

DoorDash Merchant Operations Analytics Project

Name: Jia Wen (Steven) Liu

Date: 2021-07-16

```
In [1]: # Import relevant libraries for data analysis
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import datetime as dt
import statsmodels.api as sm
import scipy.stats as ss
```

Project Goal: Provide a set of recommendations on how to improve DoorDash with provided sample delivery data

Business Recommendations:

- Utilize time-region based switchback testing to understand how many drivers are needed at each region at different hours of the day to maximize orders amount and revenue while minimizing the total delivery times. This would provide a better control of network equilibrium between each participant and further optimize the metrics for DoorDash.
- Increase the number of drivers or driver tips at each region (Mountain View, Palo Alto, and San Jose) for peak-hours to reduce the average of amount of orders per hour, average total time to customers, and less congestion for the ordering traffic. This may provide higher customer, merchant, and drivers experience overall, which leads to better experiences and branding for DoorDash that could drive up market share percentage.
- Collaborate with the merchants to increase their average order total amount for each order by providing a tier-based reward incentive system. The tier-based reward incentive system could be based on customer-merchant time diff, merchant-driver time diff, and total delivery time. The lower the time-diffs are the higher tiers/incentives the merchants are eligible for. By incorporating the merchants in the entire supply chain system with incentives, the merchants are more likely to execute the orders more efficiently and more responsibly. This would provide a higher revenue to DoorDash while providing a win-win situation to the merchants.
- Establish a more robust system to encourage customers to use scheduled orders (as ASAP=False). This allows for better preparation and relieve the supply chain and logistic burdens on all participants (merchants, drivers and DoorDash). From the t-test shown in Part F, customers spend more time on orders on average with scheduled orders. This would provide a higher revenue to DoorDash while providing a win-win situation to merchants.
- Perform location/geographical-based clustering algorithm on Restaurant ID in each region based on Order Total. This would further allow DoorDash to understand popular neighborhoods and businesses that customers frequently order from. DoorDash can plan drivers route more accurately based on these clusters and traffic. Additionally, DoorDash should leverage this data to promote and work with local merchants/businesses that are not frequently ordered by the customers. This helps the improve the DoorDash relationship with the local businesses and helping them out. Actions could be working with the advertising, sales, and marketing team to promote those businesses with highlights or discounts on DoorDash.
- Perform A/B or A/A testing to understand how much discounts are needed to see a positive impact on the order total for each restaurant.

Part A: Data Cleaning

Using the 'Sample deliveries data - 1 month.xlsx', perform exploratory data analysis and data cleaning. The plan is to:

- Convert the first four date/time columns into Datetime object. Assume the date is in the month of 2014 January for easier computation of time difference and timezone conversion because DoorDash founding date is in the early 2013. This assumption does not impact further analysis in this notebook.
- Convert UTC timezone to PT timezone.
- Check for missing values for each feature column. Either drop missing values or perform missing value imputation.

```
In [2]: # Generated from the provided 'Sample deliveries data - 1 month' file
# Store provided data in dataframe for ease of access
file_name = 'Sample deliveries data - 1 month.xlsx'
deliveries_data = pd.read_excel(file_name)

Out[2]: C:\Users\Steven\Anaconda3\Lib\site-packages\openpyxl\worksheet\_reader.py:312: UserWarning: Unknown extension 'x' is not supported, which is likely a typo.
warn(msg)
```

```
In [3]: deliveries_data.head()
```

```
Out[3]:
```

