

# Capstone Project 1: MuscleHub AB Test

## Step 1: Get started with SQL

Like most businesses, Janet keeps her data in a SQL database. Normally, you'd download the data from her database to a csv file, and then load it into a Jupyter Notebook using Pandas.

For this project, you'll have to access SQL in a slightly different way. You'll be using a special Codecademy library that lets you type SQL queries directly into this Jupyter notebook. You'll have pass each SQL query as an argument to a function called `sql_query`. Each query will return a Pandas DataFrame. Here's an example:

```
In [4]: # This import only needs to happen once, at the beginning of the notebook
from codecademySQL import sql_query
```

```
In [5]: # Here's an example of a query that just displays some data
sql_query('''
SELECT *
FROM visits
LIMIT 5
''')
```

```
Out[5]:
```

|   | index | first_name | last_name | email                         | gender | visit_date |
|---|-------|------------|-----------|-------------------------------|--------|------------|
| 0 | 0     | Karen      | Manning   | Karen.Manning@gmail.com       | female | 5-1-17     |
| 1 | 1     | Annette    | Boone     | AB9982@gmail.com              | female | 5-1-17     |
| 2 | 2     | Salvador   | Merritt   | SalvadorMerritt12@outlook.com | male   | 5-1-17     |
| 3 | 3     | Martha     | Maxwell   | Martha.Maxwell@gmail.com      | female | 5-1-17     |
| 4 | 4     | Andre      | Mayer     | AndreMayer90@gmail.com        | male   | 5-1-17     |

```
In [6]: # Here's an example where we save the data to a DataFrame
df = sql_query('''
SELECT *
FROM applications
LIMIT 5
''')
```

## Step 2: Get your dataset

Let's get started!

Janet of MuscleHub has a SQLite database, which contains several tables that will be helpful to you in this investigation:

- `visits` contains information about potential gym customers who have visited MuscleHub
- `fitness_tests` contains information about potential customers in "Group A", who were given a fitness test
- `applications` contains information about any potential customers (both "Group A" and "Group B") who filled out an application. Not everyone in `visits` will have filled out an application.
- `purchases` contains information about customers who purchased a membership to MuscleHub.

Use the space below to examine each table.

```
In [7]: # Examine visits here
sql_query('''SELECT * FROM visits LIMIT 5''')
```

```
Out[7]:
```

|   | index | first_name | last_name | email                         | gender | visit_date |
|---|-------|------------|-----------|-------------------------------|--------|------------|
| 0 | 0     | Karen      | Manning   | Karen.Manning@gmail.com       | female | 5-1-17     |
| 1 | 1     | Annette    | Boone     | AB9982@gmail.com              | female | 5-1-17     |
| 2 | 2     | Salvador   | Merritt   | SalvadorMerritt12@outlook.com | male   | 5-1-17     |
| 3 | 3     | Martha     | Maxwell   | Martha.Maxwell@gmail.com      | female | 5-1-17     |
| 4 | 4     | Andre      | Mayer     | AndreMayer90@gmail.com        | male   | 5-1-17     |

```
In [9]: # Examine fitness_tests here
sql_query('''SELECT * FROM fitness_tests LIMIT 5''')
```

```
Out[9]:
```

|   | index | first_name | last_name | email                  | gender | fitness_test_date |
|---|-------|------------|-----------|------------------------|--------|-------------------|
| 0 | 0     | Kim        | Walter    | KimWalter58@gmail.com  | female | 2017-07-03        |
| 1 | 1     | Tom        | Webster   | TW3857@gmail.com       | male   | 2017-07-02        |
| 2 | 2     | Marcus     | Bauer     | Marcus.Bauer@gmail.com | male   | 2017-07-01        |
| 3 | 3     | Roberta    | Best      | R86305@hotmail.com     | female | 2017-07-02        |
| 4 | 4     | Carrie     | Francis   | CF1896@hotmail.com     | female | 2017-07-05        |

```
In [10]: # Examine applications here
sql_query('''SELECT * FROM applications LIMIT 5''')
```

```
Out[10]:
```

