

Introduction to Machine Learning: Naïve Bayes, Linear and Logistic Regression

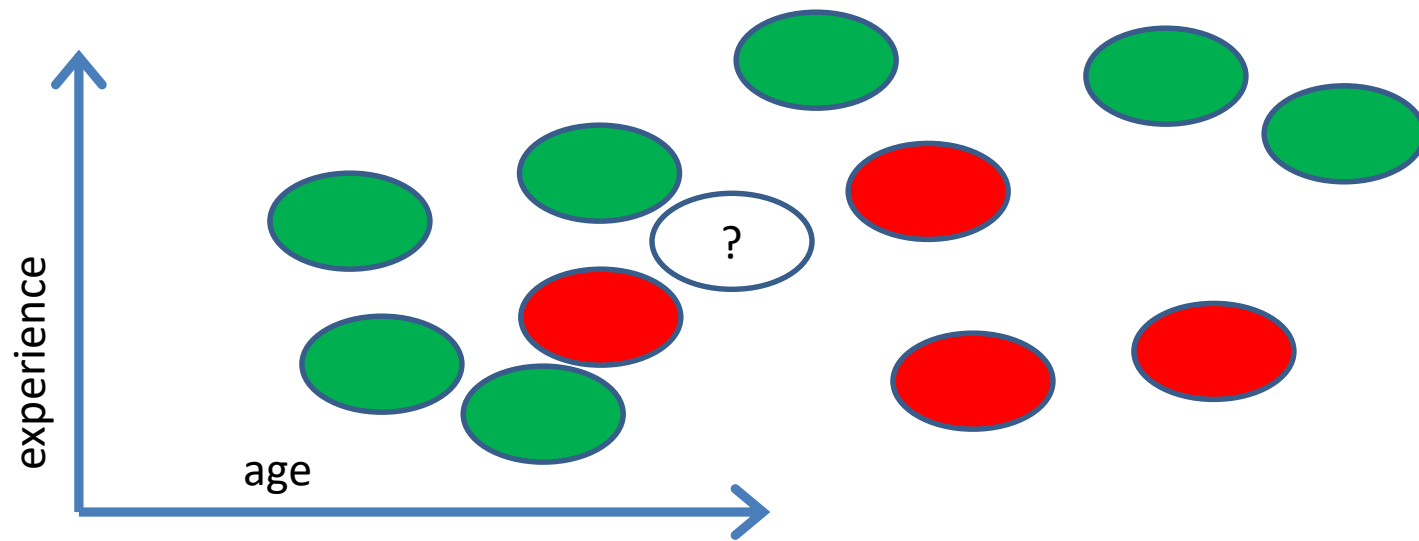
Amos Azaria

Databases and Machine Learning

- So far we have dealt with retrieving data that was previously stored.
- Many times we want to retrieve missing data, that is, data that we do not have. What if we want to know what will happen on new, unobserved data?
- For example, suppose some of our users are tagged as employed or unemployed, and our system needs to decide whether to show an ad which is relevant only to unemployed users.

Machine Learning (cont.)

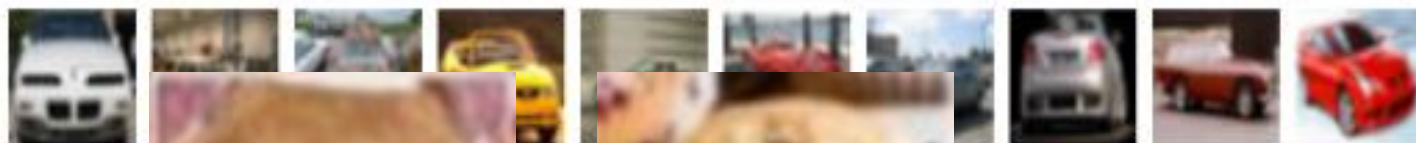
- Assume we are given a training set (e.g. data on people and an indicator whether they are employed or not), and then we are given a new example (e.g. a new person), and we want to predict a value (e.g. employment).



airplane



automobile



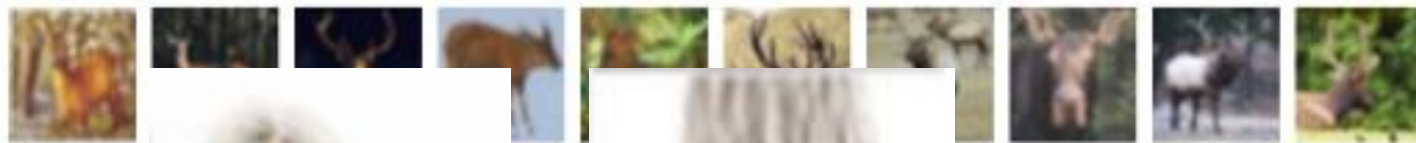
bird



cat



deer



dog



frog



horse



Video...

