

```
[15]: #Import Packages

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import plotly.offline as plt
import datetime
from datetime import datetime

#makes plotly plots render for github
from plotly.offline import download_plotlyjs, init_notebook_mode, iplob
init_notebook_mode(connected=True)

For this assignment, we are looking at customer behavior data on an app like Harvest that tracks company
hours, and tracks invoices sent and payments received.

First we will do an EDA on the data, and generate insights using the customer data, and then we will look
at some business problems like customer conversion and churn.

In [2]: #read the csvs
df_customers = pd.read_csv('customers.csv')
df_invoices = pd.read_csv('invoices_and_payments.csv')
df_expenses = pd.read_csv('expenses.csv')

df_customers['converted_at_date'] = pd.to_datetime(df_customers['converted_at'], infer
df_invoices['invoice_created_at_date'] = pd.to_datetime(df_invoices['invoice_created_

In [3]: #Print size/shape of dataframes
print(df_customers.shape)
print(df_invoices.shape)
print(df_expenses.shape)

(15028, 27)
(5214607, 16)
(5659162, 7)

In [4]: ##### CUSTOMERS EDA #####

display(df_customers.columns)

#check for duplicate rows
print('Number of Duplicate Rows:')
print(df_customers.duplicated().sum())

#check for duplicate company_id rows
print('Number of Duplicate Rows with company_id:')
print(df_customers['company_id'].duplicated().sum())

categorical_cols = ['billing_frequency',
                    'primary_industry_grouped', 'purpose',
                    'asana_enabled',
                    'slack_oauths', 'xero_enabled',
                    'quickbooks_enabled', 'stripe_enabled', 'tdd_enabled',
                    'team_module_enabled', 'estimate_module_enabled',
                    'client_dash_module_enabled', 'approval_module_enabled',
                    'majority_platform', 'tracking_frequency']

#Print unique values in each categorical column
print('UNIQUE VALUES IN EACH CATEGORICAL COLUMN:')
for col in df_customers:
    if col in categorical_cols:
        print(col,':', df_customers[col].unique())

#Print columns that contain Null values
print('Columns that contain Null values:')
print(df_customers.isnull().sum()[(df_customers.isnull().sum())>0])

#Find inactive users
print('Number of Rows Returning 0')
print('Invoices:', len(df_customers[df_customers['invoices']!=0]))
print('Payments:', len(df_customers[df_customers['payments']!=0]))
print('Total hours of payments received:', df_customers['payments'].sum())
print('Total hours of projects:', df_customers['projects'].sum())
print('Total hours of users:', df_customers['total_hours'].sum())
print('Projects:', len(df_customers[df_customers['projects']!=0]))

df_customers

Index(['company_id', 'converted_at', 'billing_frequency',