	Customer placed order datetime	Placed order with restaurant datetime	Driver at restaurant datetime	Delivered to consumer datetime	Driver ID	Restaurant ID	Consumer ID	Delivery Region	Is ASAP	Order total	Amount of discount	Amount of tip	Refunded amount
0	01/02/2012	01/03/2025	01/03/2009	01/03/2020	203	90	64736	Mountain View	True	16.33	0.0	0.85	0.0
1	13/15/2017	13/17/2045	13/18/2439	13/18/2701	309	56	64746	Palo Alto	True	76.14	0.0	6.42	0.0
2	17/19/2017	17/19/2151	17/19/1959	17/19/2109	212	190	12484	San Jose	True	16.77	6.0	2.52	0.0
3	12/04/2019	12/04/2032	12/04/2102	12/04/2157	352	194	13920	San Jose	True	25.03	0.0	5.00	0.0
4	08/23/4638	08/23/4932	08/23/4629	09/03/7153	313	9	7037	Palo Alto	True	51.57	0.0	5.16	0.0

```
In [4]: deliveries_data['Customer placed order datetime'] = '2014-01-' + deliveries_data['Customer placed order datetime'].dt.strftime('%m-%d-%Y')
deliveries_data['Placed order with restaurant datetime'] = '2014-01-' + deliveries_data['Placed order with restaurant datetime'].dt.strftime('%m-%d-%Y')
deliveries_data['Driver at restaurant datetime'] = '2014-01-' + deliveries_data['Driver at restaurant datetime'].dt.strftime('%m-%d-%Y')
deliveries_data['Delivered to consumer datetime'] = '2014-01-' + deliveries_data['Delivered to consumer datetime'].dt.strftime('%m-%d-%Y')
```

```
In [5]: deliveries_data.info()
Out[5]:
```

	Customer placed order datetime	Placed order with restaurant datetime	Driver at restaurant datetime	Delivered to consumer datetime	Driver ID	Restaurant ID	Consumer ID	Delivery Region	Is ASAP	Order total	Amount of discount	Amount of tip	Refunded amount
0	2014-01-01	2014-01-01	2014-01-01	2014-01-01	203	90	64736	Mountain View	True	16.33	0.0	0.85	0.0
1	2014-01-13	2014-01-17	2014-01-18	2014-01-18	309	56	64746	Palo Alto	True	76.14	0.0	6.42	0.0
2	2014-01-17	2014-01-19	2014-01-19	2014-01-19	212	190	12484	San Jose	True	16.77	6.0	2.52	0.0
3	2014-01-12	2014-01-04	2014-01-02	2014-01-05	352	194	13920	San Jose	True	25.03	0.0	5.00	0.0
4	2014-01-08	2014-01-23	2014-01-06	2014-01-03	313	9	7037	Palo Alto	True	51.57	0.0	5.16	0.0

```
In [6]: deliveries_data = deliveries_data.dropna(subset=['Placed order with restaurant datetime'])
```

```
In [7]: deliveries_data.info()
```

	Customer placed order datetime	Placed order with restaurant datetime	Driver at restaurant datetime	Delivered to consumer datetime	Driver ID	Restaurant ID	Consumer ID	Delivery Region	Is ASAP	Order total	Amount of discount	Amount of tip	Refunded amount
0	2014-01-01	2014-01-01	2014-01-01	2014-01-01	203	90	64736	Mountain View	True	16.33	0.0	0.85	0.0
1	2014-01-13	2014-01-17	2014-01-18	2014-01-18	309	56	64746	Palo Alto	True	76.14	0.0	6.42	0.0
2	2014-01-17	2014-01-19	2014-01-19	2014-01-19	212	190	12484	San Jose	True	16.77	6.0	2.52	0.0
3	2014-01-12	2014-01-04	2014-01-02	2014-01-05	352	194	13920	San Jose	True	25.03	0.0	5.00	0.0
4	2014-01-08	2014-01-23	2014-01-06	2014-01-03	313	9	7037	Palo Alto	True	51.57	0.0	5.16	0.0

```
Out[7]:
```

