Exploring US Flights Data The dataset, used for this project, can be found in <u>Airline On-Time Performance Data</u> link. **Burcin Sarac** M.Sc. Business Analytics FT-18 In [1]: import os import pandas as pd import matplotlib.pyplot as plt In [11]: os.chdir("E:/dersler/Practicum 1/assignment") In [24]: flights = pd.read_csv("flights.csv", parse_dates=['FL_DATE', 'DEP_TIME','ARR_TIME']) flights.shape In [14]: Out[14]: (5674621, 20) flights.head() In [19]: Out[19]: FL_DATE TAIL_NUM CARRIER ORIGIN ORIGIN_CITY_NAME DEST DEST_CITY_NAME DEP_TIME DEP_DELAY ARR_TIME ARR_DELAY CANCELLED C 2017-01-N3CGAA DEN Denver, CO PHX Phoenix, AZ -10.0 1328 -17.0 0.0 AA 1135 2017-01-N3CGAA Phoenix, AZ Portland, OR -8.0 1653 PHX PDX 1502 -8.0 0.0 AΑ 2017-01-N3EVAA AA DCA Washington, DC MIA Miami, FL 0646 -13.0 0930 -14.0 0.0 01 2017-01-N3NAAA 1646 AA MIA Miami, FL LGA New York, NY 1402 -3.0 -14.0 0.0 2017-01-N3FCAA PHX Phoenix, AZ DEN Denver, CO 2310 16.0 0107 28.0 0.0 AΑ **Question 1** Since the DEP_DELAY column shows the difference in minutes between scheduled and actual departure time and early departures shown with negative numbers, I calculate the probability of delay in a particular airport mentioned at "ORIGIN" column, by using DEP_DELAY column data. For calculating this probability, I first group by the data into Departure airports(ORIGIN) by using groupby() function and I count all positive values(occured delays) and divide it to the number of all delay data by airport. When I use DEP_DELAY, I saw that there were some Null values, it seems they were mostly caused by cancelled flights. I tried calculations both dropping and keeping this Null values and it causes slightly change in mean, median and probabity values, which may lead me to draw wrong conclusions. So I decided to drop all Null values before calculations with using notnull() at all questions. After that, I dropped %1 of airports in terms of the lowest number of departures by using quantile function. 4 of total 320 locations dropped according to this criteria. After all, I merged these two calculations by matching ORIGIN data using merge() function in pandas. I use left join at this merging stage to avoid to get back the outlier data, which was dropped at the previous step. In other words, I keep ORIGIN data from dropped dataset and match the values in this column with the probabilities. With this, unmatched ORIGIN values, which are still kept in probability dataset, will drop automatically. At last, I again use groupby() function in ORIGIN to split the row data into groups and calculate mean and median of each airport(ORIGIN). And I merge these values with the previously created data frame to collect probability, mean and median of delays in every single airport. delayprob = (flights.loc[(flights.DEP_DELAY>0) &(flights.DEP_DELAY.notnull())].groupby("ORIGIN").DEP_DELAY.count()/flights.loc[f In [428]: lights.DEP_DELAY.notnull()].groupby("ORIGIN").DEP_DELAY.count()).sort_values(ascending=False) In [439]: delayprob Out[439]: ORIGIN SWO 1.000000 1.000000 TKI 0.833333 GGG MVY 0.668085 UST 0.600000 ACK 0.562986 HYA 0.559140 BQN 0.548326 0.526197 DAL HOU 0.501605 BPT 0.500000 OTH 0.484211 MDW 0.479777 OAK 0.479249 GUM 0.450549 WYS 0.449782 