      'primary_industry_grouped', 'total_users', 'devices', 'purpose',
      'invoices', 'projects', 'payments', 'asana_enabled',
      'google_calendar_users', 'slack_oauths', 'xero_enabled',
      'quickbooks_enabled', 'stripe_enabled', 'tdd_enabled',
      'team_module_enabled', 'expense_module_enabled',
      'invoice_module_enabled', 'estimate_module_enabled',
      'client_dash_module_enabled', 'approval_module_enabled',
      'majority_platform', 'tracking_frequency', 'converted_at_date'],
      dtype=object)

DUPLICATE ROWS:
0

DUPLICATE ROWS WITH company_id:
0

UNIQUE VALUES IN EACH CATEGORICAL COLUMN:
billing_frequency: ['yearly' 'monthly']
primary_industry_grouped: [''other' 'unknown' 'tech' 'marketing' 'architecture' 'non
profit']
management consulting: ['engineering' 'design' 'healthcare' 'legal'
visual arts: ['construction' 'writing' 'accounting' 'events']
real estate: ['education' 'manufacturing' 'virtual assistant' 'music']
entertainment: ['ecommerce' 'retail' 'hospitality' 'research'
'architecture']
purpose: ['none' 'team' 'personal']
asana_enabled: [0 1]
slack_oauths: [0 1 2]
xero_enabled: [0 1]
quickbooks_enabled: [0 1]
stripe_enabled: [0 1]
tdd_enabled: [1 0]
team_module_enabled: [1 0]
expense_module_enabled: [1 0]
invoice_module_enabled: [1 0]
estimate_module_enabled: [0 1]
client_dash_module_enabled: [0 1]
approval_module_enabled: [1 0]
majority_platform: ['web' 'native' nan 'mix' 'other' 'integration' 'api']
tracking_frequency: ['weekly' 'daily' 'infrequently' 'monthly' 'work_daily' 'never']

COLUMNS THAT CONTAIN NULL VALUES:
majority_platform      811
dtype: int64

NUMBER OF ROWS RETURNING 0
invoices: 4514
payments: 6109
total_hours: 811
total_users: 8
projects: 151

Out[4]:
company_id  converted_at  billing_frequency  primary_industry_grouped  total_users  devices  purpose
0          315309      2014-02-11 16:25:08 UTC          yearly          other          0          14      none
1          438147      2015-08-27 19:58:50 UTC          yearly          unknown          0          0      none
2          1156438      2020-09-04 07:18:11 UTC          yearly          tech          0          0      team
3          262929      2014-01-27 15:50:30 UTC          yearly          marketing          0          7      none
4          194227      2013-04-02 03:26:27 UTC          yearly          unknown          0          2      none
...          ...          ...          ...          ...          ...          ...          ...          ...
15023      223993      2019-08-21 08:57:29 UTC          monthly          marketing          27          32      team
15024      203838      2013-04-30 12:50:06 UTC          monthly          tech          27          90      non
15025      526538      2016-05-16 21:36:57 UTC          monthly          tech          27          41      non
