

“AIR TRAFFIC PASSENGER DATA” DATASET - NAIVE BAYES MODEL

PROBLEM: USING THE NAÏVE BAYES MODEL TO IDENTIFY THE ASSOCIATION BETWEEN EACH FEATURE WITH THE NUMBER OF PASSENGERS

“HOUSING PRICES DATA” DATASET - LINEAR REGRESSION MODEL

PROBLEM: USING THE LINEAR REGRESSION MODEL TO IDENTIFY THE ASSOCIATION BETWEEN THE FEATURES AND THE HOUSING PRICE.

GROUP 13

PUN YI JIE (151625)

SEET ZHI NIE (151502)

TAN XIN YI (152804)

NAIVE BAYES MODEL

Data Preparation

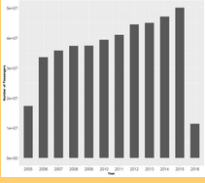
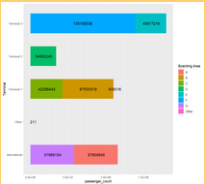
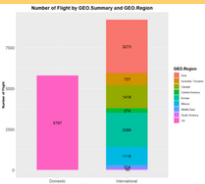
```
> str(df)
'data.frame':   15007 obs. of  16 variables:
 $ Activity.Period      : int  200507 200507 200507 200507 200507 200507 200507 200507 200507 200507
 $ Operating.Airline     : chr  "ATA Airlines" "ATA Airlines" "ATA Airlines" "Air Canada" ...
 $ Operating.Airline.IATA.Code : chr  "T2" "T2" "T2" "T2" ...
 $ Published.Airline     : chr  "ATA Airlines" "ATA Airlines" "ATA Airlines" "Air Canada" ...
 $ Published.Airline.IATA.Code : chr  "T2" "T2" "T2" "T2" ...
 $ GEO.Summary          : chr  "domestic" "domestic" "domestic" "International" ...
 $ GEO.Region           : chr  "US" "US" "US" "Canada" ...
 $ Activity.Type.Code    : chr  "Depanned" "Enplaned" "Thru / Transit" "Depanned" ...
 $ Price.Category.Code   : chr  "Low Fare" "Low Fare" "Low Fare" "Other" ...
 $ Terminal             : chr  "Terminal 1" "Terminal 1" "Terminal 1" "Terminal 1" ...
 $ Boarding.Area        : chr  "B" "B" "B" "B" ...
 $ Passenger.Count       : int  27271 29131 5415 35156 34090 6263 5500 12050 11638 4998 ...
 $ Adjusted.Activity.Type.Code : int  27271 29131 10830 35156 34090 6263 5500 12050 11638 4998 ...
 $ Adjusted.Passenger.Count : int  2005 2005 2005 2005 2005 2005 2005 2005 2005 2005 ...
 $ Year                 : chr  "July" "July" "July" "July" ...
```

```
> summary(df)
Activity.Period      Operating.Airline.IATA.Code   Published.Airline
200807 : 128   United Airlines - Pre 07/01/2013:2154   UA :3046   United Airlines - Pre 07/01/2013:2645
200808 : 127   Skywest Airlines                       :963   OO : 963   United Airlines :1107
201510 : 126   United Airlines                         :892   AS : 751   Alaska Airlines :969
201505 : 125   Alaska Airlines                         :751   DL : 386   Delta Air Lines :803
201603 : 125   Delta Air Lines                         :386   AC : 366   American Airlines :416
200710 : 123   Air Canada                             :366   VX : 362   US Airways      :407
(Other):14253 (Other)                                :9495 (Other):9133 (Other)                                :1860

Published.Airline.IATA.Code   GEO.Summary   GEO.Region   Activity.Type.Code   Price.Category.Code
UA :3752   Domestic :15797   US :5797   Depanned :17071   Low Fare: 1320
DL :803   International:9210   Asia :3273   Enplaned :17016   Other :13087
AA :416   Europe :2089   Canada :1418
US :407   Mexico :1115
AC :380   Australia / Oceania: 737
(Other):8280 (Other)                                :578

Terminal   Boarding.Area   Passenger.Count   Adjusted.Activity.Type.Code   Adjusted.Passenger.Count   Year
International:9107   A :2525   Min. : 1   Depanned :5707   Min. : 1   2015 :1460
Other : 27   G :3992   1st Qu.: 5374   Enplaned :17016   1st Qu.: 5496   2008 :1433
Terminal 1 :1241   B :1393   Median : 9354   Thru / Transit * 2: 920   Median : 9354   2007 :1469
Terminal 2 :324   C :1337   Mean : 29243   Mean : 29332   2009 :1393
Terminal 3 :12218   F :1228   3rd Qu.: 21159   3rd Qu.: 21182   2011 :1380
E : 841   Max. :459837   Max. :459837   2010 :1383
(Other):351 (Other)                                :6539

Month
August :1310
July :1303
September:1297
October :1295
January :1268
November :1263
(Other) :17271
```