ADK 0.438776 SUX 0.428571 0.427414 LAX MMH 0.426667 0.422838 LAS STL 0.412243 FLL 0.410455 STS 0.410305 BUR 0.401379 SWF 0.398292 SF0 0.397707 BWI 0.397331 IAG 0.396040 0.389439 • • • 0.174267 COU GCC 0.173173 PIH 0.172979 CLL 0.172734 0.172377 BET 0.171334 LAR FLG 0.170862 0.169838 MFE RKS 0.169440 HYS 0.166392 0.165767 DRO 0.163952 COD 0.161631 HIB 0.160109 RHI FSM0.159449 GFK 0.155942 ROW 0.150655 0.143939 BTM0.141777 GTF TXK 0.138889 BIL 0.135160 INL 0.132901 GCK 0.132597 LCH 0.130471 HLN 0.125895 0.125612 HOB CPR 0.123288 GJT 0.117829 0.112069 BRO LWS 0.086551 Name: DEP_DELAY, Length: 320, dtype: float64 In [446]: origin = flights['ORIGIN'].value_counts() In [341]: origincutted = origin.loc[origin>origin.quantile(0.01)] In [342]: origincutted.count() Out[342]: 316 In [343]: origincutted = origincutted.reset_index() In [344]: origincutted.columns = ['ORIGIN','TotalFlight'] In [429]: | delaysmerged = pd.merge(origincutted,delayprob.to_frame('probofdelay'),how='left', on='ORIGIN') In [430]: | delaysmerged.head() Out[430]: ORIGIN TotalFlight probofdelay ATL 364655 0.359352 ORD 266460 0.334790 DEN 223165 0.366954 LAX 214297 0.427414 DFW 181208 0.360143 In [440]: | delaysmerged.shape Out[440]: (316, 3) In [431]: flights.loc[flights.DEP_DELAY.notnull()].groupby("ORIGIN").DEP_DELAY.mean().sort_values().head(3) Out[431]: ORIGIN -6.954738 YAK -5.388651 BET -5.285714 DLG Name: DEP_DELAY, dtype: float64 In [432]: flights.loc[flights.DEP_DELAY.notnull()].groupby("ORIGIN").DEP_DELAY.median().sort_values().head(3) Out[432]: ORIGIN YAK -14.0 -12.0 WRG -12.0 CDV Name: DEP_DELAY, dtype: float64 In [434]: groupmean = flights.loc[flights.DEP_DELAY.notnull()].groupby("ORIGIN").DEP_DELAY.mean() In [435]: delaysmerged2 = pd.merge(delaysmerged,groupmean.to_frame('meanofdelay'),how='left', on='ORIGIN') groupmedian = flights.loc[flights.DEP_DELAY.notnull()].groupby("ORIGIN").DEP_DELAY.median() In [436]: delaysmerged2 = pd.merge(delaysmerged2,groupmedian.to_frame('medianofdelay'),how='left', on='ORIGIN').sort_values(by=['probofdel ay'], ascending=False) In [438]: delaysmerged2.head() Out[438]: ORIGIN TotalFlight probofdelay meanofdelay medianofdelay MVY 293 241 0.668085 34.293617 11.0 UST 310 66 0.600000 59.369231 14.0 **ACK** 0.562986 30.203733 665 4.0 266 0.559140 17.182796 305 HYA 2.0 BQN 1710 0.548326 21.129501 3.0 182 It seems from results that, MVY airport has the highest probability with the ratio of 66% of flights delayed in departure. However it can be said from means of delay for example, the flight departures from UST airport has 60% probability of delay, however if it delays, the possible delay time is longer than MVY. On the other hand if HYA airport is taken into account, there is 55% posibility of departure delay occured, but the delay times were shorter. **Question 2** This time, I use same codes only changing groupby column into Carrier for checking probability that a flight operated by the airline has a delay. delayprobbycr = (flights.loc[(flights.DEP_DELAY>0)&(flights.DEP_DELAY.notnull())].groupby("CARRIER").DEP_DELAY.count()/flights.g In [293]: roupby("CARRIER").DEP_DELAY.count()).sort_values(ascending=False) In [294]: delayprobbycr Out[294]: CARRIER 0.468958 VX 0.413506 