



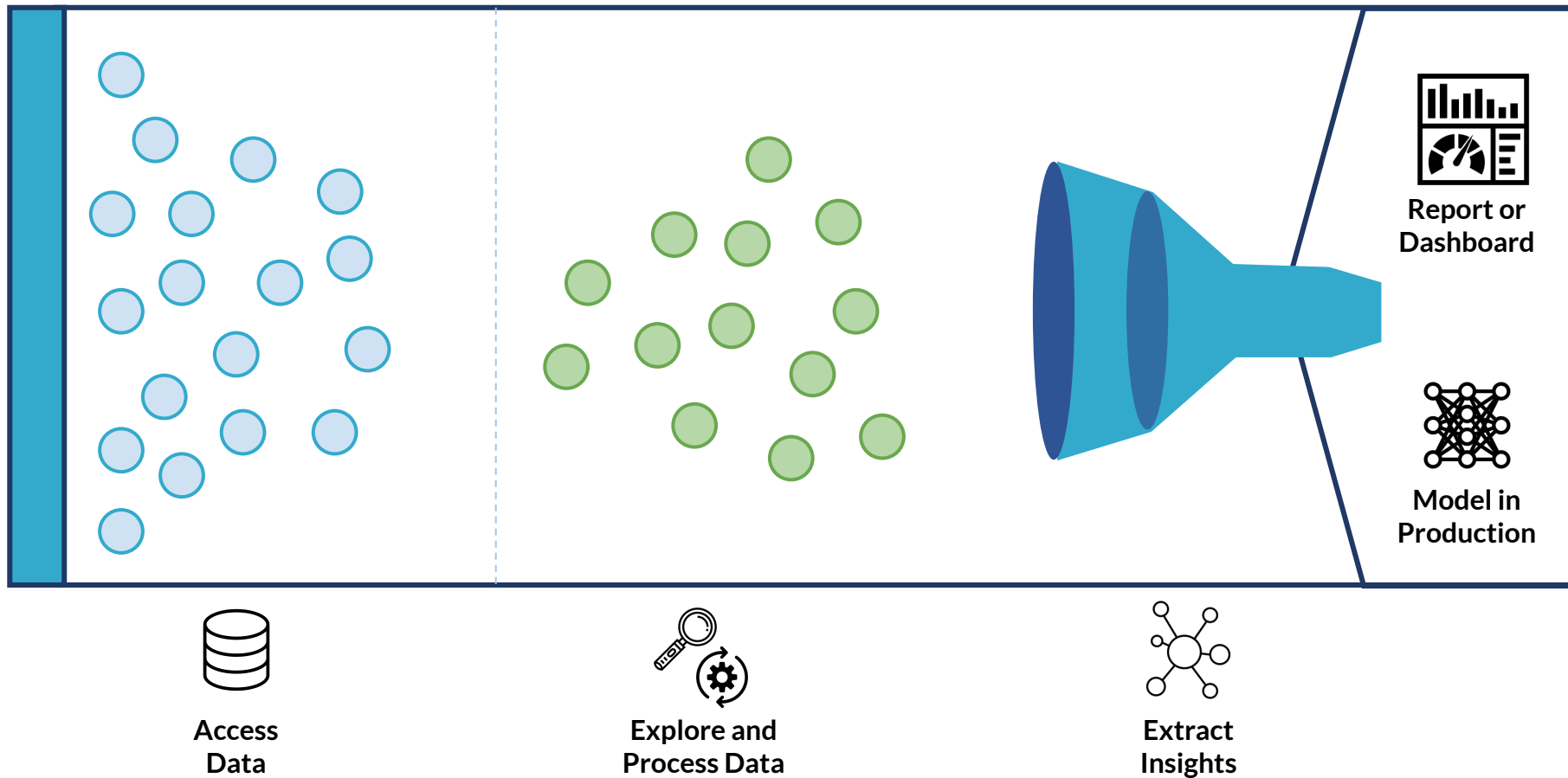
Live training: Cleaning Data in Python

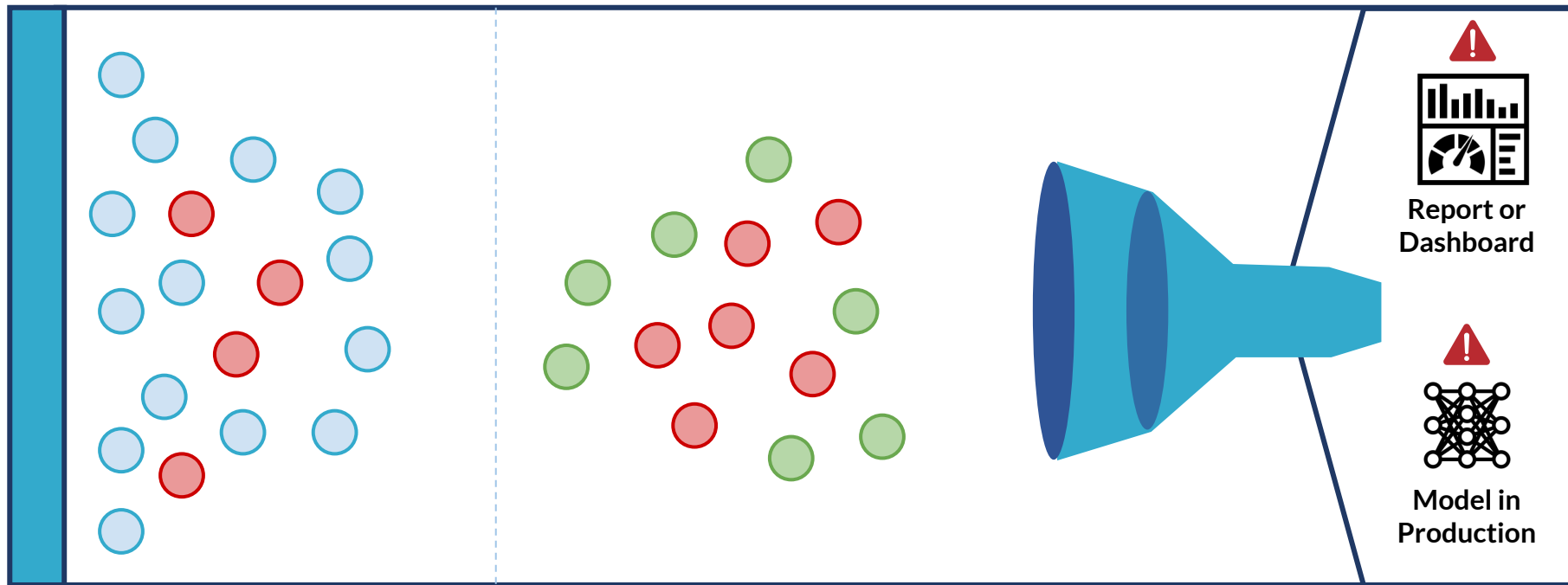




ADEL NEHME