	Customer placed order datetime	Placed order with restaurant datetime	Driver at restaurant datetime	Delivered to consumer datetime	Driver ID	Restaurant ID	Consumer ID	Delivery Region	Is ASAP	Order total	Amount of discount	Amount of tip	Refunded amount
8	2014-01-14	2014-01-15	2014-01-15	2014-01-15	196	225	2469	Mountain View	False	35.85	0.0	2.00	0.0
11	2014-01-20	2014-01-21	2014-01-21	2014-01-21	36	176	57150	San Jose	True	40.69	0.0	4.78	0.0
13	2014-01-27	2014-01-27	2014-01-27	2014-01-27	343	259	14776	Palo Alto	False	58.09	0.0	2.90	0.0
18	2014-01-12	2014-01-12	2014-01-12	2014-01-12	354	8	5783	Palo Alto	True	27.21	0.0	2.92	0.0
21	2014-01-05	2014-01-05	2014-01-05	2014-01-05	115	53	4339	Palo Alto	True	49.45	0.0	1.99	0.0

```
Out[8]:
```

	Customer placed order datetime	Placed order with restaurant datetime	Driver at restaurant datetime	Delivered to consumer datetime	Driver ID	Restaurant ID	Consumer ID	Delivery Region	Is ASAP	Order total	Amount of discount	Amount of tip	Refunded amount
18044	2014-01-18	2014-01-18	2014-01-18	2014-01-18	357	96	36970	Palo Alto	True	62.01	0.0	7.72	0.0
18055	2014-01-11	2014-01-11	2014-01-11	2014-01-11	216	3	198226	Palo Alto	True	16.82	6.0	2.00	0.0
18067	2014-01-29	2014-01-29	2014-01-29	2014-01-29	216	20	199347	Palo Alto	True	16.33	6.0	0.82	0.0
18072	2014-01-30	2014-01-30	2014-01-30	2014-01-30	290	347	69857	Mountain View	False	37.38	0.0	1.44	0.0
18073	2014-01-29	2014-01-29	2014-01-29	2014-01-29	287	44	46332	Mountain View	True	44.99	0.0	6.00	0.0

4514 rows x 13 columns

```
In [8]: deliveries_data['Driver at restaurant datetime'] = deliveries_data.apply(lambda row: row['Delivered to consumer datetime'] - row['Placed order with restaurant datetime'], axis=1)
```

Part B: Compute for time-difference between customer, merchant, driver, and customer.

Compute the following data and store them in each of the deliveries_data dataframe for further analysis.

- Customer-Merchant Time Diff (unit: seconds):** = Place order with restaurant datetime - Customer placed order date time. It represents the amount of time delay between customer placing the order and restaurant receiving the order.
- Merchant-Driver Time Diff (unit: seconds):** = Driver at restaurant datetime - Placed order with restaurant datetime. It represents the amount of time delay between restaurant receiving the order and the driver arriving at the restaurant.
- Driver-Consumer Time Diff (unit: seconds):** = Delivered to consumer datetime - Driver at restaurant datetime. It represents the amount of time delay between the driver arriving at the restaurant and the customer receiving the order.

```
In [9]: deliveries_data['Customer-Merchant time diff'] = (deliveries_data['Placed order with restaurant datetime'] - deliveries_data['Customer placed order datetime']).dt.total_seconds()
deliveries_data['Merchant-Driver time diff'] = (deliveries_data['Driver at restaurant datetime'] - deliveries_data['Placed order with restaurant datetime']).dt.total_seconds()
deliveries_data['Driver-Consumer time diff'] = (deliveries_data['Delivered to consumer datetime'] - deliveries_data['Driver at restaurant datetime']).dt.total_seconds()
deliveries_data.head()
```

