

# DATA 301

A comprehensive cheat sheet covering essential data analysis tools and techniques, including data representation, Excel, SQL, Python, R, visualization, open data, APIs, and geospatial concepts.



# **Data Representation & Excel Basics**

## Data Representation

| Binary:               | Base-2 numeral system (0, 1).<br>Fundamental for digital data.   |  |
|-----------------------|--|--|
| Sizes:                | MB (Megabyte), GB (Gigabyte), TB<br>(Terabyte). 1 TB = 1024 GB, 1 GB =<br>1024 MB.   |  |
| Data Types:           | Strings (text), Dates (YYYY-MM-<br>DD), Floats (decimal numbers).  |  |
| File<br>Encodings:    | UTF-8 (Unicode Transformation<br>Format - 8-bit) is the most<br>common. ASCII, UTF-16 are<br>others.   |  |
| Character<br>Encoding | A character encoding tells the<br>computer how to interpret raw<br>zeroes and ones into actual<br>characters. It determines what<br>code belongs to which character.<br>Common encodings: ASCII, UTF-<br>8, UTF-16 |  |
| Endianness            | Describes the order in which bytes<br>of a multi-byte data type (like<br>integers) are stored in computer<br>memory.   |  |
|                       | Big-Endian: Most significant byte<br>first. Little-Endian: Least<br>significant byte first.  |  |

## Essential Excel Formulas

| IF(condition,<br>value_if_true,<br>value_if_false<br>) | Returns one value if a<br>condition is TRUE and<br>another value if it's FALSE.<br><b>Example:</b> =IF(A1>10,<br>"High", "Low")                  |
|--|--|
| COUNTIF(range,<br>criteria)                            | Counts the number of cells<br>within a range that meet the<br>given criteria.<br>Example: =COUNTIF(B1:B10,<br>">50")                             |
| <pre>SUM(number1,<br/>[number2],<br/>)</pre>           | Adds all the numbers in a range of cells.<br>Example: =SUM(C1:C10)   |
| <pre>MAX(number1, [number2],)</pre>                    | Returns the largest value in a<br>set of numbers.<br>Example: =MAX(D1:D10)   |
| <pre>CONCATENATE(te xt1, [text2],)</pre>               | Joins several text strings into<br>one text string. Use &<br>operator as shortcut.<br>Example: =CONCATENATE(A1,<br>" ", B1) or =A1 & " " &<br>B1 |
| AVERAGE(number<br>1, [number2],<br>)                   | Calculates the average<br>(arithmetic mean) of the<br>numbers in a range.<br><b>Example:</b><br>=AVERAGE(E1:E10)                                 |

## Excel - Pivot Tables

Pivot tables are used to summarize and analyze data. They allow you to rearrange and group data in different ways to see patterns and trends.

#### Key Components:

- Rows: Categorical fields displayed as rows.
- **Columns:** Categorical fields displayed as columns.
- Values: Numerical fields that are aggregated (summed, averaged, etc.).
- Filters: Used to narrow down the data being displayed.

#### Creating a Pivot Table:

- 1. Select your data range.
- 2. Go to Insert > PivotTable .
- 3. Choose where to place the pivot table (new worksheet or existing location).
- 4. Drag and drop fields into the Rows, Columns, Values, and Filters areas.

#### Example Scenario:

Imagine you have sales data with columns: Date, Region, Product, Sales. You can create a pivot table to:

- Show total sales by region.
- Show average sales by product over time.
- Filter the data to show sales for a specific month.

# **SQL Fundamentals**

# Basic SQL Commands

| CREATE TABLE<br>table_name<br>(column1<br>datatype, column2<br>datatype,);                | Creates a new table in the<br>database. Define column<br>names and their data<br>types.<br>Example: CREATE TABLE<br>Employees (ID INT, Name<br>VARCHAR(255), Salary  | WHERE<br>condition<br>GROUP BY<br>column1. |
|---|--|--|
| INSERT INTO<br>table_name<br>(column1,<br>column2,)<br>VALUES (value1,<br>value2,);       | DECIMAL(10, 2));<br>Inserts a new row into a<br>table.<br>Example: INSERT INTO<br>Employees (ID, Name,<br>Salary) VALUES (1,<br>'John Doe', 60000.00);               | COLUMN2,<br><br>ORDER BY<br>COLUMN1        |
| UPDATE<br>table_name SET<br>column1 = value1,<br>column2 = value2,<br>WHERE<br>condition; | Modifies existing data in a<br>table. Use WHERE clause<br>to specify which rows to<br>update.<br>Example: UPDATE<br>Employees SET Salary =<br>65000.00 WHERE ID = 1; | [ASC DESC]<br>column2<br>[ASC DESC]<br>    |
| DELETE FROM<br>table_name WHERE<br>condition;   | Deletes rows from a table.<br>Always use a WHERE<br>clause to avoid deleting all<br>rows.<br>Example: DELETE FROM<br>Employees WHERE ID =<br>1;                      | HAVING                                     |
| SELECT column1,<br>column2, FROM<br>table_name WHERE<br>condition;                        | Retrieves data from one or more tables.  |  |