Content Developer





Access
Data

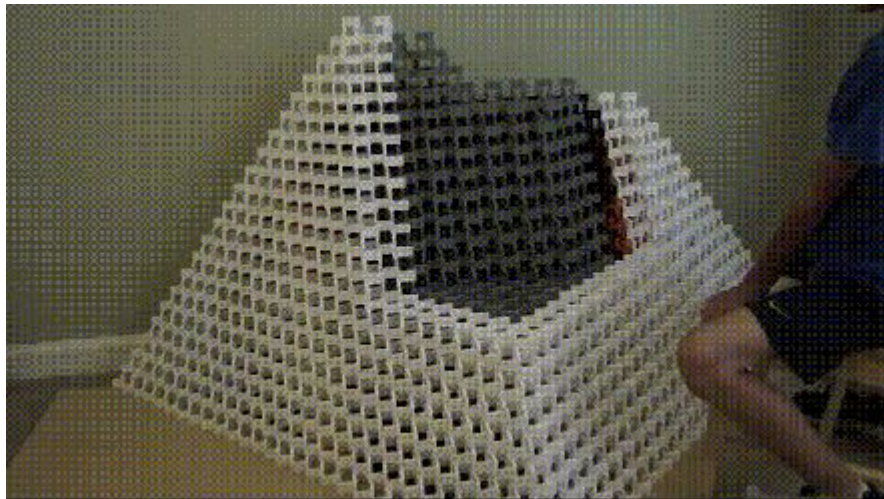


Explore and
Process Data



Extract
Insights

Why do we have dirty data?