|   | index | first_name | last_name | email                   | gender | application_date |
|---|-------|------------|-----------|-------------------------|--------|------------------|
| 0 | 0     | Roy        | Abbott    | RoyAbbott32@gmail.com   | male   | 2017-08-12       |
| 1 | 1     | Agnes      | Acevedo   | AgnesAcevedo1@gmail.com | female | 2017-09-29       |
| 2 | 2     | Roberta    | Acevedo   | RA8063@gmail.com        | female | 2017-09-15       |
| 3 | 3     | Darren     | Acosta    | DAcosta1996@hotmail.com | male   | 2017-07-26       |
| 4 | 4     | Vernon     | Acosta    | VAcosta1975@gmail.com   | male   | 2017-07-14       |

```
In [11]: # Examine purchases here
sql_query('''SELECT * FROM purchases LIMIT 5''')
```

```
Out[11]:
```

|   | index | first_name | last_name | email                   | gender | purchase_date |
|---|-------|------------|-----------|-------------------------|--------|---------------|
| 0 | 0     | Roy        | Abbott    | RoyAbbott32@gmail.com   | male   | 2017-08-18    |
| 1 | 1     | Roberta    | Acevedo   | RA8063@gmail.com        | female | 2017-09-16    |
| 2 | 2     | Vernon     | Acosta    | VAcosta1975@gmail.com   | male   | 2017-07-20    |
| 3 | 3     | Darren     | Acosta    | DAcosta1996@hotmail.com | male   | 2017-07-27    |
| 4 | 4     | Dawn       | Adkins    | Dawn.Adkins@gmail.com   | female | 2017-08-24    |

We'd like to download a giant DataFrame containing all of this data. You'll need to write a query that does the following things:

- Not all visits in `visits` occurred during the A/B test. You'll only want to pull data where `visit_date` is on or after 7-1-17.
- You'll want to perform a series of `LEFT JOIN` commands to combine the four tables that we care about. You'll need to perform the joins on `first_name`, `last_name`, and `email`. Pull the following columns:

- `visits.first_name`
- `visits.last_name`
- `visits.gender`
- `visits.email`
- `visits.visit_date`
- `fitness_tests.fitness_test_date`
- `applications.application_date`
- `purchases.purchase_date`

Save the result of this query to a variable called `df`.

Hint: your result should have 5004 rows. Does it?

```
In [33]: # Since column names are all the same across all tables, USING(column_1, column_2) can be
# used instead.
df = sql_query('''
SELECT v.first_name, v.last_name, v.gender, v.email, v.visit_date, f.fitness_test_date,
a.application_date, p.purchase_date
FROM visits AS v
LEFT JOIN fitness_tests AS f
USING (first_name, last_name, email)
LEFT JOIN applications AS a
ON v.first_name = a.first_name AND v.last_name = a.last_name AND v.email = a.email
LEFT JOIN purchases AS p
ON v.first_name = p.first_name AND v.last_name = p.last_name AND v.email = p.email
WHERE visit_date >= '7-1-17'
''')
df
```

```
Out[33]:
```

