# To Get Started

- Paper sheet and online at:
  http://www.eecs.qmul.ac.uk/~william/CAS-London-2020.html

- Download sample notebooks and data
  - Create directory (N:\session3\**fullname**)and unzip notebooks and data

- Login to Google Colab: https://colab.research.google.com/
  - Create a new notebook
  - Use the file 'upload' menu to upload the 'example' notebook

# Introduction to Data Analysis

William Marsh

Queen Mary University of London

institute of CODING

COMPUTING AT SCHOOL
EDUCATE · ENGAGE · ENCOURAGE
Part of BCS –The Chartered Institute for IT

# How This Session Works

- Introduce concepts
- Practical work
  - Collaboration: *practice teaching!*
- Repeat

- Conclude

**Probably Not Enough Time**

# Outline

- Aims
- Introducing the Python notebook using Google Colab
- Part 1: the dataframe
- Part 2: transforming data and the pivot table
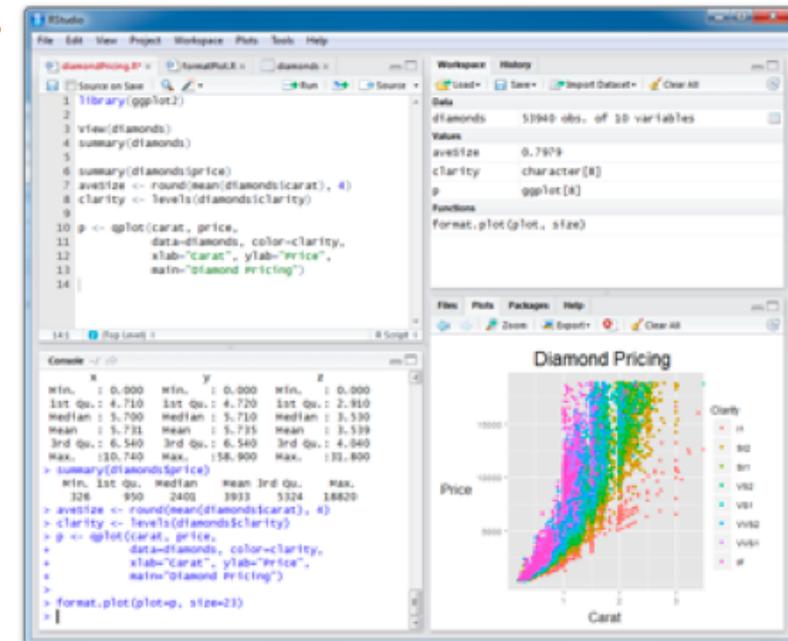- Part 3: adding columns
- Conclusion and discussion

# Session Aims

# Aims

- Introduce Python Pandas
  - Popular library for data
  - Interactive notebook – Google's Colab
- Introduce key concepts
  - Dataframe
  - Filter and select
  - Pivot table
  - Visualisation

- Pandas is very complex
- Concepts common in other environments
  - Excel
  - RStudio

- Help develop pedagogy

# The Data Life Cycle

- Using data to answer a question
  - What is the problem?
  - .. do we have the data?

- Data analysis not just technical

- We focus on the technical

# Introducing the Interactive Notebook

# IPython Notebook

- IPython – interactive Python with graphics (2001)
- Jupyter – web based interface for IPython (2014)
  - Also supports other languages
- Google Colab
  - Hosted support for Jupyter on Google Drive
  - Better interface

- Concept: program as an executable document

New, upload, save

Text: written in 'markdown'

TOC: from 'sections'

Run button

Code cell and output

Example.ipynb

File   Edit   View   Insert   Runtime   Tools   Help   All changes saved

Comment     Share

+ Code   + Text

RAM
Disk

Editing

Table of contents

Introducing the Notebook

Cells in the Notebook

Write a simple program in the notebook

Section

# Introducing the Notebook

Google colab is a version of the 'iPython' notebook. This is an example note book. Edit this note book to learn how the note book system works.

## Cells in the Notebook

The notebook is a sequence if cells.