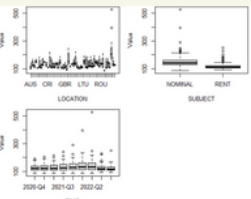


LINEAR REGRESSION MODEL

Data Preparation

```
> str(df1)
'data.frame':   698 obs. of  8 variables:
 $ LOCATION : chr  "AUS" "AUS" "AUS" "AUS" ...
 $ INDICATOR : chr  "HOUSECOST" "HOUSECOST" "HOUSECOST" "HOUSECOST" ...
 $ SUBJECT : chr  "NOMINAL" "NOMINAL" "NOMINAL" "NOMINAL" ...
 $ MEASURE : chr  "IDX2015" "IDX2015" "IDX2015" "IDX2015" ...
 $ FREQUENCY : chr  "Q1" "Q1" "Q1" "Q1" ...
 $ TIME : chr  "2020-Q4" "2021-Q1" "2021-Q2" "2021-Q3" ...
 $ Value : num  116 123 131 138 144 ...
 $ Flag.Codes : logi  NA NA NA NA NA ...

> summary(new_df1)
LOCATION      SUBJECT      TIME      Value
CAN : 17   NOMINAL:343   2020-Q4: 90   Min. : 84.96
CHE : 17   RENT :355   2021-Q1: 90   1st Qu.:110.12
GBR : 17   2021-Q2: 90   Median :124.53
IRL : 17   2021-Q3: 90   Mean :132.47
ISL : 17   2021-Q4: 89   3rd Qu.:146.53
NOR : 17   2022-Q1: 86   Max. :1528.10
(Other):596 (Other):163
```