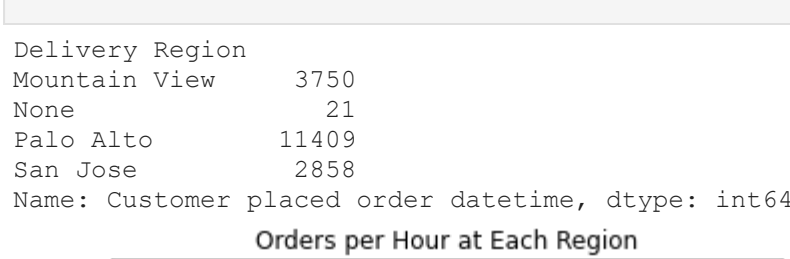
```
Out[9]:
```

	Customer placed order datetime	Placed order with restaurant datetime	Driver at restaurant datetime	Delivered to consumer datetime	Driver ID	Restaurant ID	Consumer ID	Delivery Region	Is ASAP	Order total	Amount of discount	Amount of tip	Refunded amount	Customer-Merchant time diff	Merchant-Driver time diff	Driver-Consumer time diff
0	2014-01-01	2014-01-01	2014-01-01	2014-01-01	203	90	64736	Mountain View	True	16.33	0.0	0.85	0.0	0.0	0.0	0.0
1	2014-01-13	2014-01-17	2014-01-18	2014-01-18	309	56	64746	Palo Alto	True	76.14	0.0	6.42	0.0	70.56	0.0	0.0
2	2014-01-17	2014-01-19	2014-01-19	2014-01-19	212	190	12484	San Jose	True	16.77	6.0	2.52	0.0	6.0	0.0	0.0
3	2014-01-12	2014-01-04	2014-01-02	2014-01-05	352	194	13920	San Jose	True	25.03	0.0	5.00	0.0	-12.00	0.0	0.0
4	2014-01-08	2014-01-23	2014-01-06	2014-01-03	313	9	7037	Palo Alto	True	51.57	0.0	5.16	0.0	15.84	0.0	0.0

Part C: Number of orders vs Hour of the Day

9AM-12PM and 12PM-19PM, which coincides with the time range of lunch and dinner. As a result, it is important to keep in mind that these hours are the peak-hours. Peak-hours create a greater inflow of orders and traffic to all the network members: DoorDash, customers, merchants, and drivers/dashers.

```
In [10]: deliveries_data['Customer_order_hours'] = deliveries_data['Customer placed order datetime'].dt.hour
num_order_per_hour = deliveries_data.groupby('Customer_order_hours').count().iloc[:, 0]
plt.bar(num_order_per_hour.index, num_order_per_hour.values)
plt.title('Number of orders placed per hour of the day')
plt.xlabel('Number of Orders')
plt.ylabel('Hour of the Day')
plt.xticks(np.arange(0,24,1))
plt.show()
```



Part D: Regional impact on orders and drivers

The goal of this section is to understand if specific regions have any impacts and correlations on the orders and drivers counts per hour. Business recommendations may vary a lot by regions since food delivery can be highly impacted by the hyperlocal environment variables, even if the three regions covered in this dataset are all in the same larger geographical area of California.

```
In [11]: # Define functions to plot
def plot_normalized_region_graph(df1, df2, title, index):
    plt.plot(df1['Mountain View'].index, df1['Mountain View'].values, label='Mountain View')
    plt.plot(df1['Palo Alto'].index, df1['Palo Alto'].values, label='Palo Alto')
    plt.plot(df1['San Jose'].index, df1['San Jose'].values, label='San Jose')
    plt.legend()
    plt.title(title)
    plt.xlabel('Hour of the Day')
    plt.ylabel(index)
    plt.xticks(np.arange(0,24,1))
    plt.show()

def plot_normalized_driver_graph(df1, df2, title, index):
    plt.plot(df1['Mountain View'].index, df1['Mountain View'].values/df1['total_driver_count'], label='Mountain View')
    plt.plot(df1['Palo Alto'].index, df1['Palo Alto'].values/df1['total_driver_count'], label='Palo Alto')
    plt.plot(df1['San Jose'].index, df1['San Jose'].values/df1['total_driver_count'], label='San Jose')
    plt.legend()
    plt.title(title)
    plt.xlabel('Hour of the Day')
    plt.ylabel(index)
    plt.xticks(np.arange(0,24,1))
    plt.show()
```