|      |         | first_name | last_name | gender                     | email  | visit_date | fitness_test_date | application_date | purchase_date |
|------|---------|------------|-----------|----------------------------|--------|------------|-------------------|------------------|---------------|
| 0    | Kim     | Walter     | female    | KimWalter58@gmail.com      | 7-1-17 | 2017-07-03 | None              | None             |               |
| 1    | Tom     | Webster    | male      | TW3857@gmail.com           | 7-1-17 | 2017-07-02 | None              | None             |               |
| 2    | Edw  rd | Bowen      | male      | Edward.Bowen@gmail.com     | 7-1-17 | None       | 2017-07-04        | 2017-07-04       |               |
| 3    | Marcus  | Bauer      | male      | Marcus.Bauer@gmail.com     | 7-1-17 | 2017-07-01 | 2017-07-03        | 2017-07-05       |               |
| 4    | Roberta | Best       | female    | R86305@hotmail.com         | 7-1-17 | 2017-07-02 | None              | None             |               |
| ...  | ...     | ...        | ...       | ...                        | ...    | ...        | ...               | ...              | ...           |
| 4999 | Rachel  | Hensley    | female    | RachelHensley38@gmail.com  | 9-9-17 | None       | None              | None             |               |
| 5000 | Leon    | Harmon     | male      | Leon.Harmon@gmail.com      | 9-9-17 | 2017-09-15 | None              | None             |               |
| 5001 | Andy    | Pratt      | male      | AndyPratt27@gmail.com      | 9-9-17 | 2017-09-15 | None              | None             |               |
| 5002 | Ruben   | Nielsen    | male      | RubenNielsen93@hotmail.com | 9-9-17 | None       | 2017-09-13        | None             |               |
| 5003 | Charles | Carver     | male      | CC2490@gmail.com           | 9-9-17 | 2017-09-12 | None              | None             |               |

5004 rows x 8 columns

## Step 3: Investigate the A and B groups

We have some data to work with! Import the following modules so that we can start doing analysis:

- `import pandas as pd`
- `from matplotlib import pyplot as plt`

```
In [25]: import pandas as pd
from matplotlib import pyplot as plt
```

We're going to add some columns to `df` to help us with our analysis.

Start by adding a column called `ab_test_group`. It should be `A` if `fitness_test_date` is not `None`, and `B` if `fitness_test_date` is `None`.

```
In [37]: df['ab_test_group'] = df.fitness_test_date.apply(lambda x: 'A' if pd.notnull(x) else 'B')
df
```

```
Out[37]:
```

|      | first_name | last_name | gender | email                      | visit_date | fitness_test_date | application_date | purchase_date | ab_test_group |
|------|------------|-----------|--------|----------------------------|------------|-------------------|------------------|---------------|---------------|
| 0    | Kim        | Walter    | female | KimWalter58@gmail.com      | 7-1-17     | 2017-07-03        | None             | None          | A             |
| 1    | Tom        | Webster   | male   | TW3857@gmail.com           | 7-1-17     | 2017-07-02        | None             | None          | A             |
| 2    | Edward     | Bowen     | male   | Edward.Bowen@gmail.com     | 7-1-17     | None              | 2017-07-04       | 2017-07-04    | B             |
| 3    | Marcus     | Bauer     | male   | Marcus.Bauer@gmail.com     | 7-1-17     | 2017-07-01        | 2017-07-03       | 2017-07-05    | A             |
| 4    | Roberta    | Best      | female | R86305@hotmail.com         | 7-1-17     | 2017-07-02        | None             | None          | A             |
| ...  | ...        | ...       | ...    | ...                        | ...        | ...               | ...              | ...           | ...           |
| 4999 | Rachel     | Hensley   | female | RachelHensley38@gmail.com  | 9-9-17     | None              | None             | None          | B             |
| 5000 | Leon       | Harmon    | male   | Leon.Harmon@gmail.com      | 9-9-17     | 2017-09-15        | None             | None          | A             |
| 5001 | Andy       | Pratt     | male   | AndyPratt27@gmail.com      | 9-9-17     | 2017-09-15        | None             | None          | A             |
| 5002 | Ruben      | Nielsen   | male   | RubenNielsen93@hotmail.com | 9-9-17     | None              | 2017-09-13       | None          | B             |
| 5003 | Charles    | Carver    | male   | CC2490@gmail.com           | 9-9-17     | 2017-09-12        | None             | None          | A             |

5004 rows x 9 columns

Let's do a quick sanity check that Janet split her visitors such that about half are in A and half are in B.

