_2, ..., x_n|C) \approx P(x_1|C) * P(x_2|C) * ... * P(x_n|C)$$

Q: Given that we can compute this value, what do we do with it?

A: In our training phase, we 'learn' the probability of seeing our training examples under each class.

$$P(\{x_i\}|C) = P(x_1, x_2, ..., x_n|C) \approx P(x_1|C) * P(x_2|C) * ... * P(x_n|C)$$

Q: Given that we can compute this value, what do we do with it?
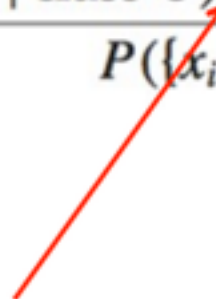
A: In our training phase, we 'learn' the probability of seeing our training examples under each class.

Then we use Bayes Theorem to compute P( class | inputs)

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*Maximum a posteriori estimate (MAP):*

*What **LABEL maximizes** the likelihood function **AND** prior?*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*Example: Text Classification*

**Does this news article talk about politics?**

*Training Set: Collection of New Articles*

Example: Text Classification

**Does this news article talk about politics?**

Training Set: Collection of New Articles

Article 1: The computer contractor who exposed....
Article 2: The parents of a missing U.S. journalist in Syria...

*Q: What are my features?*

Q: What are my features?

A: The text in the documents.

Q: What are my features?

A: The text in the documents.

Q: How to I represent them?

*Q: What are my features?*

*A: The text in the documents.*

*Q: How to I represent them?*
*A: Binary occurrence? Word counts?*

*the, computer, contractor, exposed, parents, missing, Syria, U.S.*

| the | computer | contractor | exposed | parents | missing | Syria | U.S. |
| --- | --- | --- | --- | --- | --- | --- | --- |
| 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |

| computer | contractor | exposed | parents | missing | Syria | U.S. |
|----------|-----------|---------|---------|---------|-------|------|
| 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 | 1 | 1 |

We can make some alterations
1) Drop stop words (commonly occurring words that don't have meaning)

| computer, | contractor, | exposed, | parents, | missing, | Syria, | U.S., | **POL** |
|-----------|-------------|----------|----------|----------|--------|-------|---------|
| 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 |

Our goal is to compute compute $P$ ( POL = $T$ | words in the text)

We need to **learn** $P$( word | POL ) i.e. $P$ ( Syria | POL )

| computer, | contractor, | exposed, | parents, | missing, | Syria, | U.S., | **POL** |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 |

Once we've learned P(computer | POL), P(U.S. | POL) on our training set, we want to label our test set

| computer, | contractor, | exposed, | parents, | missing, | Syria, | U.S., | **POL** |
|-----------|-------------|----------|----------|----------|--------|-------|---------|
| 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 |

The correct label, POL = True or POL = False is the one that maximize our posterior.

computer, contractor, exposed, parents, missing, Syria, U.S., **POL**

| computer | contractor | exposed | parents | missing | Syria | U.S. | POL |
|----------|-----------|---------|---------|---------|-------|------|-----|
| 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 |

Compute probability in each class:

$$P ( POL = T \mid \{x\} ) = P ( \{x\} \mid POL = T) * P(POL=T)$$

$$P ( POL = F \mid \{x\} ) = P ( \{x\} \mid POL = F) * P(POL=F)$$

computer, contractor, exposed, parents, missing, Syria, U.S., **POL**

| computer | contractor | exposed | parents | missing | Syria | U.S. | POL |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 |

Article 2: The parents of a missing U.S. journalist in Syria...

$$P ( POL = T \mid \{x\} ) = P ( \{x\} \mid POL = T ) * P(POL=T)$$

$$= P(Syria \mid POL=T) * P(journalist \mid POL=T) * P(parents \mid POL=T) \dots$$

$$* P( POL=T )$$

# LAB