Human error



Technical error



Airbnb data

listing_id: Unique identifier for a listing

Name: Description used for a listing

Host_id: Unique identifier for each host

Host_name: Name of host

Neighbourhood_full: Burrough and neighbourhood

Coordinates: Latitude, Longitude

Room_type: Type of room

Price: Price per night

Number of reviews: Number of reviews so far

Last_review: Date of last review

Reviews_per_month: # of reviews per month

Availability_365: Days available per year

Rating: Average rating (0 to 5)

Number_of_stays: Number of stays so far

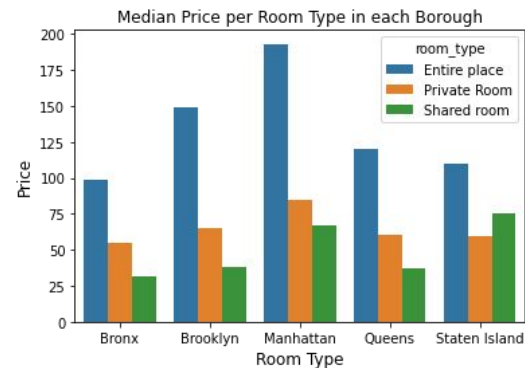
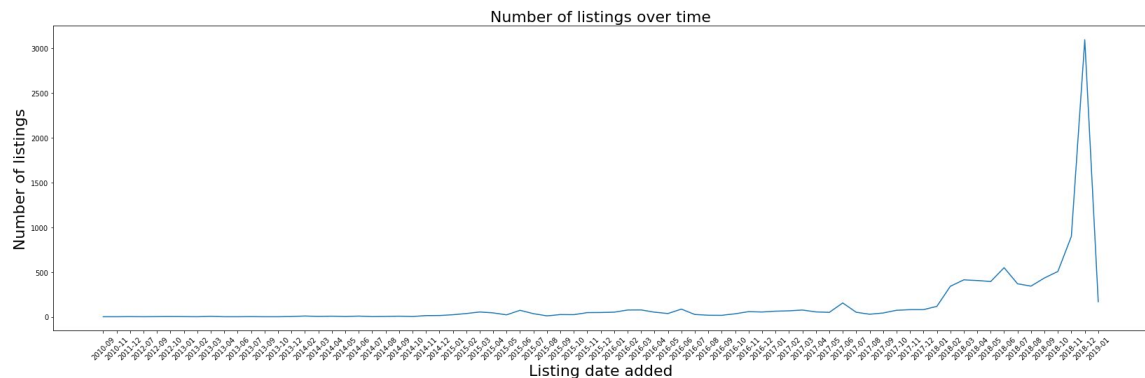
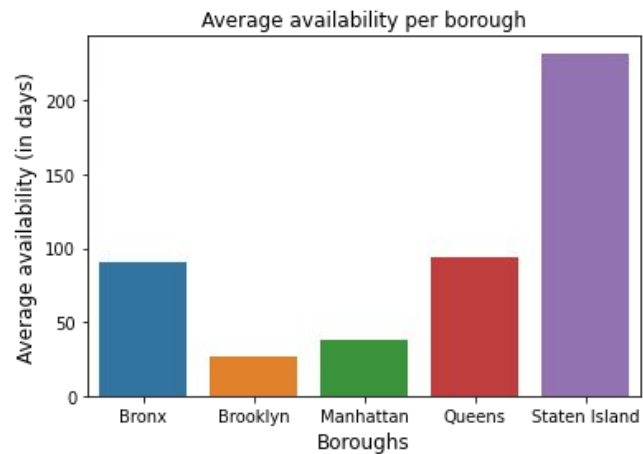
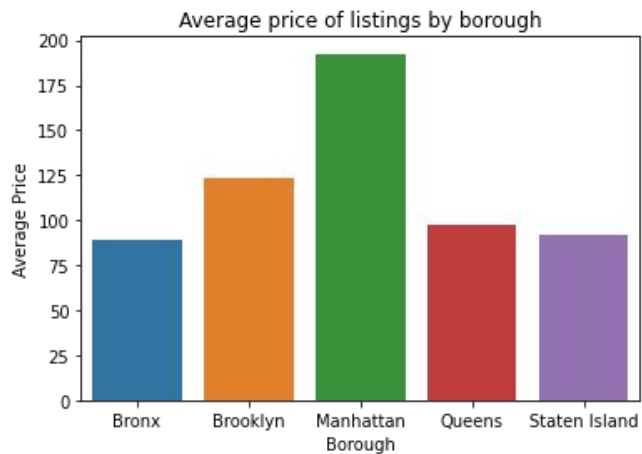
5_stars: Percentage of ratings that is 5_stars