Start by using `groupby` to count how many users are in each `ab_test_group`. Save the results to `ab_counts`.

```
In [39]: ab_counts = df.groupby('ab_test_group').email.count().reset_index()
ab_counts
```

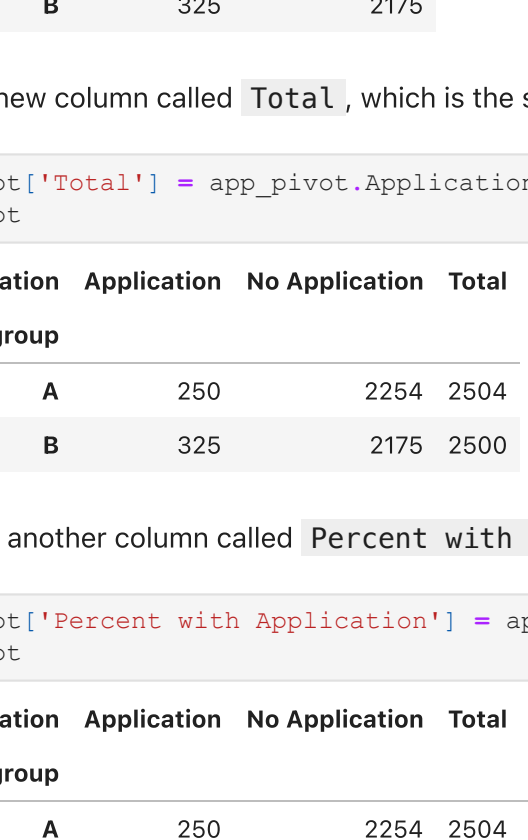
```
Out[39]:
```

|   | ab_test_group | email |
|---|---------------|-------|
| 0 | A             | 2504  |
| 1 | B             | 2500  |

We'll want to include this information in our presentation. Let's create a pie chart using `plt.pie`. Make sure to include:

- Use `plt.axis('equal')` so that your pie chart looks nice
- Add a legend labeling 'A' and 'B'
- Use `autopct` to label the percentage of each group
- Save your figure as `ab_test_pie_chart.png`

```
In [42]: plt.pie(ab_counts.email.values, autopct='%d%%')
plt.axis('equal')
plt.legend(ab_counts.ab_test_group.values)
plt.savefig('ab_test_pie_chart.png')
```



## Step 4: Who picks up an application?

Recall that the sign-up process for MuscleHub has several steps:

- Take a fitness test with a personal trainer (Only Group A)
- Fill out an application for the gym
- Send in their payment for their first month's membership

Let's examine how many people make it to Step 2, filling out an application.

Start by creating a new column in `df` called `is_application` which is `Application` if `application_date` is not `None` and `No Application`, otherwise.

```
In [44]: df['is_application'] = df.application_date.apply(lambda x: 'Application' if pd.notnull(x) \
else 'No Application')
df
```

```
Out[44]:
```

|      | first_name | last_name | gender | email                      | visit_date | fitness_test_date | application_date | purchase_date | ab_test_group |
|------|------------|-----------|--------|----------------------------|------------|-------------------|------------------|---------------|---------------|
| 0    | Kim        | Walter    | female | KimWalter58@gmail.com      | 7-1-17     | 2017-07-03        | None             | None          | A             |
| 1    | Tom        | Webster   | male   | TW3857@gmail.com           | 7-1-17     | 2017-07-02        | None             | None          | A             |
| 2    | Edward     | Bowen     | male   | Edward.Bowen@gmail.com     | 7-1-17     | None              | 2017-07-04       | 2017-07-04    | B             |
| 3    | Marcus     | Bauer     | male   | Marcus.Bauer@gmail.com     | 7-1-17     | 2017-07-01        | 2017-07-03       | 2017-07-05    | A             |
| 4    | Roberta    | Best      | female | R86305@hotmail.com         | 7-1-17     | 2017-07-02        | None             | None          | A             |
| ...  | ...        | ...       | ...    | ...                        | ...        | ...               | ...              | ...           | ...           |
| 4999 | Rachel     | Hensley   | female | RachelHensley38@gmail.com  | 9-9-17     | None              | None             | None          | B             |
| 5000 | Leon       | Harmon    | male   | Leon.Harmon@gmail.com      | 9-9-17     | 2017-09-15        | None             | None          | A             |
| 5001 | Andy       | Pratt     | male   | AndyPratt27@gmail.com      | 9-9-17     | 2017-09-15        | None             | None          | A             |
| 5002 | Ruben      | Nielsen   | male   | RubenNielsen93@hotmail.com | 9-9-17     | None              | 2017-09-13       | None          | B             |
| 5003 | Charles    | Carver    | male   | CC2490@gmail.com           | 9-9-17     | 2017-09-12        | None             | None          | A             |