0.410006 В6 F9 0.363760 0.319286 AA0.318173 0.312332 NK 0.307152 DL 0.286598 EV 0.285814 HA AS 0.270324 0.266412 Name: DEP_DELAY, dtype: float64 delayprobbycr = delayprobbycr.reset_index() In [311]: delayprobbycr.columns = ['CARRIER', 'probofdelay'] In [315]: delayprobbycr In [316]: Out[316]: **CARRIER** probofdelay 0 WN 0.468958 VX 0.413506 1 2 0.410006 В6 3 F9 0.363760 0.319286 AA 5 UΑ 0.318173 0.312332 7 DL 0.307152 8 0.286598 ΕV 9 HA 0.285814 10 0.270324 AS 11 00 0.266412 crmean = flights.loc[flights.DEP_DELAY.notnull()].groupby("CARRIER").DEP_DELAY.mean().sort_values(ascending=False) In [306]: crmean.head(4) Out[306]: CARRIER 16.582982 13.219399 12.689932 EV 12.127640 Name: DEP_DELAY, dtype: float64 In [307]: crmedian = flights.loc[flights.DEP_DELAY.notnull()].groupby("CARRIER").DEP_DELAY.median().sort_values(ascending=False) In [308]: crmedian.head(3) Out[308]: CARRIER 0.0 WN VX -2.0 -2.0 Name: DEP_DELAY, dtype: float64 In [323]: delaysbycr = pd.merge(delayprobbycr,crmean.to_frame('meanofdelay'),how='left', on='CARRIER') In [325]: delaysbycr = pd.merge(delaysbycr,crmedian.to_frame('medianofdelay'),how='left', on='CARRIER') In [327]: airlines = pd.read_csv("L_UNIQUE_CARRIERS.csv") In [330]: airlines.head(2) Out[330]: Code Description 02Q Titan Airways 04Q Tradewind Aviation In [331]: delaysbycr = pd.merge(delaysbycr,airlines,how='left', left_on='CARRIER',right_on='Code') In [335]: delaysbycr Out[335]: CARRIER probofdelay meanofdelay medianofdelay Code Description 0 0.468958 10.248374 Southwest Airlines Co. WN 0.0 WN 1 0.413506 13.219399 -2.0 VX Virgin America VX 2 B6 0.410006 16.582982 -2.0 B6 JetBlue Airways 3 F9 0.363760 12.127640 -2.0 F9 Frontier Airlines Inc. AA 0.319286 8.180574 -3.0 AA American Airlines Inc. 5 0.318173 9.484484 -3.0 United Air Lines Inc. UΑ UA 6 0.312332 NK 10.171258 **-**4.0 NK Spirit Air Lines 7 DL DL 0.307152 8.497947 -2.0 Delta Air Lines Inc. 8 0.286598 EV ExpressJet Airlines Inc. ΕV 12.689932 **-**4.0 9 0.285814 Hawaiian Airlines Inc. HA 1.401808 -3.0 HA 10 0.270324 2.456161 AS Alaska Airlines Inc. AS **-**5.0 11 00 0.266412 10.325922 SkyWest Airlines Inc. **-**4.0 00 From the results it seems that Southwest Airlines has the highest delay probability in departure. Yet, according to the mean of departure delays with them shows that, delay times were shorter than the some of other airways, which have lower probability of delay. On the other hand, Alaska Airlines delayed in 27% of its departures and the average delay time was around 2.46, which was one of the best ratios. **Question 3** This histogram shows the distribution of departures for airports. That is, the x-axis is bins containing flights and the y-axis is the number of airports with departing flights in each particular bin. I created "origin" data to draw a histogram with only filtering 'ORIGIN' column via value_counts() function. So I can see the the flight numbers in terms of departure airports. And then I used this dataset to create a histogram with .hist() function. And I also add labels and title for the graph. origin = flights['ORIGIN'].value_counts() In [377]: origin.hist(color='blue', edgecolor='black', In [452]: alpha=0.5, figsize=(12, 10)) plt.xlabel('Nr of Flights') plt.ylabel('Nr