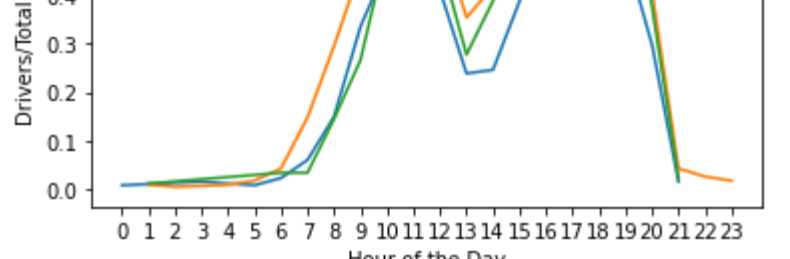
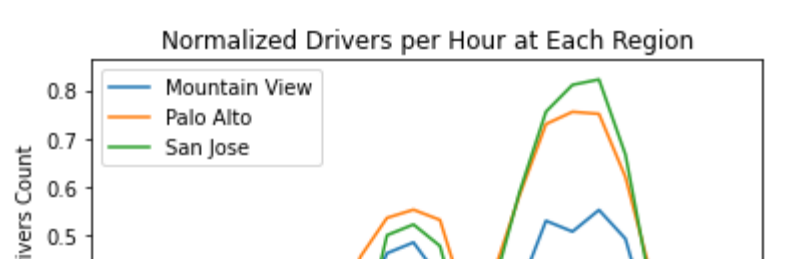
Orders per Hour at Each Region Discussion:

The trends in Orders count per hour at each region plot are similar to that in the Number of Orders placed per hour of the day plot in Part C. Palo Alto has a significantly higher amount of orders per hour at peaks hours than Mountain View and San Jose, indicating that it is a bigger market.

I further looked into how many unique customers, restaurants, and drivers are in each region. Although there are a higher number of restaurants in the San Jose area (155, compared to 148 in Palo Alto and 110 in Mountain View), they have the least number of orders among the three regions. This leads to several business proposals that need further investigation to boost the business in the San Jose area: improvement on the quality of restaurants, and promotions and marketing campaigns to further conversion for new customers.

The Normalized Orders per hour is the orders per hour divided by the total order amount at each region. It shows the normalized orders per hour are similar for all 3 regions.

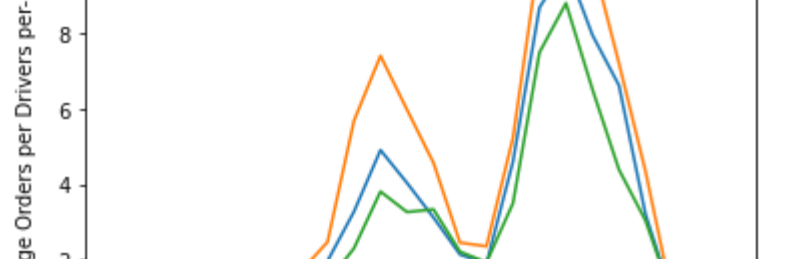
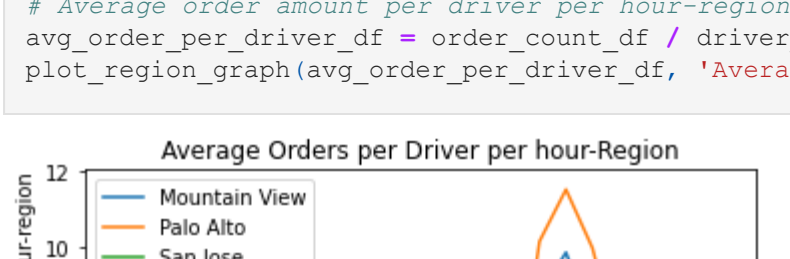
```
In [12]: # Order count
df1 = deliveries_data.groupby(['Delivery Region', 'Customer_order_hours']).count().iloc[:, 0]
total_order_count_df = deliveries_data.groupby(['Delivery Region', 'Customer_order_hours']).count().iloc[:, 0]
print('total_order_count_df')
plot_region_graph(df1, df2, title, index)
print('total_order_count_df')
```