5004 rows x 10 columns

Now, using `groupby`, count how many people from Group A and Group B either do or don't pick up an application. You'll want to group by `ab_test_group` and `is_application`. Save this new DataFrame as `app_counts`:

```
In [46]: app_counts = df.groupby(['ab_test_group', 'is_application']).email.count().reset_index()
app_counts
```

```
Out[46]:
```

|   | ab_test_group | is_application | email |
|---|---------------|----------------|-------|
| 0 | A             | Application    | 250   |
| 1 | A             | No Application | 2254  |
| 2 | B             | Application    | 325   |
| 3 | B             | No Application | 2175  |

We're going to want to calculate the percent of people in each group who complete an application. It's going to be much easier to do this if we pivot `app_counts` such that:

- The `columns` are `ab_test_group`
- The `index` are `is_application` Perform this pivot and save it to the variable `app_pivot`. Remember to call `reset_index()` at the end of the pivot!

```
In [48]: app_pivot = app_counts.pivot(index='is_application',
                                   columns='ab_test_group', values='email')
app_pivot
```

```
Out[48]:
```

|               | is_application | Application | No Application |
|---------------|----------------|-------------|----------------|
| ab_test_group |                |             |                |
| A             |                | 250         | 2254           |
| B             |                | 325         | 2175           |

Define a new column called `Total`, which is the sum of `Application` and `No Application`.

```
In [50]: app_pivot['Total'] = app_pivot.Application + app_pivot['No Application']
app_pivot
```

```
Out[50]:
```

|               | is_application | Application | No Application | Total |
|---------------|----------------|-------------|----------------|-------|
| ab_test_group |                |             |                |       |
| A             |                | 250         | 2254           | 2504  |
| B             |                | 325         | 2175           | 2500  |

Calculate another column called `Percent with Application`, which is equal to `Application` divided by `Total`.

```
In [52]: app_pivot['Percent with Application'] = app_pivot.Application / app_pivot.Total
app_pivot
```

```
Out[52]:
```

|               | is_application | Application | No Application | Total | Percent with Application |
|---------------|----------------|-------------|----------------|-------|--------------------------|
| ab_test_group |                |             |                |       |                          |
| A             |                | 250         | 2254           | 2504  | 0.09984                  |
| B             |                | 325         | 2175           | 2500  | 0.13000                  |

It looks like more people from Group B turned in an application. Why might that be?

We need to know if this difference is statistically significant.

Choose a hypothesis tests, import it from `scipy` and perform it. Be sure to note the p-value. Is this result significant?

```
In [54]: # Import Chi Square Test
from scipy.stats import chi2_contingency
contingency2 = [[250, 2254], [325, 2175]]
_, pvalue2, _, _ = chi2_contingency(contingency2)
pvalue2
```

```
Out[54]: 0.0009647827600722304
```

## Step 5: Who purchases a membership?

Of those who picked up an application, how many purchased a membership?

