# PREDICTING ON-TIME FLIGHTS

Rowan Rollman

# DATASET

- Data comes from the U.S. Department of Transportation - Bureau of Transportation Statistics
- The website provided allows you to create your own table
- 537,902 records
- 11 columns

- https://www.transtats.bts.gov/DL_SelectFields.aspx?gnoyr_VQ=FGJ&QO_fu146_anzr=b0-gvzr
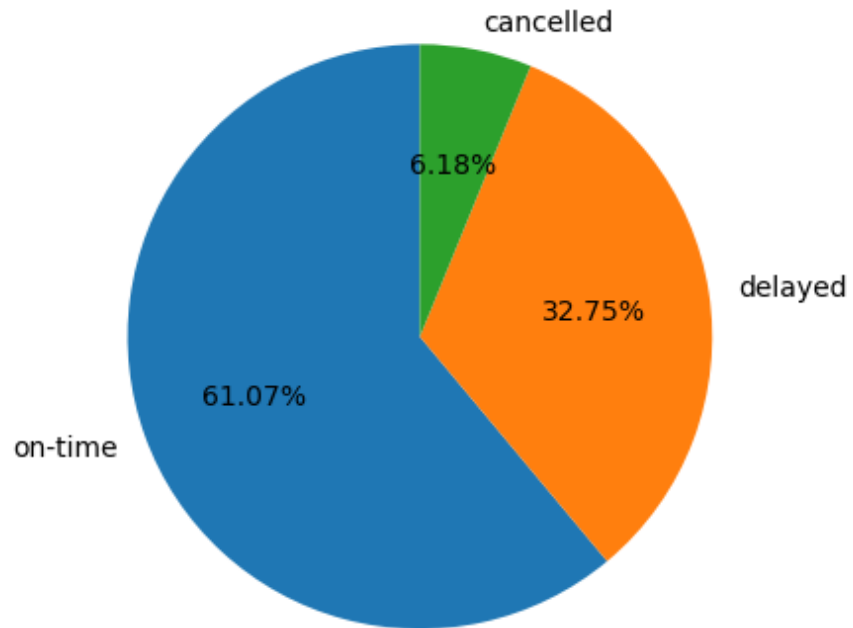
# EXPLORATORY DATA ANALYSIS

# VARIABLES

| Variable | Description |
|---|---|
| OP_UNIQUE_CARRIER | Unique carrier code, specific to a certain airline |
| ORIGIN | Code of the origin airport, where the flight departed from |
| DEST | Code of the destination airport |
| CRS_DEP_TIME | CRS Departure Time (local time: hhmm) |
| DEP_TIME | Actual Departure Time (local time: hhmm) |
| DEP_DELAY_NEW | Difference in minutes between scheduled and actual departure time. Early departures set to 0. |
| CANCELLED | Cancelled Flight Indicator (1=Yes) |
| DAY_OF_THE_MONTH | Numerical representation of the day of the month |
| DAY_OF_THE_WEEK | Numerical representation of the day of the week |

# DATA

- 17 unique airlines

- Create indicator for on-time, delayed, or cancelled

- Replace empty cells with 0

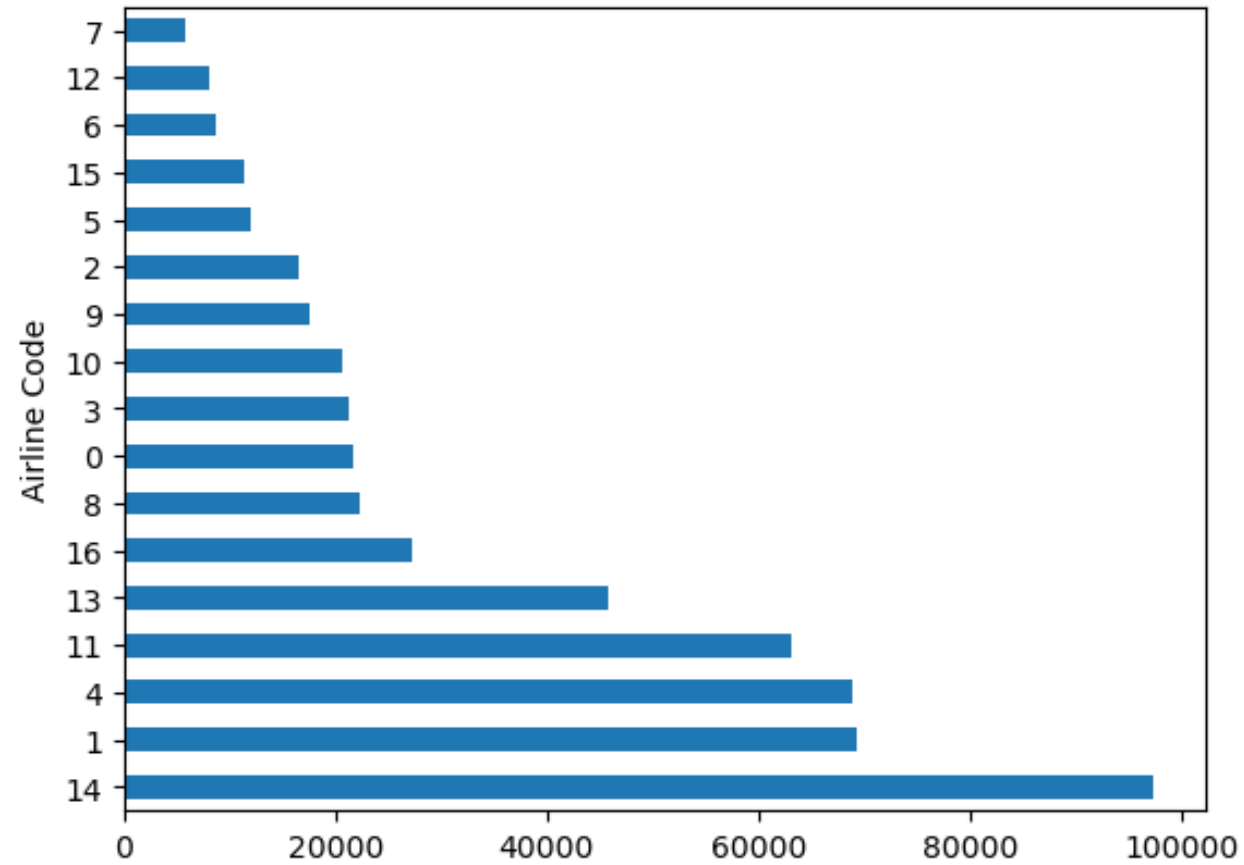- Double check for unexplained missing values