Listing_added: Date listing added to site

Airbnb data featuring listings in New York



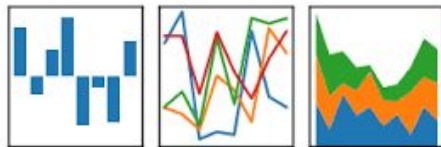
The end result





pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



Popular open source *data analysis tool* for tabular data

matplotlib



Open source *plotting* library for 2-D visualizations

Seaborn

Open source *plotting* library built on top of matplotlib



NumPy

Popular open source *computing tool* for arrays

missingno

Open source *plotting* library for missing data

datetime

Package for easy *date* data manipulation



jupyter applayout_example Last Checkpoint: 20 hours ago (autosaved)  Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

```
11         icon='backward',
12         layout=Layout(width='80%',
13                        height='30%'))
14 next_button = Button(description="Next",
15                       icon='forward',
16                       layout=Layout(width='80%',
17                                    height='30%'))
18 footer = HTML("Filename: {}".format(image_file))
19
20 AppLayout(header=header,
21           left_sidebar=prev_button,
22           center=image,
23           right_sidebar=next_button,
24           footer=footer,
25           grid_gap='20px',
26           justify_items='center',
27           align_items='center')
```

Simple Image Viewer



Filename: images/cat.jpg



!! Requires a gmail account to edit !!



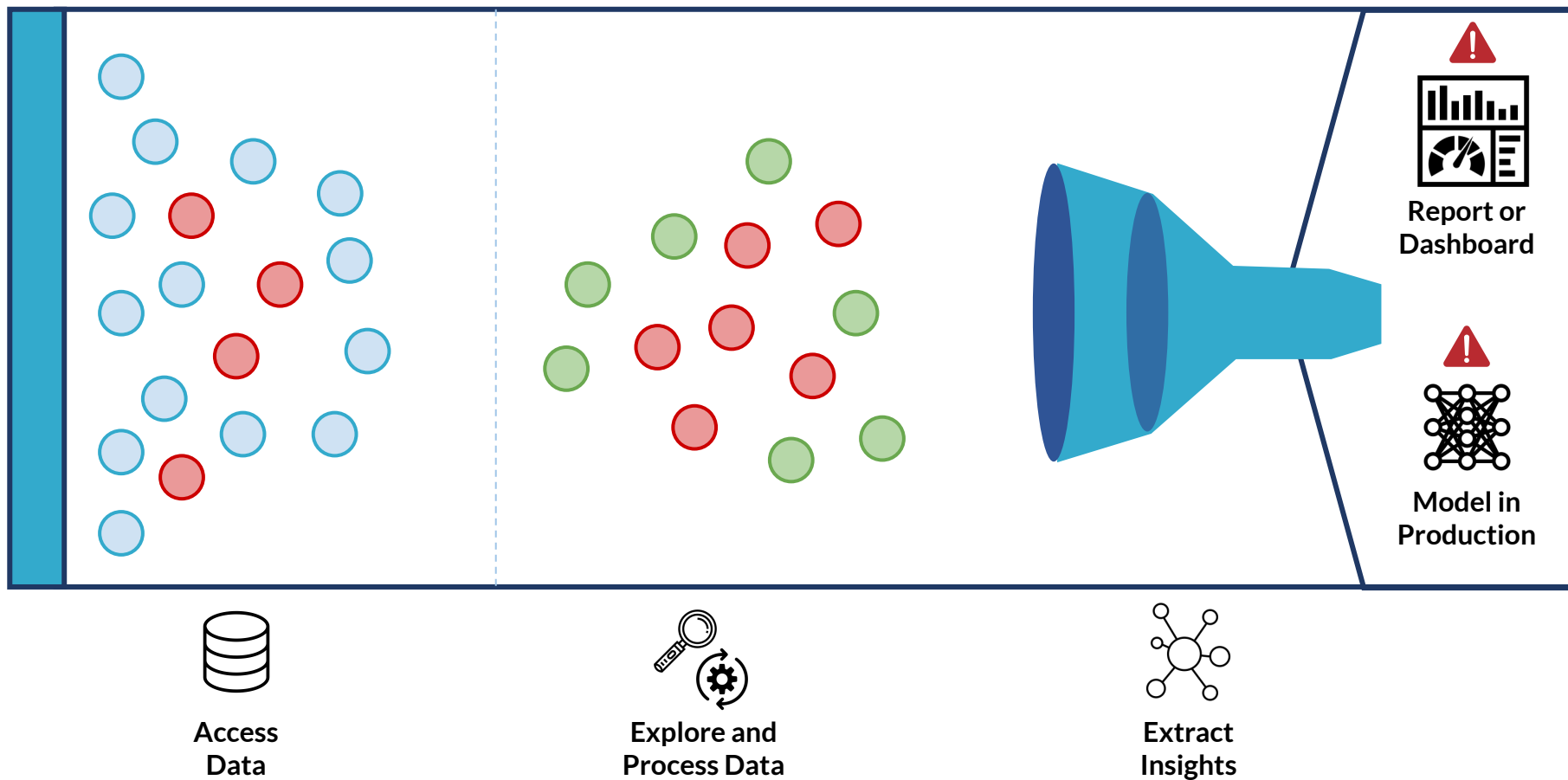
- 1** *Introduction*
- 2 *Importing our dataset*
- 3 *Diagnosing our data problems*
- 4 *Q&A*
- 5 *Our to do list*
- 6 *Data cleaning*
- 7 *Q&A*
- 8 *Recap & closing notes*
- 9 *Take home question*