Let's begin by adding a column to `df` called `is_member` which is `Member` if `purchase_date` is not `None`, and `Not Member` otherwise.

```
In [55]: df['is_member'] = df.purchase_date.apply(lambda x: \
'Member' if pd.notnull(x) else 'Not Member')
```

Now, let's create a DataFrame called `just_apps` the contains only people who picked up an application.

```
In [58]: just_apps = df[df.is_application == 'Application']
just_apps.head()
# Applied for a membership but have not completed the payment.
```

```
Out[58]:
```

|    | first_name | last_name | gender | email                   | visit_date | fitness_test_date | application_date | purchase_date | ab_test_group |
|----|------------|-----------|--------|-------------------------|------------|-------------------|------------------|---------------|---------------|
| 2  | Edward     | Bowen     | male   | Edward.Bowen@gmail.com  | 7-1-17     | None              | 2017-07-04       | 2017-07-04    | B             |
| 3  | Marcus     | Bauer     | male   | Marcus.Bauer@gmail.com  | 7-1-17     | 2017-07-01        | 2017-07-03       | 2017-07-05    | A             |
| 9  | Salvador   | Cardenas  | male   | SCardenas1980@gmail.com | 7-1-17     | 2017-07-07        | 2017-07-06       | None          | B             |
| 11 | Valerie    | Munoz     | female | VMunoz1998@gmail.com    | 7-1-17     | 2017-07-03        | 2017-07-05       | 2017-07-06    | A             |
| 35 | Michael    | Burks     | male   | MB9820@gmail.com        | 7-1-17     | None              | 2017-07-07       | 2017-07-13    | B             |

Great! Now, let's do a `groupby` to find out how many people in `just_apps` are and aren't members from each group. Follow the same process that we did in Step 4, including pivoting the data. You should end up with a DataFrame that looks like this:

```
Out[59]:
```

|   | is_member | ab_test_group | Member | Not Member | Total | Percent Purchase |
|---|-----------|---------------|--------|------------|-------|------------------|
| 0 | A         |               | ?      | ?          | ?     | ?                |
| 1 | B         |               | ?      | ?          | ?     | ?                |

Save your final DataFrame as `member_pivot`.

```
In [67]: member_groupby = just_apps.groupby(['is_member', 'ab_test_group']).email.count().reset_index()
member_pivot = member_groupby.pivot(columns='is_member',
                                   index='ab_test_group', values='email').reset_index()
member_pivot['Total'] = member_pivot.Member + member_pivot['Not Member']
member_pivot['Percent Purchase'] = member_pivot.Member / member_pivot.Total
member_pivot
```

```
Out[67]:
```

|   | is_member | ab_test_group | Member | Not Member | Total | Percent Purchase |
|---|-----------|---------------|--------|------------|-------|------------------|
| 0 | A         |               | 200    | 50         | 250   | 0.800000         |
| 1 | B         |               | 250    | 75         | 325   | 0.769231         |

It looks like people who took the fitness test were more likely to purchase a membership if they picked up an application. Why might that be?

Just like before, we need to know if this difference is statistically significant. Choose a hypothesis tests, import it from `scipy` and perform it. Be sure to note the p-value. Is this result significant?

```
In [69]: # A: Choose fitness test, B: no fitness test
# Choose Chi Square Test
contingency2 = [[200, 50], [250, 75]]
_, pvalue3, _, _ = chi2_contingency(contingency2)
pvalue3
```

```
Out[69]: 0.43258646051083327
```

Previously, we looked at what percent of people who **picked up applications** purchased memberships. What we really care about is what percentage of **all visitors** purchased memberships. Return to `df` and do a `groupby` to find out how many people in `df` are and aren't members from each group. Follow the same process that we did in Step 4, including pivoting the data. You should end up with a DataFrame that looks like this:

```
Out[70]:
```

|   | is_member | ab_test_group | Member | Not Member | Total | Percent Purchase |
|---|-----------|---------------|--------|------------|-------|------------------|
| 0 | A         |               | ?      | ?          | ?     | ?                |
| 1 | B         |               | ?      | ?          | ?     | ?                |

Save your final DataFrame as `final_member_pivot`.