Class Distribution

cancelled

6.18%

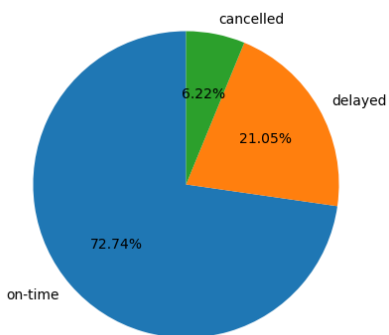delayed

32.75%

on-time

61.07%

Beginning distribution of classes is showed in the pie-chart

- 328470 on-time
- 176176 delayed
- 33256 cancelled

Distribution of Airline Records

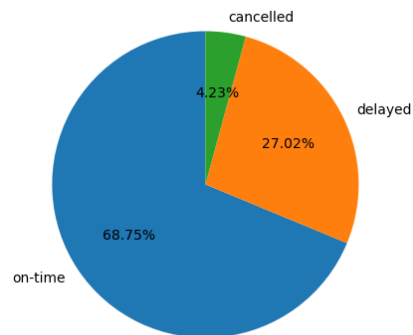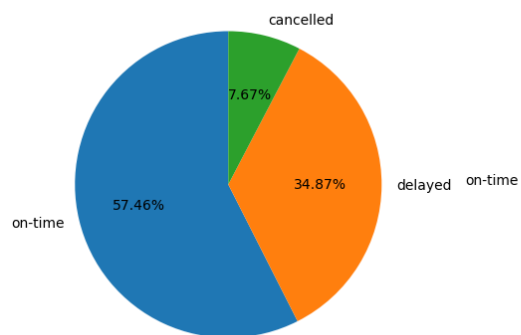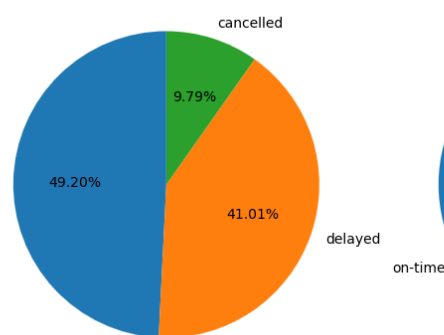Class Distribution for Airline 0 — cancelled 6.22%, delayed 21.05%, on-time 72.74%

Class Distribution for Airline 1 — cancelled 4.23%, delayed 27.02%, on-time 68.75%

Class Distribution for Airline 2 — cancelled 7.67%, delayed 34.87%, on-time 57.46%

Class Distribution for Airline 3 — cancelled 9.79%, delayed 41.01%, on-time 49.20%

Class Distribution for Airline 4 — cancelled 3.24%, delayed 30.31%, on-time 66.45%

Class Distribution for Airline 5 — cancelled 4.45%, delayed 38.43%, on-time 57.11%

Class Distribution for Airline 6 — cancelled 8.31%, delayed 35.78%, on-time 55.91%

Class Distribution for Airline 7 — cancelled 3.70%, delayed 39.52%, on-time 56.78%

Class Distribution for Airline 8 — cancelled 3.01%, delayed 25.05%, on-time 71.94%