<https://www.youtube.com/watch?v=1eBxt9HUfh8>

Building a Machine Learning Classifier

- To solve the employment problem (or similar problems), we build a machine learning model (called a classifier) that, when given a new instance (a user), predicts a label (whether she is employed or not).
- There are many machine learning classification algorithms:
 - Naïve Bayes
 - Logistic regression
 - Neural networks
 - SVM
 - Decision Trees
 - KNN

Train / Test Split

- In order to train our classifier, we must have some labeled data (for example, we must know for some users whether they are employed or not).
- In order to evaluate the performance of our classifier, we split our data into train and test sets, train our classifier on the train set, and then test it on the test set.

Naïve Bayes

Amos Azaria

Spam Classifier

Spam:

- Buy it, pay later! Click me!
- You Won 10000 Dollars! Click here!

Real (Non-Spam):

- Will See you later.
- Will you want to meet later?
- I'm waiting for you.

Test:

- Are you paying too much? Click now!

Counting Messages / Documents

$$y^* = \operatorname{argmax}_{k \in \{1, \dots, K\}} p(y = k) \prod_{i=1}^n p(x_i | y = k)$$

If we perform stemming (or lemmatization): paying=pay

	Examples	are	you	paying	too	much	?	click	now	!
Spam	3	1	2	2	1	1	1	3	1	3
Real	4	1	4	1	1	1	2	1	1	1

- Laplace's Smoothing (rule of succession)

$$p(y = 0 | x) = \frac{3}{7} \cdot \frac{1}{4} \frac{2}{4} \frac{2}{4} \frac{1}{4} \frac{1}{4} \frac{1}{4} \frac{3}{4} \frac{1}{4} \frac{3}{4} = 5.$$

$$p(y = 1 | x) = \frac{4}{7} \cdot \frac{1}{5} \frac{4}{5} \frac{1}{5} \frac{1}{5} \frac{1}{5} \frac{2}{5} \frac{1}{5} \frac{1}{5} \frac{1}{5} = 2.$$

- For smoothing, we may take into account the number of words and length of every message.
- In practice, we many times maximize the log-likelihood (to deal with such small numbers).
- Note that when using the log, we sum rather than multiply.

Spam:

- Buy it, pay later!

- You Won 100!

Real (Non-Spam):

- I'm waiting for you.

- Will you want to buy it?

- I'm waiting for you.

Test:

- Are you paying too much? Click now!

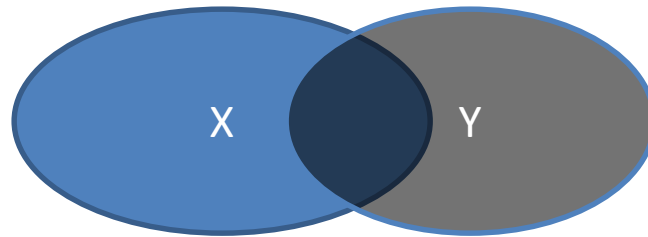
We assume that there is a message that includes the word and one that doesn't

Naïve Bayes (Why? How?)

- Notation:
 - The instances (messages in our case) are denoted by: x_1, x_2, \dots, x_m
 - m in our case is 5
 - (so x_1 = Buy it, pay later! Click me!)
 - Note that each of the x s is a vector: $x_t = x_{t1}, x_{t2}, \dots, x_{tn}$
 - The classes (Spam/Real in our case) are denoted by: y_1, y_2, \dots, y_m
 - (so y_1 = spam)

Bayes Rule

$$p(X, Y) = p(X)p(Y | X)$$



$$p(X)p(Y | X) = p(X, Y) = p(Y)p(X | Y)$$

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

This is an important formula which is always true!

Naïve Bayes (Cont.)

- Suppose we have a new message x_t (a vector) and we want to find the most likely class (y), that is, we want to know if x_t is more likely to present class 0 (e.g. Spam) or class 1 (not-spam), etc.
- For that we want to compute for every k :
 $p(y=k | x_t)$
- We then pick the k that maximizes the above.
- Denote that k as y^* .

Naïve Bayes


- Naïve Bayes makes the conditionally independent assumption:

Always true $p(x_t, y_t) = p(y_t)p(x_{t1} | y_t)p(x_{t2} | y_t, x_{t1})p(x_{t3} | y_t, x_{t1}, x_{t2}) \cdot \dots$

Naïve Bayes assumption $p(x_t, y_t) = p(y_t)p(x_{t1} | y_t)p(x_{t2} | y_t)p(x_{t3} | y_t) \cdot \dots$

Bayes Rule
$$p(y = k | x_t) = \frac{p(y=k) \prod_{i=1}^n p(x_{ti} | y=k)}{p(x_t)}$$

$$y^* = \operatorname{argmax}_{k \in \{1, \dots, K\}} p(y = k) \prod_{i=1}^n p(x_{ti} | y = k)$$



Independent of K (so we can ignore when looking for the maximum)

Naïve Bayes (Cont.)