15026      409087      2015-05-06 08:27:42 UTC          monthly          architecture  27          5      non
15027      232338      2015-07-29 18:16:17 UTC          monthly          tech          27          90      non

15028 rows x 27 columns
```

```
In [5]: #SUMMARY OF TOTAL RESULTS

#lifetime totals
print('Total number of customers:', df_customers['company_id'].count())
print('Total number of invoices sent:', df_customers['invoices'].sum())
print('Total number of payments received:', df_customers['payments'].sum())
print('Total number of projects:', df_customers['projects'].sum())
print('Total number of hours logged:', df_customers['total_hours'].sum())
print('Total number of users:', df_customers['total_users'].sum())

#Print mean, max, and min for continuous columns
col_columns = ['total_users', 'invoices', 'projects', 'payments', 'total_hours']
display(df_customers[['total_users', 'invoices', 'projects', 'payments', 'total_hou

Total number of customers: 15028
Total number of invoices sent: 5054179
Total number of payments received: 4275499
Total number of projects: 4511062
Total number of hours logged: 419087684.2500009
Total number of users: 113558

TABLE SUMMARY
total_users  invoices  projects  payments  total_hours
count  15028.000000    15028.000000    15028.000000    15028.000000    1.502800e+04
mean      7556428      336.317474      300.177136      284.502196      2.788712e+04
std      20082364      1232.339742      980.476202      1137391476      9.549754e+04
min         0.000000         0.000000         0.000000         0.000000         0.000000e+00
25%         1.000000         0.000000         16.000000         0.000000         7.269975e+02
50%        2.000000         2.700000         62.000000         10.000000         3.823510e+03
75%        7.000000        235.000000        227.000000        183.000000         1.799153e+04
max       700.00000     75885.000000     58470.000000     71180.000000     5.676045e+06

Using the customer data, we will plot some insights related to the customer's and the type of accounts
they hold and their industry, invoices, and customer conversions over time.
```

```
In [10]: #Plot number of customers vs purpose
purpose_counts = df_customers.groupby('purpose')[['company_id']].count().reset_index()
fig = px.bar(purpose_counts, x='purpose', y='customer counts', title='Count of custo

fig.update_layout(
    autosize=False,
    width=800,
    height=400)
# fig.show()
plt.plot(fig)

#Plot number of customers vs primary industry grouped
industry_counts = df_customers.groupby(['primary_industry_grouped', 'purpose'])[['comp
industry_counts.columns = ['industry', 'purpose', 'customer counts']
fig = px.bar(industry_counts, x='industry', y='customer counts', color='purpose', tit
fig.update_layout(
    autosize=False,