Class Distribution for Airline 9 — cancelled 3.42%, delayed 33.53%, on-time 63.05%

Class Distribution for Airline 10 — cancelled 10.91%, delayed 28.95%, on-time 60.13%

Class Distribution for Airline 11 — cancelled 6.78%, delayed 30.05%, on-time 63.17%

Class Distribution for Airline 12 — cancelled 5.26%, delayed 38.46%, on-time 56.29%

Class Distribution for Airline 13 — cancelled 6.35%, delayed 39.99%, on-time 53.66%

Class Distribution for Airline 14 — cancelled 6.62%, delayed 40.66%, on-time 52.72%

Class Distribution for Airline 15 — cancelled 12.77%, delayed 31.87%, on-time 55.36%

Class Distribution for Airline 16 — cancelled 10.52%, delayed 23.32%, on-time 66.16%

# LOGISTIC REGRESSION

Model score: 92.413%
Cross validation scores:
[0.92064299 0.9212617 0.92081552 0.92079073
0.92007188]
Mean of cross validation scores: 0.920716564

Model score: 92.237%
Cross validation scores:
[0.92020921 0.92090227 0.920233 0.92029497
0.91974964]
Mean of cross validation scores: 0.920277816

# KNN CLASSIFIER

# KNN SCORES FROM 1-29

| K-Value | Train Score | Test Score |
|---------|-------------|------------|
| 1 | 0.977529956 | 0.863529552 |
| 2 | 0.913952001 | 0.861826646 |
| 3 | 0.914001576 | 0.858093637 |
| 4 | 0.888095958 | 0.849846813 |
| 5 | 0.887706791 | 0.851081234 |
| 6 | 0.871966606 | 0.843563164 |
| : | : | : |
| : | : | : |

- Best Test Score: 0.8635295517415748
- From K Value 1

# BEST K VALUE MODEL

Using the best K value from the previous slide, I created another KNN Model, with a score of 86.35%



Best KNN Classifier Model - Confusion Matrix

# CROSS VALIDATION ON KNN

The array of scores is as shown

Max Score:.85.646%

[0.8564594217462498, 0.8560603395076043, 0.852914783678l534, 0.84549830ll59704, 0.846554259669882l, 0.838932042872ll28, 0.840297844569l09, 0.834223797424044, 0.835l34566693872, 0.82994650229l8342, 0.830638080037699l2, 0.826220l23230466, 0.82687778684346Z7, 0.822953894805l658, 0.823538885l437371, 0.82008l00l5448632, 0.820638Z235098248, 0.817570007600639l, 0.8179195l37458823, 0.8l5433307679392S, 0.8l53787742500693, 0.8l33858450273384, 0.8l3383364869396Z, 0.8ll330945372961S, 0.8ll204526993Z571, 0.809347930084Z62l, 0.80927852443262Z3, 0.807588003422Z968, 0.80Z29Z98666034O9]

# XGBCLASSIFIER

Used multiple learning rates to find the optimal one,

learning_rates = [0.05, 0.1, 0.25, 0.5, 0.75, 1]

xgb.XGBClassifier(max_depth=10, n_estimators=100, learning_rate=lr)

| Learning Rate | Model Score |
| --- | --- |
| 0.05 | 0.9474999256372885 |
| 0.1 | 0.9483179154644695 |
| 0.25 | 0.9496638805437402 |
| 0.5 | 0.9521773401945328 |
| 0.75 | 0.9523483744311253 |
| 1.0 | 0.9521029774829709 |

Best learning rate is 0.75 with a score of 0.9523483744311253

# LIGHT GBM

lgbm.LGBMClassifier(num_leaves=100,n_estimators=100,max_depth=10,learning_rate=lr, bagging_fraction=.8, bagging_freq=5)

| Learning Rate | Model Score |
|---|---|
| 0.05 | 0.9518798893482852 |
| 0.1 | 0.9533373984948987 |
| 0.25 | 0.9535233052738035 |
| 0.5 | 0.882767185222642 |
| 0.75 | 0.545041494393930515 |
| 1.0 | 0.8444629524970999 |

Best learning rate is 0.25 with a score of 0.9535233052738035

# HISTGRADIENTBOOSTINGCLASSIFIER

HistGradientBoostingClassifier(learning_rate=lr)

| Learning Rate | Model Score |
|---|---|
| 0.05 | 0.9496936256283649 |
| 0.1 | 0.9515006395193194 |
| 0.25 | 0.5756640590142479 |
| 0.5 | 0.9450905737826825 |
| 0.75 | 0.9329248341711532 |
| 1.0 | 0.8978553793985544 |

Best learning rate is 0.1 with a score of 0.9515006395193194

# LAST TWO MODELS

## SGDClassifier

- Running on all processors

- Using log_loss

- Score 0.8893259763824028

## ADABoostClassifier

- N_estimators 100

- Random_state = 42

- Score 0.937341979237931