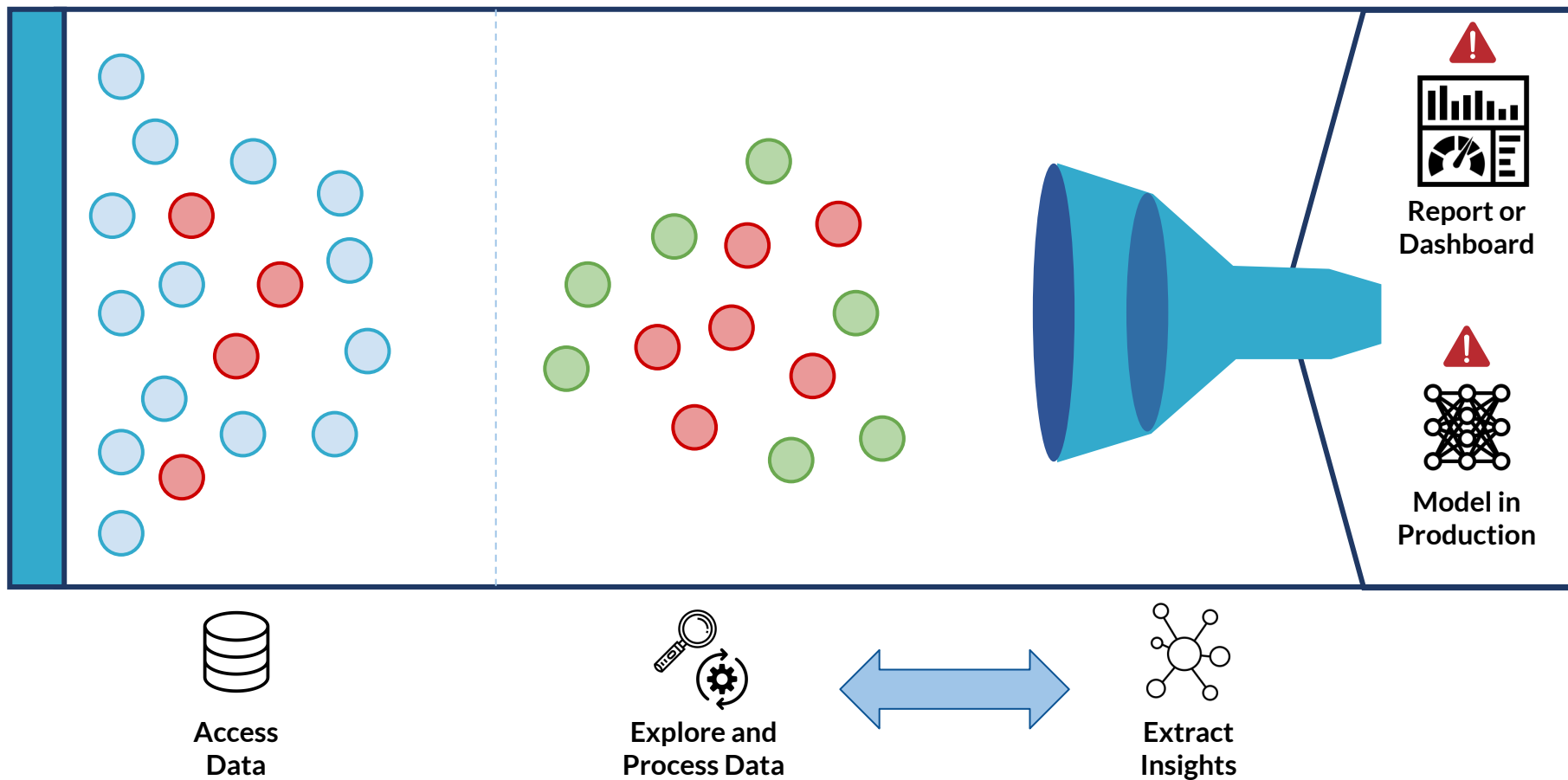
Notebook



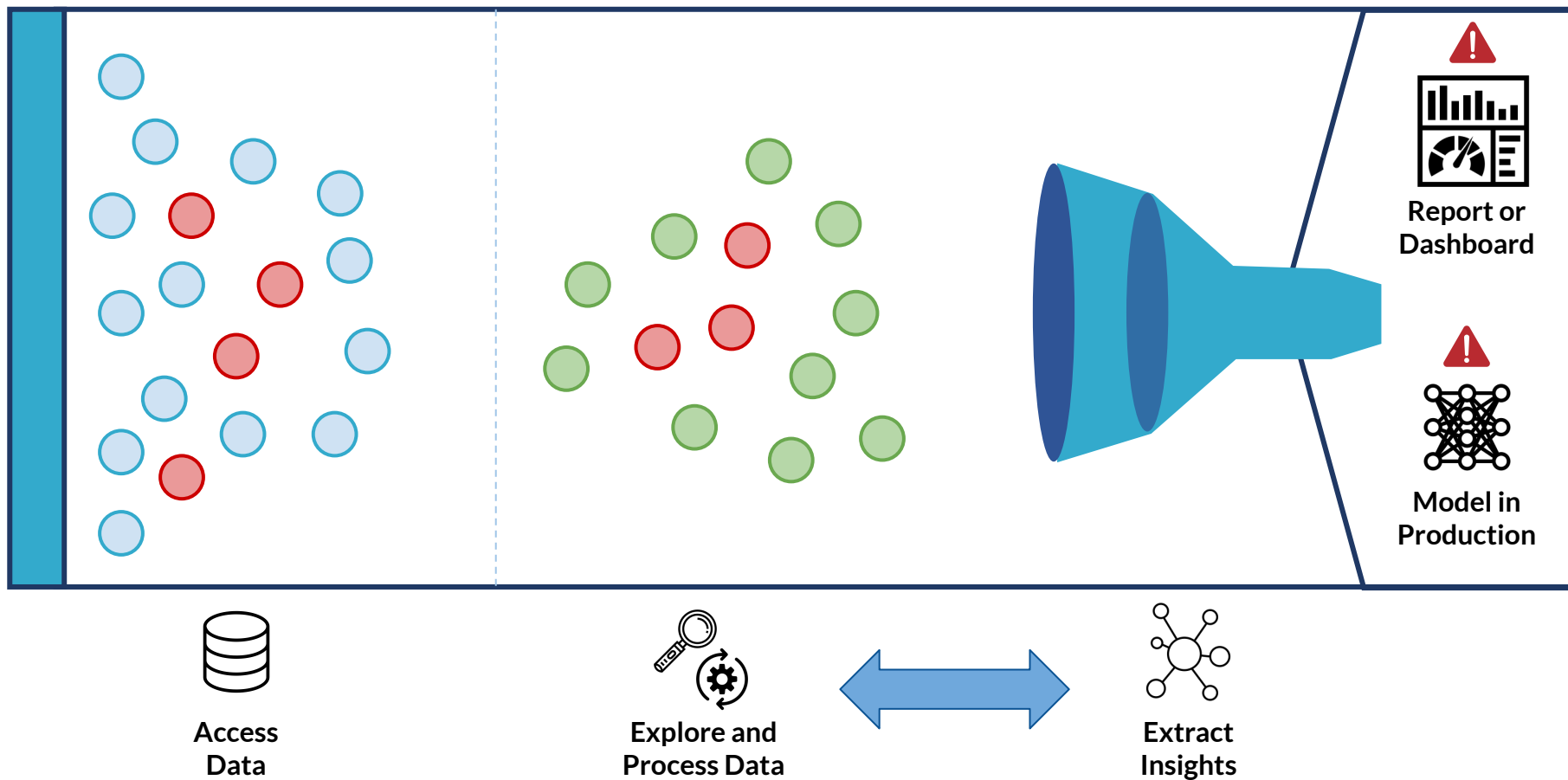
- 1 *Introduction*
- 2 *Importing our dataset*
- 3 *Diagnosing our data problems*
- 4 *Q&A*
- 5 *Our to do list*
- 6 *Data cleaning*
- 7 *Q&A*
- 8 *Recap & closing notes***
- 9 *Take home question*