    width=900,
    height=400)
# fig.show()
plt.plot(fig)

#Plot total number of hours vs primary industry grouped
df = df_customers.groupby(['primary_industry_grouped'])[['total_users', 'invoices', 't
fig = px.bar(df, x='primary_industry_grouped', y='total_hours', title='Total hours lo
plt.plot(fig)

#Plot mean number of hours vs primary industry grouped
df = df_customers.groupby(['primary_industry_grouped'])[['total_users', 'invoices', 't
df.columns = ['primary_industry_grouped', 'avg_no_users', 'avg_invoices', 'avg_no_proj
fig = px.bar(df, x='primary_industry_grouped', y='avg_hours', title='Average hours lo
plt.plot(fig)

#Plot histogram of total users per company
plt.figure(figsize=(10,8))
bins = np.arange(0,40,1)
plt.hist(np.clip(df_customers['total_users'], bins[0], bins[-1]), bins=bins)
plt.xlabel('Total users in each company')
plt.ylabel('Counts')
plt.title('Histogram of total users per company')
plt.show()

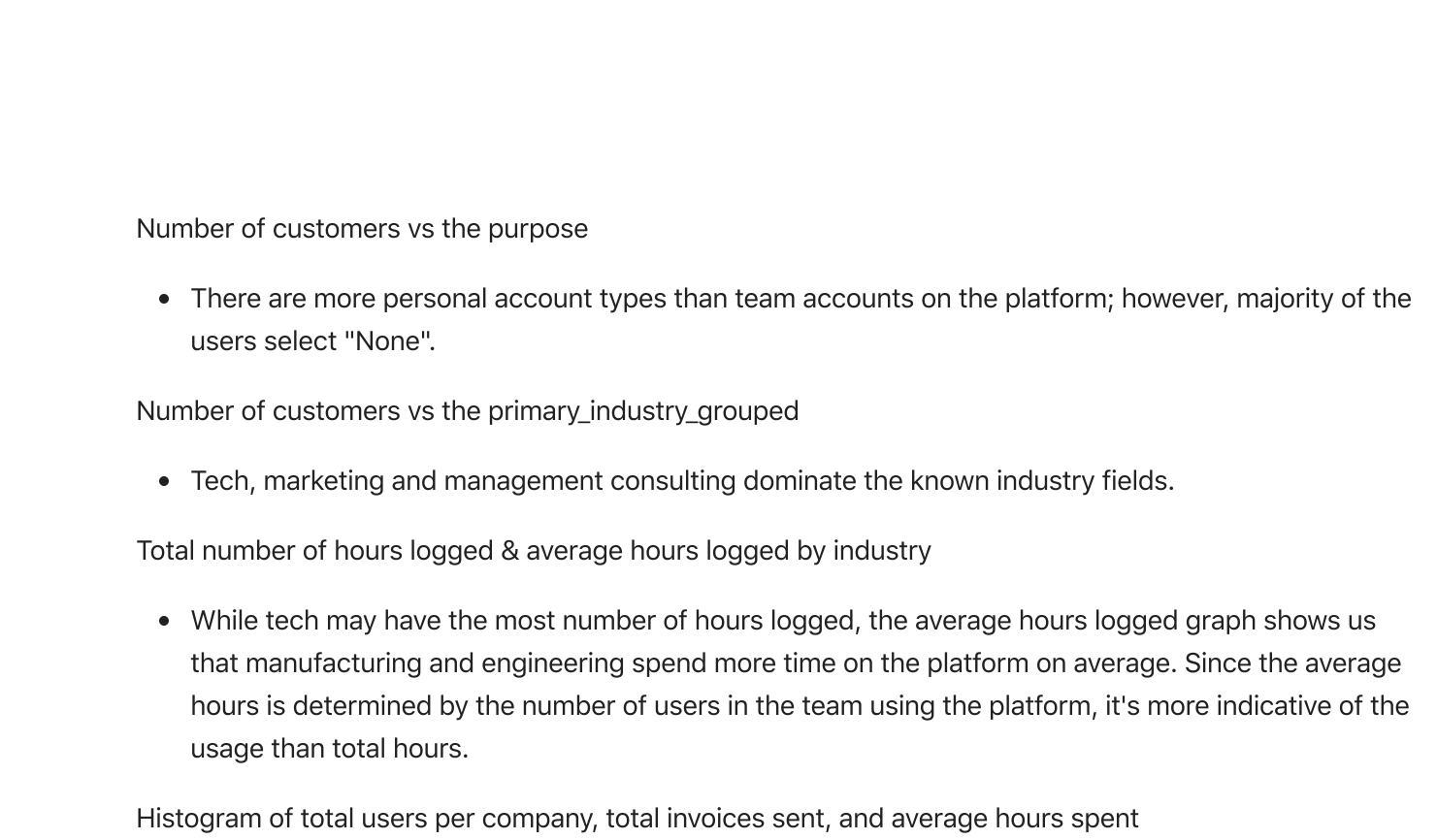
#Plot histogram of number of invoices sent per company
plt.figure(figsize=(10,8))
bins1=np.linspace(df_customers['invoices'].min(),3000, 8) #outlier in this column skew
plt.hist(df_customers['invoices'],bins=bins1, edgecolor='k',linewidth=1)
plt.xlabel('Number of invoices')
plt.ylabel('Counts')
plt.title('Number of invoices created across an account's entire lifetime')
plt.show()

#Plot histogram of average hours spent per user per company
plt.figure(figsize=(10,8))
bins = np.arange(0,4000,1000)
df_customers['average_hours'] = df_customers['total_hours']/df_customers['total_users']
plt.hist(np.clip(df_customers['average_hours'], bins[0], bins[-1]), bins=bins)
plt.xlabel('Average number of hours')
plt.ylabel('Counts')
plt.title('Average hours per user, logged across an account's lifetime')
plt.show()

#Plot a time series of the number of companies that converted to paying customers
df = df_customers.groupby(pd.Grouper(key='converted_at_date', freq='M'))[['company_id'
df.columns = ['converted_at_date', 'number_of_companies']
fig = px.line(df, x='converted_at_date', y='number_of_companies', title='Customer con
plt.plot(fig)

#Plot a time series of the number of companies that converted to paying customers (ag
df = df_customers.groupby(pd.Grouper(key='converted_at_date', freq='M'))[['company_id'
df.columns = ['converted_at_date', 'purpose', 'number_of_companies']
fig = px.line(df, x='converted_at_date', y='number_of_companies', color='purpose', tit
plt.plot(fig)

#Plot a time series of the number of companies that converted to yearly vs monthly pay
df = df_customers.groupby(pd.Grouper(key='converted_at_date', freq='M'))[['company_id'
df.columns = ['converted_at_date', 'billing_frequency', 'number_of_companies']
fig = px.line(df, x='converted_at_date', y='number_of_companies', color='billing_freq
plt.plot(fig)
```



