

CASE STUDY: FASTEST AIRLINES (PART 2)

1. OBJECTIVES

The purpose of this study is to develop a very popular predictive model called Naive-Bayes. You will be using Python as a tool to help you analyze data of flight delays in order to determine the likelihood of delays.

The objectives are to:

- Use Python to read and analyze data
- Apply the concept of conditional probability and Bayes' Theorem
- Develop a simple Naive-Bayes classifier model using m-estimates

2. PREDICTING DELAYS

As any frequent flier knows, flight delays are an inevitable part of commercial air travel. In the previous case study, you were able to determine which airline is the fastest for a particular route. In doing so, you had to take into account departure and arrival delays of each airline. While the cause for delays are largely unpredictable, we can make some insights on which flights are more likely to be delayed without taking weather conditions into consideration.

An obvious predictor of delays is the airline itself. For example, in 2014, only 74.34% of American Airline's flights were on-time while Delta had a 85.21% on-time record [1]. How can we use historical data to make predictions about delays? You may have already learned about several predictive models such as ordinary least-squares. We will study a very popular model in statistics and machine learning called *Naive-Bayes*.

2.1. Naive-Bayes Classifier. Recall that *Bayes' Theorem* tells you the probability of an event based on conditions possibly related to that event. That is, given some event A and B ,

$$(1) \quad P(A|B) = \frac{P(A)P(B|A)}{P(B)}$$

The Naive-Bayes classifier is based on Bayes' Theorem with the assumption of independence between attributes. Suppose we have a vector $X = (x_1, x_2, \dots, x_n)$ representing n **attributes** (independent variables). We wish to compute

$$(2) \quad P(C_k|X) = \frac{P(C_k)P(X|C_k)}{P(X)}$$

where C_k represents possible outcomes (**classes**). If we are trying to decide whether a particular flight is likely to be delayed, the outcome would be from the set $C_k \in \{Y, N\}$, corresponding to the flight being delayed (Y) or not (N). The vector X would contain information about the flight (e.g., airline, origin airport, destination airport) that could give some insight into which class will occur. One approach to picking a class is to pick the value C_k that is most probable, given that we know the values of the attributes. In other words, we pick the C_k that maximizes $P(C_k|X)$. Since the vector X is known, and $P(C_k|X) = P(C_k \cap X)/P(X)$ by the definition of conditional probability, then $P(X)$ will be the same for all possible classes, and we can instead focus on selecting the class that maximizes $P(C_k \cap X)$.

Noting that the attribute vector X typically contains values of several attributes, so we can rewrite $P(C_k \cap X) = P(C_k \cap x_1 \cap x_2 \cap \dots \cap x_n)$, where event x_i tells us the particular value of attribute i . In practice, the symbol \cap is often replaced with a comma when intersecting many events. By repeatedly applying the definition of conditional probability, we have

$$\begin{aligned}
 P(C_k \cap X) &= P(C_k)P(X|C_k) \\
 &= P(C_k)P(x_1, \dots, x_n|C_k) \\
 &= P(C_k)P(x_1|C_k)P(x_2, \dots, x_n|C_k, x_1) \\
 &\quad \vdots \\
 &= P(C_k)P(x_1|C_k)P(x_2|C_k, x_1) \cdots P(x_n|C_k, x_1, x_2, \dots, x_{n-1})
 \end{aligned}
 \tag{3}$$

This decomposes $P(C_k \cap X)$ into a product of the probability of outcome C_k with conditional probabilities of each attribute. We assume **conditional independence** between each attributes (hence, the “naive” portion of the name). That is, $P(x_i|C_k, x_j) = P(x_i|C_k)$ where $i \neq j$, which allows us to discard all conditions except C_k in Equation (2). Hence, this equation simplifies to

$$P(C_k|X) \propto P(C_k) \prod_{i=1}^n P(x_i|C_k) \tag{4}$$

With this definition, we can define a **classifier** (a decision rule). A common classifier is to pick the most probable class, as discussed above. This is called the **maximum a posterior (MAP) rule** which assigns a label $\hat{y} = C_k$ for some k :

$$\hat{y} = \underset{k}{\operatorname{argmax}} P(C_k) \prod_{i=1}^n P(x_i|C_k) \tag{5}$$

2.2. M-Estimates. The conditional densities $P(x_i|C_k)$ can be defined in different ways (e.g., Normal, Poisson) depending on the problem. If our dataset is small, we often do not know the underlying distributions. An intuitive way to estimate $P(x_i|C_k)$ is to define

$$P(x_i|C_k) = \frac{\hat{n}}{n} \tag{6}$$

where \hat{n} is the number of observations which $C = C_k$ and $X = x_i$ and n is the number of observations $C = C_k$. However, for small data sets we may find that that $\hat{n} = 0$ or $n = 0$, which pose a problem. A way to avoid fix this problem is to use *m-estimates*:

$$P(x_i|C_k) = \frac{\hat{n} + mp}{n + m} \tag{7}$$

where m is called the equivalent sample size and p is the a priori estimate of $P(x_i|C_k)$. A typical choice for p is $p = 1/(\# \text{ of possible values for attribute } i)$, which assumes that all attribute values are equally likely. The idea behind m-estimates is to pretend we have m extra observations with class $C = C_k$, with $m \cdot p$ of them having attribute $X = x_i$. Essentially, m says how confident we are of our prior estimate p since as m increases, we have $P(x_i|C_k) \rightarrow p$.

Let us consider an example problem.