Drivers per Hour at Each Region Discussion:

The drivers count per hour at each region shows they are similar to the 'Number of Orders placed per hour of the day' plot in Part C. Palo Alto has a significantly higher amount of drivers per hour at peaks hours than Mountain View and San Jose. The Normalized Drivers per hour is the drivers per hour divided by the total drivers amount at each region. It shows the normalized orders drivers per hour are similar for all 3 regions.

```
In [13]: # Driver count
df1 = deliveries_data.groupby(['Delivery Region', 'Customer_order_hours']).nunique().iloc[:, 0]
total_driver_count_df = deliveries_data.groupby(['Delivery Region', 'Customer_order_hours']).nunique().iloc[:, 0]
print('total_driver_count_df')
plot_region_graph(df1, df2, title, index)
print('total_driver_count_df')
```



Average order amount per driver per hour-region:

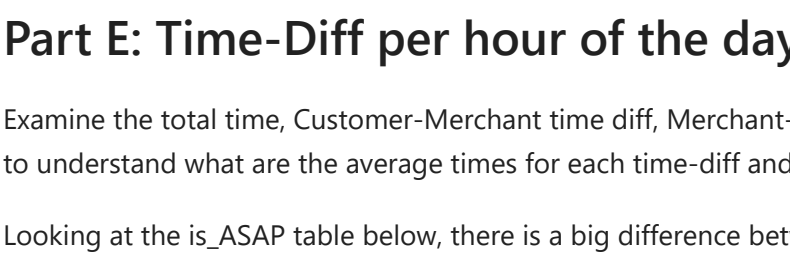
The average order count per driver per hour-region region shows they are similar to the Number of Orders placed per hour of the day plot in Part C. Mountain View has a significantly higher average order amount per driver per hour at peaks hours than Palo Alto and San Jose. With an estimate of maximum of 12 average orders per driver at peak-hour of 17PM, each driver would have to deliver each order in 5 minutes to the customer.

- Mountain View: max of average 12 orders per hour per driver with average 5 minutes delivery time to the customers.
- Palo Alto: max of average 10 orders per hour per driver with average 6 minutes delivery time to the customers.
- San Jose: max of average 9 orders per hour per driver with average 6.7 minutes delivery time to the customers.

Based on these computed metrics, there are too many orders per hour each driver needs to deliver at each region. With the expected average of 5 minutes delivery time, there are many factors like traffic, unexpected delays at the merchants, and other etc. in real life that most likely delay the drivers to deliver the orders successfully within 5 minutes.

In later part of this project, I computed the average delivery time from drivers getting to the merchants until customers getting the orders with 'Driver-Consumer Time Diff' computed in Part B. This would be used to compare to the computed average delivery time in this section.

```
In [14]: # Average order amount per driver per hour-region
avg_order_per_driver_df = order_count_df / driver_count_df
plot_region_graph(avg_order_per_driver_df, title, index)
```