- Note that the calculations shown previously (the fractions of the appearances of every word) aren't an immediate corollary of the formula we derived.
- In general we try to maximize the likelihood of the data: $\prod_{i=0}^N p(x_i, y_i)$
- It turns out (and easy to show), that if we use the fraction of messages each word appears in (as we have done), we in-fact maximize the likelihood of the data.
- We end up with:

$$y^* = \operatorname{argmax}_{k \in \{1, \dots, K\}} p(y = k) \prod_{i=1}^n \frac{\text{num of messages in class } k \text{ with word } i}{\text{num of messages in class } k}$$
$$y^* = \operatorname{argmax}_{k \in \{1, \dots, K\}} p(y = k) \prod_{i=1}^n \frac{\sum_{t=1}^m 1\{x_{ti}=1 \wedge y_t=k\}}{\sum_{t=1}^m 1\{y_t=k\}}$$

Multiple predictions

- How would we prepare the data for multiple queries?
 - Determine how many documents we have in total. Denote this value by p_{tot} .
 - We first count how many documents every class has. Denote these values by: p_k .
 - [We then find all words that appear in all documents and form a dictionary.]
 - For each word in the dictionary (i) and for every class (k), we need to count how many documents it appears in. Denote these values by p_{ki} .
- For every new query we simply compute:

$$y^* = \operatorname{argmax}_{k \in \{1, \dots, K\}} \frac{p_k}{p_{\text{tot}}} \prod_{i=1}^n \frac{p_{ki}}{p_k}$$

Online Naïve Bayes: Parallel Execution (In Spark)

- Suppose the data is stored in an RDD as (message, class) tuples. E.g. (greeting / valediction):

```
input_data = sc.parallelize([("hello there", 0), ("hi there", 0), ("go home", 1),  
("see you",1), ("good bye to you", 1)])
```

- We need to find p_{tot} , p_k and all p_{ki} :

```
>>> pk = input_data.map(lambda (message, cls): (cls,  
1)).reduceByKey(lambda a,b: a+b).collectAsMap()
```

```
>>> ptot = sum(pk.values())
```

```
>>> pki = input_data \  
    .flatMap(lambda (message, cls): list(set([(cls,w) for w in message.split()]))) \  
    .map(lambda (cls, word): ((cls,word), 1)) \  
    .reduceByKey(lambda a,b: a+b).collectAsMap()
```

Answering a new query

$$y^* = \operatorname{argmax}_{k \in \{1, \dots, K\}} \frac{p_k}{p_{tot}} \prod_{i=1}^n \frac{p_{ki}}{p_k}$$

numpy is a Python library that helps dealing with matrices and other math operations.

```
>>> import numpy as np
```

```
>>> query = "hello hi"
```

Cast not required in Python 3

```
>>> class_probs = [pk[k]/float(ptot)*np.prod
(np.array([pki.get((k,i),0)/float(pk[k]) for i in query.split()])) for k
in range(0,2)]
```

0 as default

```
>>> print(class_probs)
```

No Laplace smoothing,
so we get zeros.

[0.10000000000000000000000001, 0.0]

```
>>> y_star = np.argmax(np.array(class_probs))
```

```
>>> print(y_star)
```

0

It is indeed a
greeting!

What is a good classifier?

- Accuracy?
- What if we classify cancer with a 1% chance?
 - If we classify all as healthy, we will have 99% accuracy.
 - What would you think about a classifier that:
 - catches 98% of the cancers,
 - but has only 30% precision (if you are told that you might have cancer, there is only 30% that you actually have).

Recall and Precision

Confusion Matrix	Classified as Positive	Classified as Negative
Really Positive	True Positive	False Negative
Really Negative	False Positive	True Negative

- Accuracy?
 - Trues / All
 - $(\text{True Positive} + \text{True Negative}) / (\text{True Positive} + \text{True Negative} + \text{False Negative} + \text{False Positive})$
- Recall
 - What fraction of positives did we actually find?
 - True Positive / Really Positive
 - True Positive / (True Positive + False Negative)
- Precision:
 - If we say positive, how precise are we?
 - True Positive / Classified as Positive
 - True Positive / (True Positive + False Positive)

F-Measure (F_1 -Score, F-Score)

- A combination of precision and recall:
- $2 (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$
- E.g.:
 - Precision = 1, recall = 0,
 - F-Measure = 0
 - Precision = recall = 0.5
 - F-Measure = 0.5
 - Precision = recall = x
 - F-Measure = x