- A cell (such as this one) can contain text. This makes a data analysis a doocument.
- Other cells in the notebook contain code. The analysis document becomes executable.

The next cell contaions code. Run it!

```
[ ]   x = [x for x in range(1,10)]
      x
```

    [1, 2, 3, 4, 5, 6, 7, 8, 9]

# Practical Break

Please try the 'example' notebook

# The Data Frame

- Header row
  - Shows the columns

- Rows
  - Shows individuals

- Tidy data
  - All columns have headings
  - All columns same 'type' (e.g. numbers)
  - No blanks

|   | Name | Age | Team |
|---|------|-----|------|
| 0 | John | 24 | Arsenal |
| 1 | Mary | 27 | Spurs |
| 2 | Peter | 31 | Chelsea |

```
,Name,Age,Team
0,John,24,Arsenal
1,Mary,27,Spurs
2,Peter,31,Chelsea
```

Loaded from CSV

# The Data: Country of Birth

- Taken from 2011 census
  - 67,252 row
  - Example of 'narrow' or 'tall' data

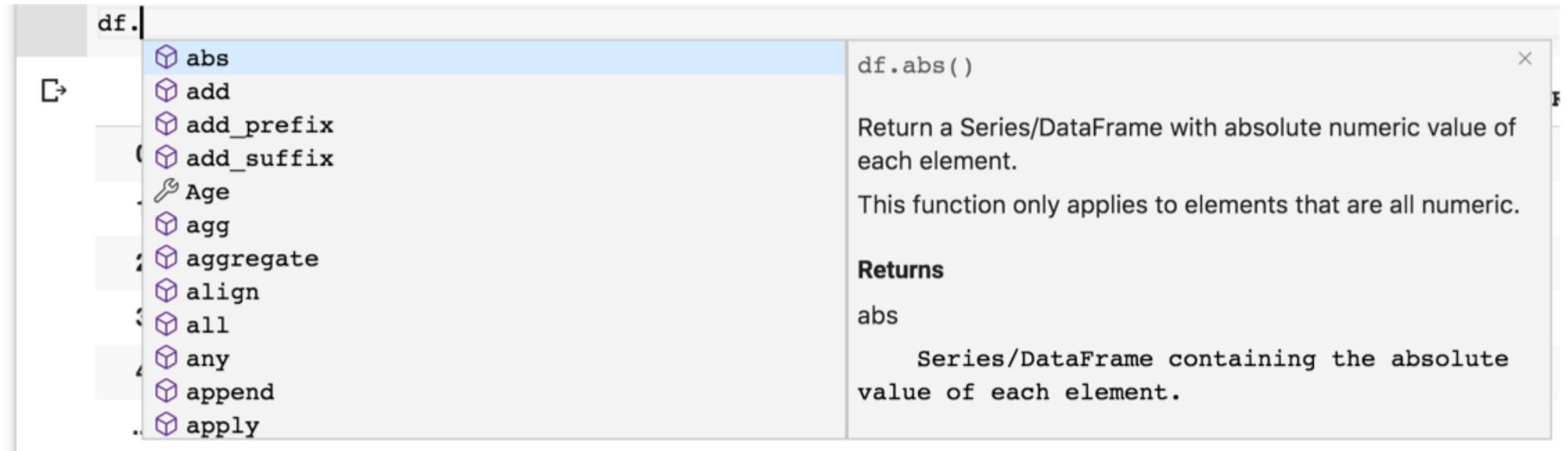| Area | Age | Sex | Usual Residents | Birth Country | Birth Region |
|---|---|---|---|---|---|
| Tower Hamlets | Age 0 to 4 | Females | 3 | Ghana | Africa |
| Tower Hamlets | Age 5 to 9 | Females | 2 | Ghana | Africa |
| Tower Hamlets | Age 10 to 15 | Females | 4 | Ghana | Africa |

# Investigating the Data

- Task 1.1: Load the data to a dataframe
- Task 1.2: Look at values the unique values in each column
  - Answer some questions about the data

| Column | Description |
|---|---|
| Area | Includes London Boroughs |
| Age | The ages in a number of bands |
| Sex | Males and Females |
| Usual Residents | An integer |
| BirthCountry | A country |
| BirthRegion | A region e.g. Africa or Europe |

# Getting Help

- Pop up help



- Pandas documentation - https://pandas.pydata.org/docs/
  - API – many optional arguments
  - User guides

# Practical Break

Part 1: 'Data analysis' notebook

# Part 2: Selecting, Transforming and visualising Data

- Data transformation: Tall to Wide

- The Pivot Table

# What Questions Can this Dataset Answer?

- How many people from the Americas live in Sutton?