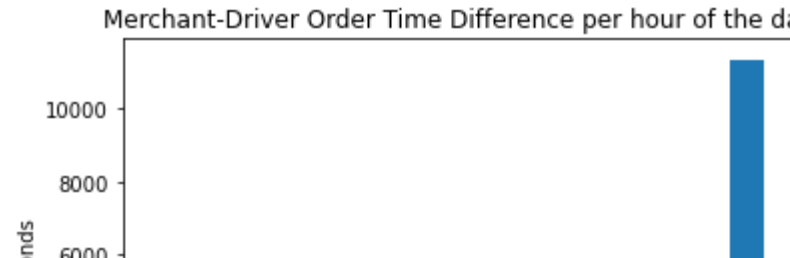
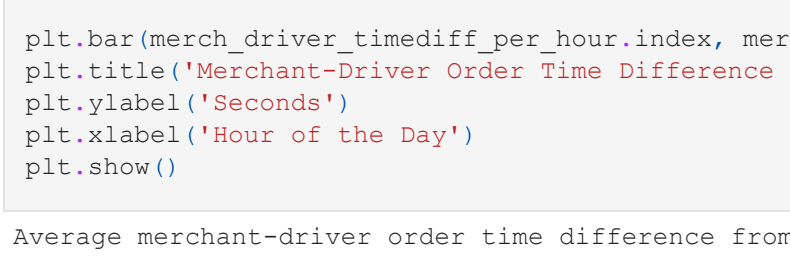
Correlation between Time and region: Computed the correlation between total time, Customer-Merchant time diff, Merchant-Driver time diff, Driver-Consumer time diff at each region for on-demand deliveries (Is ASAP=True). The following are significant if each of the times has an impact to the order count for each region. With a threshold of p-value of 0.005, the following relations are significant:

- Total Time at San Jose: corr (78.43%) and p-value(0.004)
- Customer-Merchant time diff at San Jose: corr (87.6%) and p-value (0.004)
- Driver-Consumer time diff at Mountain View: corr (63.7%) and p-value (0.048)

```
In [15]: deliveries_data['total_time'] = (deliveries_data['Delivered to consumer datetime'] - deliveries_data['Customer placed order datetime']).dt.seconds
for col in ['total_time', 'Customer-Merchant time diff', 'Merchant-Driver time diff', 'Driver-Consumer time diff']:
    print(col)
    for region in ['Mountain View', 'Palo Alto', 'San Jose']:
        print(region)
        print(ss.pearsonr(deliveries_data[deliveries_data['Is ASAP']==True].groupby([region], 'Customer_order_hours')[col].values, deliveries_data[deliveries_data['Is ASAP']==True].groupby([region], 'Customer_order_hours')[col].values))
print('')
```

```
total_time
Mountain View 3750
Palo Alto 21
San Jose 11409
Name: Customer placed order datetime, dtype: int64

Orders per Hour at Each Region
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23
0 1823 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900 1900
1 075857 095645 102439 105701 107000 111251 112231 111251 111251 111251 111251 111251 111251 111251 111251 111251 111251 111251 111251 111251 111251 111251 111251 111251
2 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237 110237
3 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157 200157
4 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638 154638
Name: Driver ID, dtype: int64
```



Correlation between Time and region: Computed the correlation between total time, Customer-Merchant time diff, Merchant-Driver time diff, Driver-Consumer time diff at each region for on-demand deliveries (Is ASAP=True). The following are significant if each of the times has an impact to the order count for each region. With a threshold of p-value of 0.005, the following relations are significant:

- Total Time at San Jose: corr (78.43%) and p-value(0.004)
- Customer-Merchant time diff at San Jose: corr (87.6%) and p-value (0.004)
- Driver-Consumer time diff at Mountain View: corr (63.7%) and p-value (0.048)

```
In [16]: deliveries_data.groupby(deliveries_data['Is ASAP']).median()['Order total']
```

```
Out[16]:
```

	Customer placed order datetime	Placed order with restaurant datetime	Driver at restaurant datetime	Delivered to consumer datetime	Driver ID	Restaurant ID	Consumer ID	Delivery Region	Is ASAP	Order total	Amount of discount	Amount of tip	Refunded amount
0	2014-01-01	2014-01-01	2014-01-01	2014-01-01	203	90	64736	Mountain View	True	16.33	0.0	0.85	0.0
1	2014-01-13	2014-01											