Number of customers vs the purpose

- There are more personal account types than team accounts on the platform; however, majority of the users select "None".

Number of customers vs the primary industry grouped

- Tech, marketing and management consulting dominate the known industry fields.

Total number of hours logged & average hours logged by industry

- While tech may have the most number of hours logged, the average hours logged graph shows that manufacturing and engineering spend more time on the platform on average. Since the average hours is determined by the number of users more the time using the platform, it's more indicative of the usage than total hours.

Histogram of total users per company, total invoices sent, and average hours spent

- In all these graphs above, we see a long-tailed distribution, which makes sense because the distribution can't be below 0 and we expect anomalies (like time spent) to skew the data.

Time series of the number of companies that converted to paying customers by purpose, and by billing frequency

- The customer conversion has substantially grown since 2006 - 2022. It appears that in Feb 2018, the company introduced the "purpose" feature which lets users determine whether their account is for personal or team uses. We can see that the personal account type is growing faster than the team account type. From the last graph, we see that the monthly billing frequency is much more popular than the yearly billing frequency.

Next, we will look at the invoices and expenses dataset & merge the two together.

```
In [7]: ##### EXPENSES EDA #####

#Re-format expenses dataframe and add a few new features
billable_expenses_aggregated=df_expenses.loc[df_expenses['is_billable']!=1,]['company_
billable_expenses_aggregated=billable_expenses_aggregated.rename(columns=['total_cost
nonbillable_expenses_aggregated=df_expenses.loc[df_expenses['is_billable']!=0,]['compa
nonbillable_expenses_aggregated=nonbillable_expenses_aggregated.rename(columns=['total
expenses_aggregated=pd.merge(billable_expenses_aggregated,nonbillable_expenses_aggre
expenses_aggregated=billable_cost_percentag']=100.0*expenses_aggregated['total_bill
expenses_aggregated['nonbillable_cost_percentag']=100.0*expenses_aggregated['total_

#Drop negative total billables/nonbillables
expenses_aggregated=expenses_aggregated.loc[(expenses_aggregated['billable_cost_perce
expenses_aggregated['nonbillable_cost_per

display(expenses_aggregated)

#We should convert all billable costs to a standardized currency (USD) in order to co
#to do this, we should merge expenses_aggregated with the invoices dataframe in order

company_id  invoice_id  total_billable_cost  total_billable_units  total_nonbillable_cost  total_nonbillable

0          1417          0          14634.19          85.00          34375.00
1          1628          0          789.54          7.00          0.00
2          2327          0          18134.42          129.00          0.00
3          2327          1          105245.49          40349.87          0.00
4          2327          953657          1280.00          1.00          0.00
...          ...          ...          ...          ...          ...
390924      1411306          0          0.00          0.00          44.82
390925      1412171          0          0.00          0.00          262.83
390926      1414059          0          0.00          0.00          172.00
390927      1414483          0          0.00          0.00          265500.00
390928      1416738          0          0.00          0.00          2855.34

390716 rows x 8 columns
```

```
In [14]: ##### INVOICES EDA #####

#Merge in the expenses_aggregated dataframe on company_id and invoice_id
df_invoices_and_expenses=pd.merge(df_invoices, expenses_aggregated, on=['company_id',
# Number of invoice data rows that merged with the expenses_aggregated dataframe
print('Percentage of invoices with merged expenses information:',str(round(100.0*
display(df_invoices_and_expenses.head())

Percentage of invoice data with merged expenses information: 8.4%

company_id  invoice_id  recurring_invoice_id  invoice_created_at  invoice_sent_at  invoice_state  invoice_amo

0          612957      25740109          NaN          2020-10-16 01:27:54 UTC          2020-10-16 01:35:40 UTC          open          384415
1          274327      14430600          NaN          2017-11-07 19:19:20 UTC          2017-11-07 19:20:06 UTC          open          22063
2          359036      33084769          NaN          2022-07-06 03:51:28 UTC          2022-07-04 05:00:00 UTC          open          28518
3          40237      10979856          NaN          2016-10-25 12:31:57 UTC          2016-10-27 09:20:33 UTC          open          7500
4          456494      3092832          NaN          2022-01-11 13:38:13 UTC          2022-01-11 13:36:16 UTC          open          180000

5 rows x 22 columns
```

Using the invoices data, we can take a look at when companies were last active by using the data they last sent an invoice (we could have also used payments). Since we don't have data with the times they logged in, we are going to assume the invoices give us an idea about their activity.

```
In [19]: #Plot a graph of companies vs the data of the last invoice sent
df = df_invoices.groupby('company_id')[['invoice_created_at']].max().reset_index()
df['invoice_created_at_date'] = pd.to_datetime(df['invoice_created_at'], utc=True)
df = df.groupby(pd.Grouper(key='invoice_created_at_date', freq='Y'))[['company_id']]
df.columns = ['invoice_created_at_date', 'number_of_companies']

fig = px.bar(df, x='invoice_created_at_date', y='number_of_companies', title='Number o
plt.plot(fig)

plt.plot(fig)
```

In the next section, I will take a look at customer conversion and churn to see if we can develop any new insights using the data.