- Which Borough has the 'most' young people?

- What age are people born outside the UK?

- ...

- *What additional information is missing?*

# Selecting and Transforming Data

- Selecting some data
  - Data for one borough
  - … or one age group
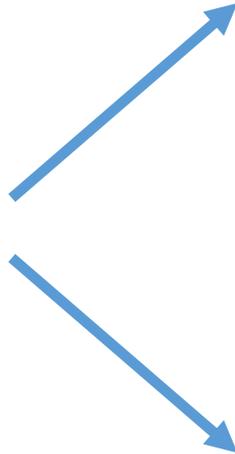  - … or one country or region

- Use conditions

- *No loops*

```
df[(df['Area']=='Tower Hamlets')]

df[(df['Area']=='Tower Hamlets') & (df['Age'] =='Age 0 to 4')]
```

# The Pivot Table

- Transform data
- Origin in spreadsheets

| Person | Genre | Rating |
|--------|--------|--------|
| Andy | Classic | Like |
| Andy | Jazz | Hate |
| Andy | Folk | Hate |
| Bill | Classic | Hate |
| Bill | Jazz | Like |
| Bill | Folk | Like |
| Charlie | Classic | Like |
| Charlie | Jazz | Like |
| Charlie | Folk | Hate |

| Person | Genre | | |
|--------|---------|------|------|
| | Classic | Jazz | Folk |
| Andy | Like | Hate | Hate |
| Bill | Hate | Like | Like |
| Charlie | Like | Like | Hate |

| Genre | Person | | |
|-------|--------|------|---------|
| | Andy | Bill | Charlie |
| Classic | Like | Hate | Like |
| Jazz | Hate | Like | Like |
| Folk | Hate | Like | Hate |

# The Pivot Table

- Transform data

| Person | Place | Purpose | Visits |
|--------|-------|---------|--------|
| Andy | Berlin | Hols | 1 |
| Andy | Berlin | Work | 2 |
| Andy | Paris | Hols | 2 |
| Andy | NY | Work | 3 |
| Andy | Madrid | Hols | 1 |
| Bill | Berlin | Work | 4 |
| Bill | Paris | Work | 3 |
| Charlie | Paris | Hols | 1 |
| Charlie | Rome | Hols | 1 |
| Charlie | Zurich | Hols | 1 |

| Purpose | Visits | | |
|---------|------|------|---------|
| | Andy | Bill | Charlie |
| Hols | 4 | 0 | 3 |
| Work | 5 | 7 | 0 |

| Purpose | Visits | | | | | |
|---------|--------|--------|----|-------|------|--------|
| | Berlin | Madrid | NY | Paris | Rome | Zurich |
| Hols | 1 | 0 | 3 | 3 | 1 | 1 |
| Work | 6 | 1 | 0 | 3 | 0 | 0 |