## Filtering and Aggregations

| WHERE  | Filters rows based on a specified<br>condition.<br>Example: SELECT * FROM<br>Employees WHERE Salary ><br>50000;  |
|--|--|
| GROUP BY<br>column1,<br>column2,<br>                             | Groups rows that have the same<br>values in specified columns.<br>Often used with aggregate<br>functions.<br>Example: SELECT Department,<br>AVG(Salary) FROM Employees<br>GROUP BY Department; |
| ORDER BY<br>column1<br>[ASC DESC],<br>column2<br>[ASC DESC],<br> | Sorts the result-set based on<br>one or more columns. ASC for<br>ascending, DESC for<br>descending.<br>Example: SELECT * FROM<br>Employees ORDER BY Salary<br>DESC;                            |
| LIMIT  | Limits the number of rows<br>returned.<br>Example: SELECT * FROM<br>Employees ORDER BY Salary<br>DESC LIMIT 10;  |
| HAVING   | Filters results after grouping,<br>used with GROUP BY.<br>Example: SELECT Department,<br>AVG(Salary) FROM Employees<br>GROUP BY Department HAVING<br>AVG(Salary) > 55000;                      |

#### Joins

**JOIN:** Returns rows when there is a match in both tables based on the join condition.

SELECT \* FROM TableA JOIN TableB ON TableA.Column = TableB.Column;

**LEFT JOIN (or LEFT OUTER JOIN):** Returns all rows from the left table (TableA), and the matched rows from the right table (TableB). If there is no match in TableB, it returns NULL values for columns from TableB.

SELECT \* FROM TableA LEFT JOIN TableB ON TableA.Column = TableB.Column;

**RIGHT JOIN (or RIGHT OUTER JOIN):** Returns all rows from the right table (TableB), and the matched rows from the left table (TableA). If there is no match in TableA, it returns NULL values for columns from TableA.

SELECT \* FROM TableA RIGHT JOIN TableB ON TableA.Column = TableB.Column;

FULL OUTER JOIN: Returns all rows when there is a match in one of the tables. If there are rows in TableA that do not match TableB, or rows in TableB that do not match TableA, the FULL OUTER JOIN will include these rows in the result set. Columns that do not have matching values will contain NULL.

SELECT \* FROM TableA FULL OUTER JOIN TableB ON TableA.Column = TableB.Column;

**INNER JOIN:** Same as JOIN. Returns rows only when there's a match in both tables.

SELECT \* FROM TableA INNER JOIN TableB ON TableA.Column = TableB.Column;