Customer Conversion

Since we don't have any data about the customer's behavior prior to the date they converted (ex. when they first sent the platform), we have to make some assumptions. We are going to assume that the date they sent an invoice (prior to their conversion date) is their first day on the platform.

By calculating the days between their first invoice and the conversion date, we can calculate the length of time before a customer converts & plot it on a histogram to observe the distribution of this data.

```
In [10]: #1) Customer Conversion

conversion_df=pd.merge(df_customers[['company_id','converted_at_date']],groupby('comp
df_invoices_and_expenses.loc[conversion_df['invoice_created_at_date']<conversion_d
converted_df['days_to_convert']=conversion_df['converted_at_date'] - converted_df['inv
# display(converted_df)

#Plot histogram of total users per company
plt.figure(figsize=(10,8))
bins = np.arange(0,200,1)
plt.hist(np.clip(converted_df['days_to_convert'], bins[0], bins[-1]), bins=bins)
plt.xlabel('Days To Convert')
plt.ylabel('Counts')
plt.title('Histogram of Days To Convert')
plt.show()

print('Numbers below are for converted customers who created an invoice before convert
print('Percentage of Users that Convert by the 7-Day mark:',round(100.0*converted_d
print('Percentage of Users that Convert by the 14-Day mark:',round(100.0*converted_d
print('Percentage of Users that Convert by the 30-Day mark:',round(100.0*converted_d
print('Percentage of Users that Convert by the 60-Day mark:',round(100.0*converted_d
print('Percentage of Users that Convert by the 21-Day mark:',round(100.0*converted_d

<ipython-input-10-22c98ccc02>:8: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.
Try using loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/ue
r_guide/indexing.html#returning-a-view-versus-a-copy
```


Numbers below are for converted customers who created an invoice before converting:
Percentage of Users that Convert by the 7-Day mark: 15.95
Percentage of Users that Convert by the 14-Day mark: 29.13
Percentage of Users that Convert by the 21-Day mark: 45.03

Percentage of Users that Convert by the 30-Day mark: 72.53
Percentage of Users that Convert by the 60-Day mark: 86.11
Percentage of Users that Convert by the 365-Day mark: 94.43