of Airports') plt.title(' Distribution of Departures for Airports') Out[452]: Text(0.5,1,' Distribution of Departures for Airports') Distribution of Departures for Airports 200 Nr of Airports 100 50 50000 100000 150000 200000 250000 300000 350000 Nr of Flights It can be seen from graph that, most of the airports, which is around 260, have less than 50 thousand departing flights in 2017, and there are a few outlier airports, which have more than 300 thousand departing flights only in 2017. **Question 4** This time, I will create a plot shows the number of flights and the number of delayed flights per month of year. With this I aim to check the distribution of delayed departures by month. For drawing this plot, I decide to create a new column which includes only month data. Since I read my csv data with using parse_dates function for FL_DATE column inside read_csv() command, it allows me to shape the dates as Months, years and even days with only a single command. And this time I used dt.to period('M') to create month column from that FL DATE column. After creating new column, I use groupby() again in the data by grouping "month" column for both delayed flights and all flights numbers. And then I concatenate them with again using merge() function. This command only works with dataframes, so I first reset indexes of that newly created datasets and rename their column names. After all I use to_frame() command to merge the datasets added into the merge command with a specific column name. In [463]: flights['month'] = pd.to_datetime(flights['FL_DATE']).dt.to_period('M') In [464]: | flights.tail() Out[464]: FL_DATE TAIL_NUM CARRIER ORIGIN ORIGIN_CITY_NAME DEST_DEST_CITY_NAME DEP_TIME DEP_DELAY ARR_TIME ... CANCELLED CANC 2017-12-5674616 N705SW WN TUS Tucson, AZ LAX 26.0 1653 ... Los Angeles, CA 1616 0.0 2017-12--7.0 5674617 N463WN WNTUS Tucson, AZ MDW Chicago, IL 0843 1251 ... 0.0 2017-12-5674618 N956WN TUS WN Tucson, AZ SAN San Diego, CA 1451 1.0 1459 ... 0.0 2017-12-5674619 N222WN WN TUS Tucson, AZ San Diego, CA 1038 13.0 1049 ... 0.0 SAN 2017-12-5674620 N488WN WN TUS Tucson, AZ SAN 0800 -5.0 0812 ... 0.0 San Diego, CA 5 rows × 21 columns delayedbymonth = flights.loc[(flights2.DEP_DELAY>0)&(flights.DEP_DELAY.notnull())].groupby("month").DEP_DELAY.count() In [467]: In [489]: delayedbymonth.reset_index() Out[489]: month DEP_DELAY **0** 2017-01 168679 **1** 2017-02 123753 **2** 2017-03 167523 **3** 2017-04 163216 **4** 2017-05 177062 197341 **5** 2017-06 **6** 2017-07 196822 **7** 2017-08 188299 **8** 2017-09 124503 9 2017-10 145727 **10** 2017-11 125099 **11** 2017-12 168280 flightbymonth = flights.loc[(flights.DEP_DELAY.notnull())].groupby("month").DEP_DELAY.count() In [470]: flightbymonth.reset_index() In [485]: flightbymonth.columns =['month', 'NumberofFlights'] In [491]: | flights_delays = pd.merge(flightbymonth.to_frame('NumberofFlights'), delayedbymonth.to_frame('NumberofDelays') , how='left', on='m onth') flights_delays In [492]: Out[492]: Number of Flights Number of Delays month 2017-01 441131 168679 123753 2017-02 404205 2017-03 480070 167523 2017-04 460931 163216 2017-05 482444 177062 2017-06 488887 197341 2017-07 503313 196822 2017-08 188299 499444 2017-09 443384 124503 