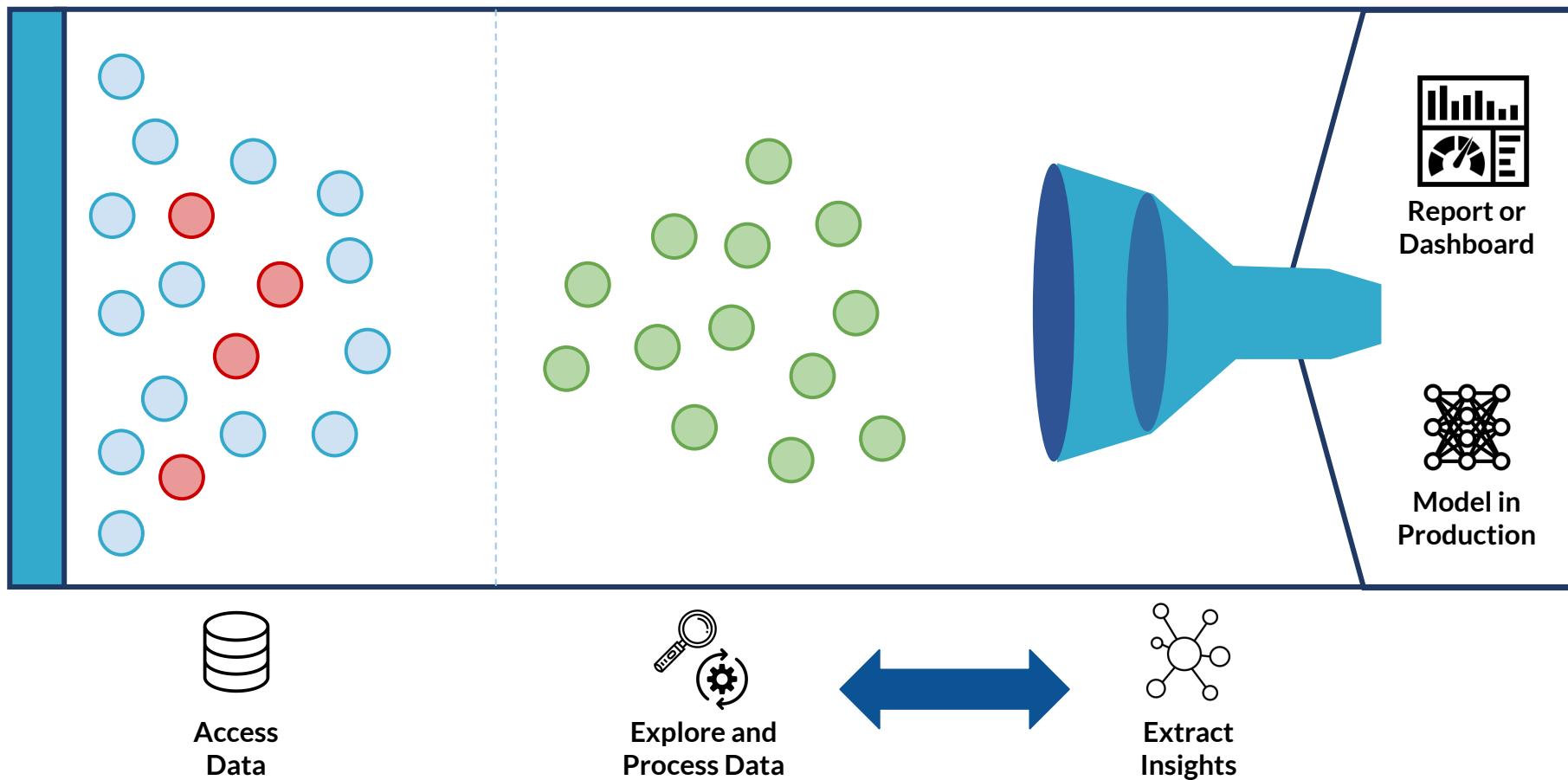


The data science workflow - revisited



The data science workflow - revisited







Check out our upcoming [webinars!](#)



Live training: Data Viz with ggplot2

Bring your data to life with the power of data viz with ggplot2, and answer real-world data questions. This session will run for three hours, allowing you time to really immerse yourself in the subject, and includes short breaks and opportunities to ask the expert questions throughout the training.

Wednesday, April 15, 2020, 11 AM EST, 4 PM BST



Register [here](#)



DCVirtual: Webinar week

Join DataCamp and industry thought leaders for a week-long virtual conference on all things data science. Learn how to roll out a data strategy that includes online training, how to measure the success of your data science initiatives, and hear from experts about how AI and machine learning are impacting industries like finance and healthcare.

Monday, April 20 to Friday, April 24, 2020, 11 AM EST

Registration Link Coming Soon





DataCamp for Enterprise: What's New in Q2 2020

Discover what's new in Q2 2020 for DataCamp Enterprise plans.

Wednesday, April 29, 2020, 11 AM EST, 4 PM BST



Register [here](#)



Pick one of the following:

- 1) *What is the average price of listings by borough? Visualize your results with a bar plot!*
- 2) *What is the average availability in days of listings by borough? Visualize your results with a bar plot!*
- 3) *What is the median price per room type in each borough? Visualize your results with a bar plot!*
- 4) *Visualize the number of listings over time.*

Functions that should/could be used:

- `.groupby()` and `.agg()`
- `sns.barplot(x = , y = , hue = , data =)`
- `sns.lineplot(x = , y = , data =)`
- `.dt.strftime()` for extracting specific dates from a `datetime` column