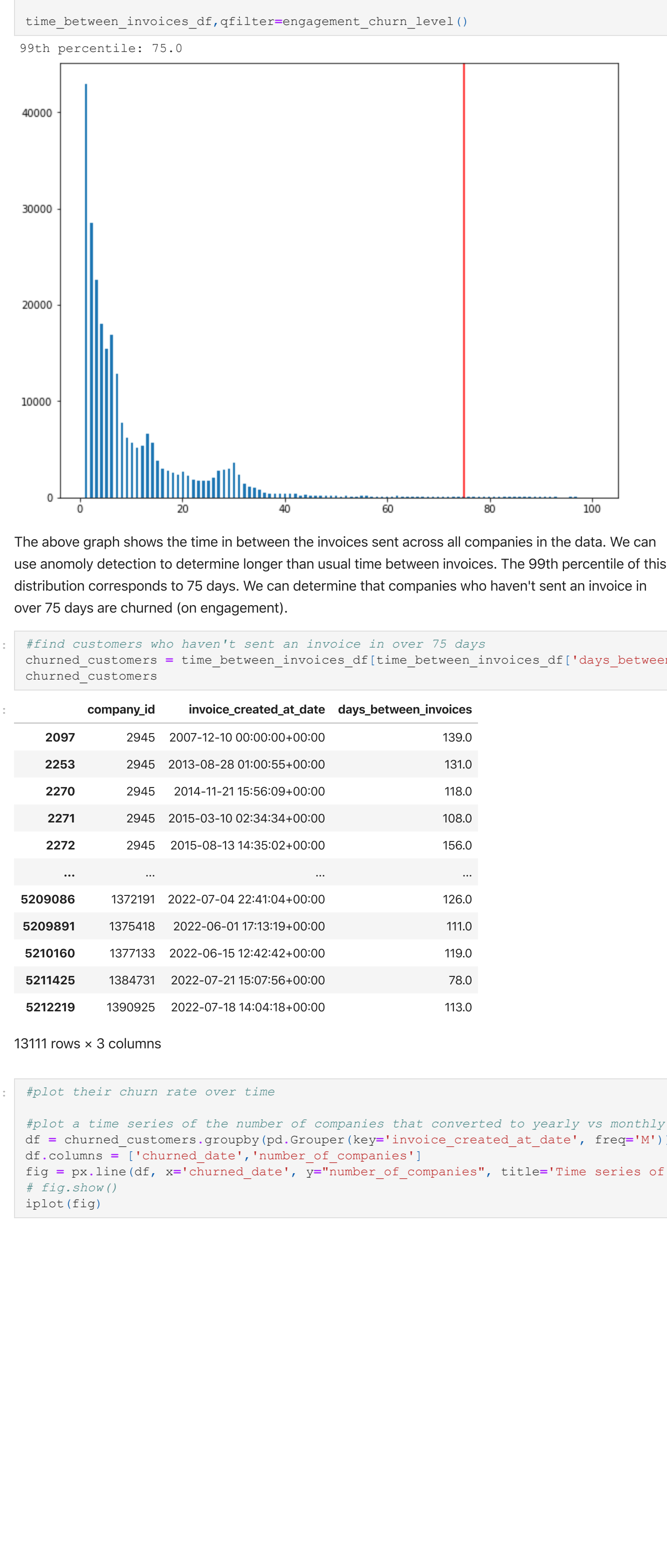
The graph above shows data from customers who converted and their days to conversion. Again, since we don't have the days when they joined the platform as a free user, we are assuming that the first invoice date is the creation date. From the graph above, we see that for users who converted, 72% of them did so before 30 days. This is an interesting result that makes sense because Harvest has a 30 day free trial!

Unfortunately, we can't determine features that predict a customer will convert because we do not have the data for customers who churned before becoming a paying customer.

Customer Engagement Churn

Technically, we can't predict "churn" in the sense that a customer stopped paying and using the platform, but we can look at engagement churn. We have to make the following assumptions:

- Assumption 1: Churn is defined as a customer who haven't created an invoice in a specific number of days (to be determined)
- Assumption 2: Customers who churn on engagement don't return



The above graph shows the time in between the invoices sent across all companies in the data. We can use anomaly detection to determine longer than usual time between invoices. The 99th percentile of this distribution corresponds to 75 days. We can determine that companies who haven't sent an invoice in over 75 days are churned (on engagement).

```
In [12]: #find customers who haven't sent an invoice in over 75 days
churned_customers = time_between_invoices_df[time_between_invoices_df['days_between_invoices'] > 75]

Out[12]:
```

	company_id	invoice_created_at_date	days_between_invoices
	2097	2945 2007-12-10 00:00:00+00:00	139.0
	2253	2945 2013-08-28 01:00:55+00:00	131.0
	2270	2945 2014-11-21 15:56:09+00:00	118.0
	2271	2945 2015-03-10 02:34:34+00:00	108.0
	2272	2945 2015-08-13 14:35:02+00:00	156.0

	5209086	1372191 2022-07-04 22:41:04+00:00	126.0
	5209891	1375418 2022-06-01 17:13:19+00:00	111.0
	5210160	1377133 2022-06-15 12:42:42+00:00	119.0
	5211425	1384731 2022-07-21 15:07:56+00:00	78.0
	5212219	1390925 2022-07-18 14:04:18+00:00	113.0

13111 rows x 3 columns

```
In [21]: #plot their churn rate over time
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```