Aggregate over the person

# Practical Break

Part 2: 'Data analysis' notebook
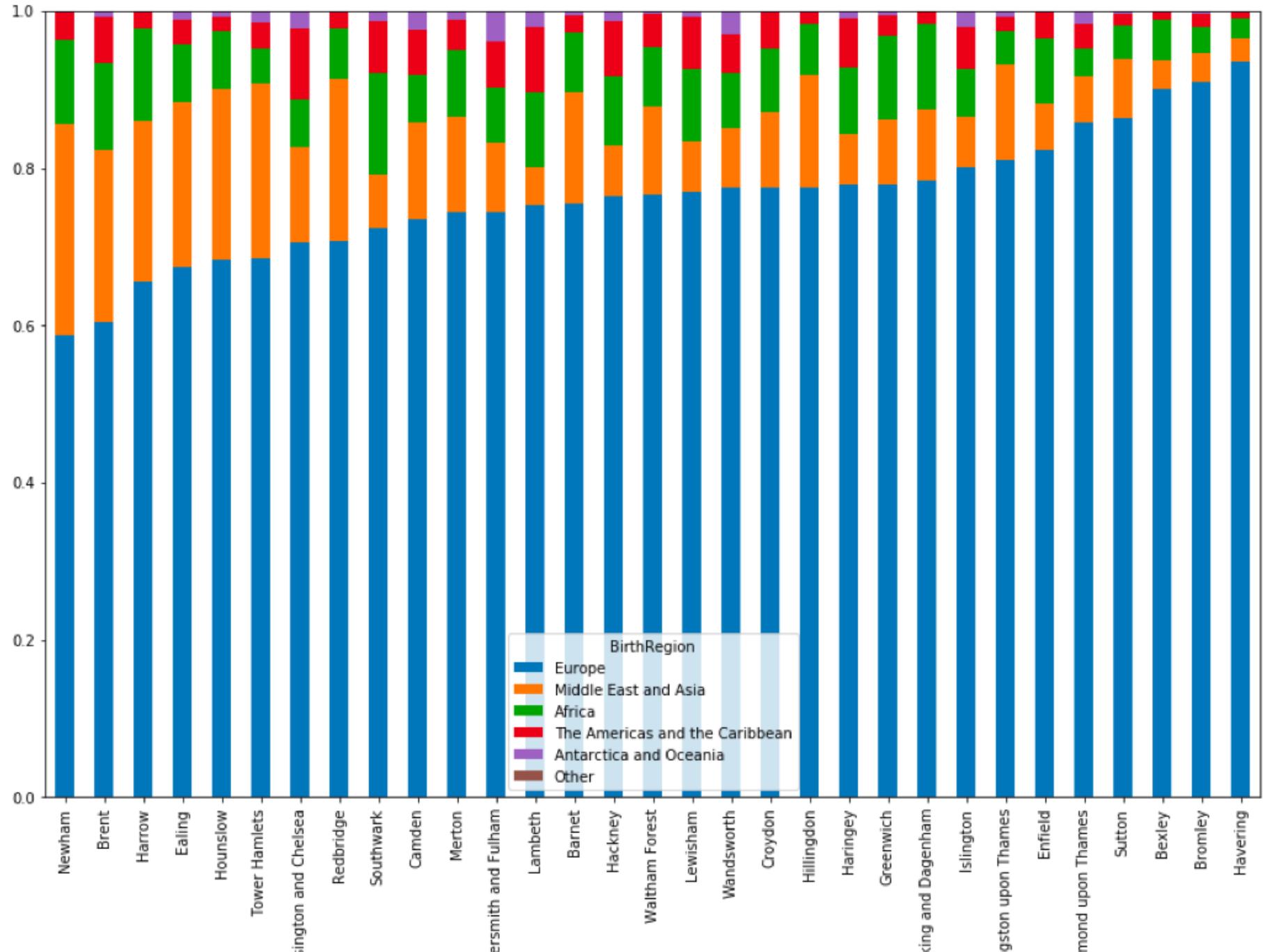
# Part 3: Adding Columns

# The Problem of Comparing Boroughs

- Boroughs vary in population

- What does this question mean?

*Where in London do most people born in Spain live?*

- Need to transform the data into 'proportion'

# Proportion of Borough Born in Each Region

# Practical Break

Part 3: 'Data analysis' notebook

# Summary and Discussion

# Summary

- The environment
  - Web-based Interactive Python
  - Does not have to be hosted

- The Pandas library
  - Comprehensive but complex

- Key concepts
  - Dataframe and selecting data
  - Transforming: pivot table
  - Visualizing: plot

**Discussion Questions**
- What is emphasis of curriculum?
- Should we use a large data file?
- Is the complexity of Pandas manageable?
- Balance of technical versus 'interpretation' of data