Example 1: Consider the following table containing information of delays.

| Obs | Origin | Destination | Airline | Delay |
|-----|--------|-------------|---------|-------|
| 1 | DFW | ORD | Delta | Y |
| 2 | DFW | ORD | Delta | N |
| 3 | DFW | ORD | Delta | Y |
| 4 | LAX | ORD | Delta | N |
| 5 | LAX | ORD | United | Y |
| 6 | LAX | LGA | United | N |
| 7 | LAX | LGA | United | Y |
| 8 | LAX | LGA | Delta | N |
| 9 | DFW | LGA | United | N |
| 10 | DFW | ORD | United | Y |

TABLE 1. Flight Delay 1

The class is the variable Delay and the attributes are Origin, Destination and Airline. Classify DFW-LGA on Delta using Naive-Bayes with MAP and m-estimates. Assume $m = 3$.

Solution: First, note that DFW-LGA on Delta is not contained in the table. However, we can compute the probability of delay for this flight using Equation (5). Using an equivalent sample size of $m = 3$, we compute the relevant conditional probabilities as follows:

$$\begin{aligned}
P(DFW|Y) &= \frac{3 + 3(0.5)}{5 + 3} = 0.5625 \\
P(DFW|N) &= \frac{2 + 3(0.5)}{5 + 3} = 0.4375 \\
P(LGA|Y) &= \frac{1 + 3(0.5)}{5 + 3} = 0.3125 \\
P(LGA|N) &= \frac{3 + 3(0.5)}{5 + 3} = 0.5625 \\
P(Delta|Y) &= \frac{2 + 3(0.5)}{5 + 3} = 0.4375 \\
P(Delta|N) &= \frac{3 + 3(0.5)}{5 + 3} = 0.5675
\end{aligned}$$

For $P(DFW|Y)$, we have five observations where $C_k = Y$, of which three observations have $x_i = DFW$. So, we have that $n = 5$, $\hat{n} = 3$ and $p = 0.5$ since there are only two values of the attribute Origin (DFW and LAX). Since Y and N each occur in five of the ten data points, we compute $P(Y) = P(N) = 0.5$, we apply Equation (5) to obtain

$$\begin{aligned}
P(Y|DFW, LGA, Delta) &\propto P(Y)P(DFW|Y)P(LGA|Y)P(Delta|Y) \\
&= 0.0385
\end{aligned}$$

and

$$\begin{aligned}
P(N|DFW, LGA, Delta) &\propto P(N)P(DFW|N)P(LGA|N)P(Delta|N) \\
&= 0.0698
\end{aligned}$$

Since $\hat{y} = \operatorname{argmax} \{0.0385, 0.0698\} = 0.0698$, we classify DFW-LGA on Delta as N .

Exercise 1: Using the Table 2, will the flight SEA-ATL on Southwest with good weather be delayed? Use m-estimates for the conditional probabilities with $m = 4$.

| Obs | Origin | Destination | Airline | Weather | Delay |
|-----|--------|-------------|-----------|---------|-------|
| 1 | SEA | ATL | Southwest | Poor | N |
| 2 | SEA | BOS | American | Good | Y |
| 3 | SEA | ATL | United | Poor | Y |
| 4 | SFO | BOS | Southwest | Poor | Y |
| 5 | SFO | ATL | American | Good | N |
| 6 | SFO | BOS | United | Poor | Y |
| 7 | SFO | BOS | United | Good | N |
| 8 | SEA | BOS | Southwest | Poor | Y |

TABLE 2. Flight Delay 2

Let us now use real historical flight data to predict flight delays. The dataset *FlightDelay.csv* contains 2,346 observations of flights departing from JFK or SFO to ATL, LAS or ORD in November 2015. There are 5 variables.

| Variable Number | Variable Name | Description |
|-----------------|-----------------|--|
| 1 | Carrier | IATA Carrier Identification |
| 2 | Origin | Origin airport |
| 3 | Destination | Destination airport |
| 4 | Departure Delay | Difference between scheduled and actual departure time (minutes) |
| 5 | Arrival Delay | Difference between scheduled and actual arrival time (minutes) |

Reference: Bureau of Transportation Statistics

Exercise 2: Using *FlightDelay.csv*, write a Python program to perform the following.

- (1) Read the dataset.
- (2) Compute the total delay time for each flight. The total delay time is the sum of departure delay and arrival delay time.
- (3) Assign any flight with total delay time greater than 15 minutes with “Y” for being delayed. Otherwise, the flight is not delayed and is assigned “N”.

Exercise 3: Using your Python program from the previous exercise, for each of the following flights, compute the probabilities of delay and not delay given the route and carrier and then classify whether the flight is delayed using Naive-Bayes (assume $m = 3$).

- (1) JFK-LAS on American Airlines (AA)
- (2) JFK-LAS on JetBlue (B6)

(3) SFO-ORD on Virgin Airlines (VX)

(4) SFO-ORD on Southwest Airlines (WN)

Your output should contain $P(\text{Delay}|\text{Ori}, \text{Dest}, \text{Carr})$, $P(\text{Notdelay}|\text{Ori}, \text{Dest}, \text{Carr})$ and a classification of whether the flight is delayed or not.

REFERENCES

- [1] “The World’s Best and Worst On-Time Airlines ” *Skift*. 18 Nov. 2014. Web. 5 Aug. 2015.
- [2] Mitchell, Tom M. “Machine Learning.” New York: McGraw-Hill. 1997.