2017-10 476493 145727 2017-11 452744 125099 2017-12 458881 168280 In [524]: flights_delays.plot() plt.xlabel('Months') plt.ylabel('Nr of Flights') plt.title('All Flights & Delayed Flights by Month') Out[524]: Text(0.5,1,'All Flights & Delayed Flights by Month') All Flights & Delayed Flights by Month 500000 450000 400000 of Flights 350000 NumberofFlights NumberofDelays 300000 ≥ 250000 200000 150000 Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec Months From the graph it can be said that, the departure delays only differentiates on December in a negative way, in other words departure delays don't affected from monthly changes except December. On December it seems that, the flights delayed more than general trend. And there is a slight change than regular trend on March and October, it means delays occured less than expected. **Question 5** Lastly, I try to show for each possible origin and destination, which airline has the best performance, in terms of mean departure delay. For creating this table, I used again groupby() function, but this time with three columns(ORIGIN,DEST,CARRIER) inside this command. This filter shows me best possible airlines in terms of average delay times in specific departure and arrival destinations. flights.loc[(flights.DEP_DELAY.notnull())].groupby(["ORIGIN","DEST","CARRIER"]).DEP_DELAY.mean().sort_values() In [534]: Out[534]: ORIGIN DEST CARRIER RHI IMT 00 -33.676471 PDX OKC 00 -32.000000 HIB INL 00 -25.250000 ISP F9 -16.000000 PHL MSN SF0 UA -16.000000 SAT MCO F9 -14.000000 AS SIT ANC -13.375000 BTVΕV -13.000000 IAD OMA SF0 UΑ -12.000000 MSN IAD UA -11.500000 SDF BOS DL -11.000000 IAH BOS AA -11.000000 00 MOB IAH-11.000000 PIA DTW ΕV -11.000000 GJT ATL 00 -11.000000 STL IAD ΕV -11.000000 MYR LGA ΕV -10.000000 DCA DTW 00 -10.000000 ΕV CLE MSN-10.000000 MKE UA -10.000000 IAH CLE 00 -10.000000 LNK MKE ΕV -10.000000 MHT ΕV -9.200000 LGA DTW DCA 00 -9.000000 ΕV MSNCLE -9.000000 IAH MSN UA -9.000000 LGA ORF DL -9.000000 SF0 UΑ -9.000000 OMA JNU YAK AS -8.870056 AS OAK PDX -8.779221 AGS UA ATL 140.000000 JAC MSP 00 145.000000 F9 PHX LAS 149.000000 BOS DL 150.000000 MIA SJC MSP 00 155.000000 WN 158.750000 **EWR** LAS ΕV 164.000000 MDT IAD 00 169.000000 PHX OMA AABDL 178.000000 DCA 00 **RNO** ORD 183.500000 00 MS0 BTM188.000000 ΕV PVD ERI 195.000000 00 210.000000 ATL SGF SMF DL 216.250000 SEA 00 STL UIN 219.000000 UA 227.000000 BTR DCA HSV DTW ΕV 228.000000 AA IND ORD 230.142857 00 PWM DTW 291.250000

 DL

ΕV

00

00

UΑ

00

ΕV

ΕV

00

00

00

Name: DEP_DELAY, Length: 8445, dtype: float64

LEX

MDT

RDM

OTH

EWR

BUR

JFK

IAD

IMT

MS0

MQT

LGA

IAD

SMF

EUG

OMA

PSP

CVG

JAX

TVC

PIH

TVC

307.000000

374.000000

387.000000

420.666667

482.333333

800.000000

888.000000

939.000000

1470.000000

1475.000000

1480.000000

can be filtered by using .loc function again to check best possible option in a specific route.

It is hard to tell something for all destinations but for example, for flights, which depart from RHI airport and arrive to IMT directly, best option seems Skywest Airlines(OO) if only the departure delay times taken into account. From this table best possible airline option for all destinations can be interpreted. And also it