Model Technical Details

```
Call:
lm(formula = Value ~ LOCATION + SUBJECT + TIME, data = new_df1)

Residuals:
    Min       1Q   -55.945    -7.466    -0.985     3Q    6.479     Max   259.507

Coefficients:
(Intercept) 127.8642  5.0872 25.135 < 2e-16 ***
LOCATIONAUS    20.0822  6.3861  3.145 0.001740 **
LOCATIONGBR    16.0042  6.3861  0.257 0.797100
LOCATIONGBR    16.0042  6.0592  1.986 0.047479 **
LOCATIONGBR    -32.2010  7.7342 -4.163 3.56e-05 ***
LOCATIONCAN    18.0139  6.3040  2.861 0.004365 **
LOCATIONGBR    -3.6278  6.3040 -0.575 0.565170
LOCATIONGBR    27.1900  6.3861  4.258 2.18e-05 ***
LOCATIONGBR    16.4013  6.0372  1.815 0.070034
LOCATIONGBR    15.4856  6.4794  2.390 0.017138 *
LOCATIONGBR    14.7221  9.0174  1.633 0.103040
LOCATIONGBR    33.0336  6.3861  5.172 3.10e-07 ***
LOCATIONGBR    13.2443  6.3861  2.074 0.038487 *
LOCATIONGBR    4.6688  6.3861  0.731 0.464994
LOCATIONGBR    3.8195  6.3861  0.598 0.549994
LOCATIONGBR    0.0592 -0.041 0.967584
LOCATIONGBR    1.3570  6.3861  0.212 0.831799
LOCATIONGBR    32.9525  6.3861  1.160 0.30e-07 ***
LOCATIONGBR    -5.5882  6.3861 -0.875 0.381878
LOCATIONGBR    -4.1556  6.3861 -0.651 0.511458
LOCATIONGBR    2.6102  6.3040  0.414 0.678979
LOCATIONGBR    -12.1193  6.3861 -1.898 0.058180 **
LOCATIONGBR    34.6766  6.3861  8.562 < 2e-16 ***
LOCATIONGBR    -21.0943  6.0592 -2.617 0.009069 **
LOCATIONGBR    -3.4124  6.0592 -0.426 0.670413
LOCATIONGBR    18.9528  6.3040  3.006 0.002747 **
LOCATIONGBR    39.5771  6.3040  6.278 6.13e-10 ***
LOCATIONGBR    1.3746  6.3861  0.215 0.829639
LOCATIONGBR    -13.9480  6.3861 -2.184 0.029317 *
LOCATIONGBR    -10.8175  7.0120 -1.543 0.123194
LOCATIONGBR    -5.1365  6.3861 -0.866 0.386127
LOCATIONGBR    36.3830  6.3861  5.697 1.86e-08 ***
LOCATIONGBR    20.4175  6.3861  3.197 0.001456 **
LOCATIONGBR    16.3483  6.3861  2.560 0.010697 *

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 17.09 on 639 degrees of freedom
Multiple R-squared:  0.7583, Adjusted R-squared:  0.7364
F-statistic: 34.57 on 58 and 639 DF, p-value: < 2.2e-16
```

Model Planning and Development

- determine the category from each features that contributes the highest conditional probability
- classify the passenger count based on the defined predictor variables

Model Technical Details

```
P(C_i | x_1, x_2, ..., x_n) \propto \left( \prod_{j=1}^n P(x_j | C_i) \right) \cdot P(C_i) \text{ for } 1 \leq i \leq k
```

```
> cm <- confusionMatrix(test_data$passenger_prediction)
print(cm)
Confusion Matrix and Statistics

          Reference
Prediction Low High
Low 2117 298
High 118 601

Accuracy : 0.8752
95% CI : (0.8635, 0.8863)
No Information Rate : 0.7304
P-Value (Acc > NRI) : < 2.2e-16

Kappa : 0.6619
McNemar's Test P-value : < 2.2e-16

Sensitivity : 0.9515
Specificity : 0.6685
Pos Pred Value : 0.8860
Neg Pred Value : 0.8339
Prevalence : 0.7304
Detection Prevalence : 0.7843
Balanced Accuracy : 0.8100

"Positive" Class : Low
```

FINDINGS:

|   |   |
|---|---|
| Naive Bayes Model   | Linear Regression Model   |
| Only US provided domestic flights while Asia have the most international flights. | The features that give higher housing price is the location "TUR", subject "NOMINAL" and during the time "2022-Q2". |
| There is an increasing number of passengers from 2005 to 2016.                    | The housing price increasing from "2020-Q4" until "2022-Q2", decrease on "2022-Q3" and then increase on "2022-Q4".  |
| The data is from one of the airports in the US.                                   |   |

RECOMMENDATIONS:

- The air traffic control might have to allocate some flights using terminal 3 or the international terminal to use terminal 2 to increase the air traffic efficiency.
- The housing price can always be compared using linear regression to ensure a suitable and reasonable price with the features provided.