# Python and R Programming

| Python Fundamentals               |   | R Programming Fundamentals        |  |
|-----------------------------------|---|-----------------------------------|--|
| Input Handling:                   | <pre>Using try/except blocks<br/>to handle potential errors<br/>when taking user input.<br/>try:<br/>age =<br/>int(input("Enter your<br/>age: "))<br/>except ValueError:<br/>print("Invalid<br/>input. Please enter a</pre> | Vectors:<br>Subsetting:           | Creating and manipulating<br>vectors.<br>my_vector <- c(1, 2, 3,<br>4, 5)<br>print(my_vector[1]) #<br>Accessing vector element<br>Using subset() and indexing<br>to filter data.<br>subset_data <- |
| Loops:                            | number.")<br>for and while loops for<br>iterating over data   |                                   | <pre>subset(data, column &gt; 10) print(data[data\$column &gt; 10, ]) # Indexing</pre>   |
|                                   | <pre>structures. for i in range(5):     print(i) while True:     response =     input("Type 'exit' to</pre>   | Data Frames:                      | <pre>Working with data frames. df &lt;- data.frame(Name =   c("John", "Jane"), Age =   c(30, 25)) print(df\$Name) # Accessing column</pre>   |
|                                   | <pre>quit: ")     if response ==     'exit':         break</pre>  | Statistical<br>Tests:             | Performing t-tests and other<br>statistical analyses.<br>t.test(data\$column1,<br>data\$column2)   |
| List/Dictionary<br>Logic:         | Creating, accessing, and<br>manipulating lists and<br>dictionaries.<br>my_list = [1, 2, 3]  | Summary<br>Statistics:            | Calculating descriptive<br>statistics.<br>summary(data\$column)  |
|                                   | <pre>my_dict = {'a': 1,     'b': 2} print(my_list[0]) #</pre>   | Data Processing                   | g - Python & R   |
|                                   | Accessing list element<br>print(my_dict['a']) #<br>Accessing dictionary<br>value  | Max, Filter:                      | operations.<br>data = [1, 2, 3, 4, 5]<br>sum_data = sum(data)<br>max_data = max(data)  |
| File Reading &<br>String Parsing: | Reading data from files and manipulating strings.   |                                   | <pre>filtered_data = list(filter(lambda x: x &gt; 2, data))</pre>  |
|                                   | 'r') as f:<br>content = f.read()<br>words =<br>content.split()  | R - Filtering &<br>Summary Stats: | Data manipulation and<br>summary statistics.<br>data <- c(1, 2, 3, 4,<br>5)  |
| Function<br>Definitions:          | <pre>Defining reusable blocks of<br/>code.<br/>def greet(name):<br/>print(f"Hello,<br/>{name}!")</pre>  |                                   | filtered_data <-<br>data[data > 2]<br>summary_data <-<br>summary(data)   |
|                                   | areet("World")  |                                   |  |

# Visualization, Open Data, and Geospatial

Data Visualization Tools

| Python:  | Using (matplotlib) and (seaborn) for creating visualizations.   |
|----------|---|
|          | <pre>import matplotlib.pyplot as plt<br/>import seaborn as sns<br/>sns.histplot(data['column'])</pre>   |
|          | <pre>plt.show()</pre>   |
| R:       | Using ggplot2 for creating visualizations.  |
|          | library(ggplot2)<br>ggplot(data, aes(x = column)) +<br>geom_histogram()   |
| Tableau: | A data visualization tool for creating<br>interactive dashboards and reports.<br>Dimensions (blue pills) are categorical,<br>and Measures (green pills) are<br>numerical. |
|          | Drag dimensions and measures to<br>rows, columns, and marks cards to<br>create visualizations.  |
|          |   |

| CSV/JSON Data<br>Loading (Python): | Loading data from CSV and<br>JSON files.<br>import pandas as pd<br>data_csv =<br>pd.read_csv('data.csv'<br>)<br>data_json =<br>pd.read_json('data.jso<br>n')                          |
|------------------------------------|---|
| Accessing APIs<br>(Python):        | <pre>Accessing data from APIs,<br/>e.g., Google Maps.<br/>import requests<br/>response =<br/>requests.get('https://<br/>maps.googleapis.com/<br/>.')<br/>data = response.json()</pre> |
| CSV/JSON Data<br>Loading (R):      | Loading data from CSV and<br>JSON files.<br>data_csv <-<br>read.csv('data.csv')<br>library(jsonlite)<br>data_json <-<br>fromJSON('data.json')   |
| Accessing APIs<br>(R):             | Accessing data from APIs<br>library(httr)<br>response <-<br>GET('https://api.examp<br>le.com/data')<br>data <-<br>content(response,<br>"parsed")                                      |

Open Data & APIs

# Geospatial (GIS) Fundamentals

| <b>Coordinate Data:</b><br>Latitude and Longitude are used to specify<br>locations on Earth. Latitude ranges from -90 to<br>+90 (degrees North/South), and Longitude<br>ranges from -180 to +180 (degrees East/West). |
|---|
| <ul> <li>Geospatial File Types:</li> <li>CSV with lat/lon: Simple text file where columns represent latitude and longitude coordinates.</li> </ul>  |
| • <b>KML (Keyhole Markup Language):</b> XML-<br>based file format for representing<br>geographic data in Google Earth, Google<br>Maps, and other GIS software.  |
| Simple Mapping (Python):<br>Using libraries like folium to create interactive<br>maps.  |
| <pre>import folium m = folium.Map(location=[40.7128, -74.0060], zoom_start=10) m.save('map.html')</pre>   |
| Simple Mapping (R):<br>Using libraries like leaflet to create interactive<br>maps.  |
| library(leaflet)  |
| addTiles() %>%  |

addTiles() %>% setView(lng = -74.0060, lat = 40.7128, zoom = 10)