Bonus points if you finish more than one question

Submission details:

- Share with us a code snippet with your output on LinkedIn, Twitter or Facebook
- Tag us on ``@DataCamp`` with the hashtag ``#datacamplive``



Recap of the functions used

Diagnosis functions	Description
<code>import pandas as pd</code>	Imports the <code>pandas</code> package with the alias <code>pd</code>
<code>.head()</code>	Prints the header of a DataFrame
<code>.dtypes</code>	Gets the data types of each column in a DataFrame
<code>.info()</code>	Returns a # observations, data types and missing values per column
<code>.describe()</code>	Returns statistical distribution of numeric value in a DataFrame
<code>.isna().sum()</code>	Returns # of missing values per column
<code>sns.distplot()</code>	Plots distribution of one variable
<code>msno.matrix()</code>	Visualizes missingness matrix
<code>msno.barplot()</code>	Visualizes missingness barplot
<code>.duplicated(subset = , keep =)</code>	Lets you find duplicates in a DataFrame based on all or subset of columns

Treatment functions	Description
<code>.str.replace("", "")</code>	Replaces one string with another for each row of a <code>str</code> column
<code>.str.split("", expand = True)</code>	Splits a string column into two based on input
<code>.astype()</code>	Converts a column to a datatype of choice
<code>pd.to_datetime()</code>	Converts a date column to datetime
<code>.str.lower()</code>	Lowercases each row in a <code>str</code> column
<code>.str.strip("")</code>	Removes a pattern from each row of an <code>str</code> column
<code>.replace()</code>	Replace values for others in a column
<code>.fillna()</code>	Fills missing values of a column with a value of your choice
<code>.drop_duplicates()</code>	Drops duplicates