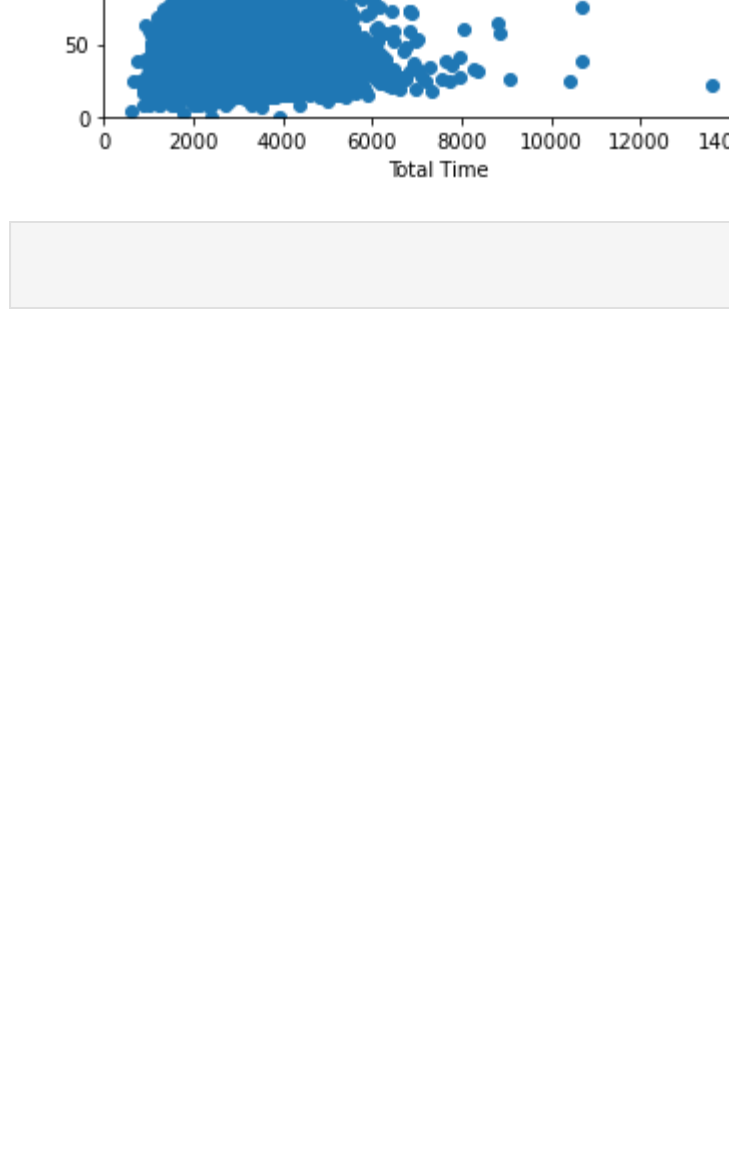
Using the t-test to compare the means of the two systems, on-demand orders (is ASAP=True) and scheduled orders (is ASAP=False), there is a significant difference in the mean values between the two as shown above. This is important for DoorDash and merchants to understand how can we incentive and encourage users to use scheduled orders more frequently.

Total time from customer placing the order to receiving order is significant to the amount of orders. As shown below, there is a 25.2% coreolation with a very low p-value. This is important to keep in mind as we are extracting insights from overall trends of all merchants. The majority of merchants have similar total time, but the goal of DoorDash is to work with the merchants to reduce the order time and maxmizing the experiences for everyone in the networks. Further recommendations and improvements are discussed in the business recommendations.

In [24]:

```
print(ss.pearsonr(deliveries_data_is_asap_df['total_time'],deliveries_data_is_asap_df['Order total']))

plt.scatter(deliveries_data_is_asap_df['total_time'],deliveries_data_is_asap_df['Order total'])
plt.xlim((0,15000))
plt.ylim((0,300))
plt.title('Average Order total dollars vs Total Time')
plt.ylabel('Order total')
plt.xlabel('Total Time')
plt.show()
```



In []: