

Part X

Data and Empirics

PANDAS FOR PANEL DATA

Contents

- *Pandas for Panel Data*
 - *Overview*
 - *Slicing and Reshaping Data*
 - *Merging Dataframes and Filling NaNs*
 - *Grouping and Summarizing Data*
 - *Final Remarks*
 - *Exercises*
 - *Solutions*

62.1 Overview

In an [earlier lecture on pandas](#), we looked at working with simple data sets.

Econometricians often need to work with more complex data sets, such as panels.

Common tasks include

- Importing data, cleaning it and reshaping it across several axes.
- Selecting a time series or cross-section from a panel.
- Grouping and summarizing data.

`pandas` (derived from ‘panel’ and ‘data’) contains powerful and easy-to-use tools for solving exactly these kinds of problems.

In what follows, we will use a panel data set of real minimum wages from the OECD to create:

- summary statistics over multiple dimensions of our data
- a time series of the average minimum wage of countries in the dataset
- kernel density estimates of wages by continent

We will begin by reading in our long format panel data from a CSV file and reshaping the resulting `DataFrame` with `pivot_table` to build a `MultiIndex`.

Additional detail will be added to our DataFrame using pandas' merge function, and data will be summarized with the groupby function.

62.2 Slicing and Reshaping Data

We will read in a dataset from the OECD of real minimum wages in 32 countries and assign it to `realwage`.

The dataset can be accessed with the following link:

```
url1 = 'https://raw.githubusercontent.com/QuantEcon/lecture-python/master/source/_
static/lecture_specific/pandas_panel/realwage.csv'
```

```
import pandas as pd

# Display 6 columns for viewing purposes
pd.set_option('display.max_columns', 6)

# Reduce decimal points to 2
pd.options.display.float_format = '{:,.2f}'.format

realwage = pd.read_csv(url1)
```

Let's have a look at what we've got to work with

```
realwage.head() # Show first 5 rows
```

Unnamed: 0	Time	Country	Series	\
0	2006-01-01	Ireland	In 2015 constant prices at 2015 USD PPPs	
1	2007-01-01	Ireland	In 2015 constant prices at 2015 USD PPPs	
2	2008-01-01	Ireland	In 2015 constant prices at 2015 USD PPPs	
3	2009-01-01	Ireland	In 2015 constant prices at 2015 USD PPPs	
4	2010-01-01	Ireland	In 2015 constant prices at 2015 USD PPPs	

Pay period	value
0 Annual	17,132.44
1 Annual	18,100.92
2 Annual	17,747.41
3 Annual	18,580.14
4 Annual	18,755.83

The data is currently in long format, which is difficult to analyze when there are several dimensions to the data.

We will use `pivot_table` to create a wide format panel, with a `MultiIndex` to handle higher dimensional data.

`pivot_table` arguments should specify the data (values), the index, and the columns we want in our resulting dataframe.

By passing a list in columns, we can create a `MultiIndex` in our column axis

```
realwage = realwage.pivot_table(values='value',
                                index='Time',
                                columns=['Country', 'Series', 'Pay period'])
realwage.head()
```

```

Country          Australia \
Series      In 2015 constant prices at 2015 USD PPPs
Pay period          Annual Hourly
Time
2006-01-01          20,410.65  10.33
2007-01-01          21,087.57  10.67
2008-01-01          20,718.24  10.48
2009-01-01          20,984.77  10.62
2010-01-01          20,879.33  10.57

Country          ... \
Series      In 2015 constant prices at 2015 USD exchange rates ...
Pay period          Annual ...
Time
2006-01-01          23,826.64 ...
2007-01-01          24,616.84 ...
2008-01-01          24,185.70 ...
2009-01-01          24,496.84 ...
2010-01-01          24,373.76 ...

Country          United States \
Series      In 2015 constant prices at 2015 USD PPPs
Pay period          Hourly
Time
2006-01-01          6.05
2007-01-01          6.24
2008-01-01          6.78
2009-01-01          7.58
2010-01-01          7.88

Country          ... \
Series      In 2015 constant prices at 2015 USD exchange rates ...
Pay period          Annual Hourly
Time
2006-01-01          12,594.40  6.05
2007-01-01          12,974.40  6.24
2008-01-01          14,097.56  6.78
2009-01-01          15,756.42  7.58
2010-01-01          16,391.31  7.88

[5 rows x 128 columns]

```

To more easily filter our time series data, later on, we will convert the index into a `DatetimeIndex`

```

realwage.index = pd.to_datetime(realwage.index)
type(realwage.index)

```

```

pandas.core.indexes.datetimes.DatetimeIndex

```

The columns contain multiple levels of indexing, known as a `MultiIndex`, with levels being ordered hierarchically (`Country > Series > Pay period`).

A `MultiIndex` is the simplest and most flexible way to manage panel data in pandas

```

type(realwage.columns)

```

```
pandas.core.indexes.multi.MultiIndex
```

```
realwage.columns.names
```

```
FrozenList(['Country', 'Series', 'Pay period'])
```

Like before, we can select the country (the top level of our MultiIndex)

```
realwage['United States'].head()
```

```
Series      In 2015 constant prices at 2015 USD PPPs      \
Pay period      Annual Hourly
Time
2006-01-01      12,594.40      6.05
2007-01-01      12,974.40      6.24
2008-01-01      14,097.56      6.78
2009-01-01      15,756.42      7.58
2010-01-01      16,391.31      7.88

Series      In 2015 constant prices at 2015 USD exchange rates
Pay period      Annual Hourly
Time
2006-01-01      12,594.40      6.05
2007-01-01      12,974.40      6.24
2008-01-01      14,097.56      6.78
2009-01-01      15,756.42      7.58
2010-01-01      16,391.31      7.88
```

Stacking and unstacking levels of the MultiIndex will be used throughout this lecture to reshape our dataframe into a format we need.

.stack() rotates the lowest level of the column MultiIndex to the row index (.unstack() works in the opposite direction - try it out)

```
realwage.stack().head()
```

```
Country      Australia      \
Series      In 2015 constant prices at 2015 USD PPPs
Time      Pay period
2006-01-01 Annual      20,410.65
           Hourly      10.33
2007-01-01 Annual      21,087.57
           Hourly      10.67
2008-01-01 Annual      20,718.24

Country      \
Series      In 2015 constant prices at 2015 USD exchange rates
Time      Pay period
2006-01-01 Annual      23,826.64
           Hourly      12.06
2007-01-01 Annual      24,616.84
           Hourly      12.46
2008-01-01 Annual      24,185.70

Country      Belgium      ...      \
Series      In 2015 constant prices at 2015 USD PPPs      ...
```

(continues on next page)

(continued from previous page)

```

Time      Pay period
2006-01-01 Annual      21,042.28 ...
              Hourly      10.09 ...
2007-01-01 Annual      21,310.05 ...
              Hourly      10.22 ...
2008-01-01 Annual      21,416.96 ...

Country      United Kingdom \
Series      In 2015 constant prices at 2015 USD exchange rates
Time      Pay period
2006-01-01 Annual      20,376.32
              Hourly      9.81
2007-01-01 Annual      20,954.13
              Hourly      10.07
2008-01-01 Annual      20,902.87

Country      United States \
Series      In 2015 constant prices at 2015 USD PPPs
Time      Pay period
2006-01-01 Annual      12,594.40
              Hourly      6.05
2007-01-01 Annual      12,974.40
              Hourly      6.24
2008-01-01 Annual      14,097.56

Country      \
Series      In 2015 constant prices at 2015 USD exchange rates
Time      Pay period
2006-01-01 Annual      12,594.40
              Hourly      6.05
2007-01-01 Annual      12,974.40
              Hourly      6.24
2008-01-01 Annual      14,097.56

[5 rows x 64 columns]

```

We can also pass in an argument to select the level we would like to stack

```
realwage.stack(level='Country').head()
```

```

Series      In 2015 constant prices at 2015 USD PPPs      \
Pay period      Annual Hourly
Time      Country
2006-01-01 Australia      20,410.65  10.33
              Belgium      21,042.28  10.09
              Brazil      3,310.51  1.41
              Canada      13,649.69  6.56
              Chile      5,201.65  2.22

Series      In 2015 constant prices at 2015 USD exchange rates
Pay period      Annual Hourly
Time      Country
2006-01-01 Australia      23,826.64  12.06
              Belgium      20,228.74  9.70
              Brazil      2,032.87  0.87
              Canada      14,335.12  6.89
              Chile      3,333.76  1.42

```

Using a `DatetimeIndex` makes it easy to select a particular time period.

Selecting one year and stacking the two lower levels of the `MultiIndex` creates a cross-section of our panel data

```
realwage['2015'].stack(level=(1, 2)).transpose().head()
```

```
Time
Series      In 2015 constant prices at 2015 USD PPPs
Pay period
Country
Australia    21,715.53  10.99
Belgium      21,588.12  10.35
Brazil        4,628.63   2.00
Canada       16,536.83   7.95
Chile         6,633.56   2.80

Time
Series      In 2015 constant prices at 2015 USD exchange rates
Pay period
Country
Australia    25,349.90  12.83
Belgium      20,753.48   9.95
Brazil        2,842.28   1.21
Canada       17,367.24   8.35
Chile         4,251.49   1.81
```

For the rest of lecture, we will work with a dataframe of the hourly real minimum wages across countries and time, measured in 2015 US dollars.

To create our filtered dataframe (`realwage_f`), we can use the `xs` method to select values at lower levels in the multiindex, while keeping the higher levels (countries in this case)

```
realwage_f = realwage.xs(('Hourly', 'In 2015 constant prices at 2015 USD exchange_
↪rates'),
                        level=('Pay period', 'Series'), axis=1)
realwage_f.head()
```

```
Country      Australia  Belgium  Brazil  ...  Turkey  United Kingdom  \
Time
2006-01-01      12.06     9.70    0.87  ...    2.27             9.81
2007-01-01      12.46     9.82    0.92  ...    2.26            10.07
2008-01-01      12.24     9.87    0.96  ...    2.22            10.04
2009-01-01      12.40    10.21    1.03  ...    2.28            10.15
2010-01-01      12.34    10.05    1.08  ...    2.30             9.96

Country      United States
Time
2006-01-01         6.05
2007-01-01         6.24
2008-01-01         6.78
2009-01-01         7.58
2010-01-01         7.88

[5 rows x 32 columns]
```

62.3 Merging Dataframes and Filling NaNs

Similar to relational databases like SQL, pandas has built in methods to merge datasets together.

Using country information from [WorldData.info](https://raw.githubusercontent.com/QuantEcon/lecture-python/master/source/_static/lecture_specific/pandas_panel/countries.csv), we'll add the continent of each country to `realwage_f` with the `merge` function.

The dataset can be accessed with the following link:

```
url2 = 'https://raw.githubusercontent.com/QuantEcon/lecture-python/master/source/_static/lecture_specific/pandas_panel/countries.csv'
```

```
worlddata = pd.read_csv(url2, sep=';')
worlddata.head()
```

```

   Country (en) Country (de)      Country (local)  ... Deathrate  \
0  Afghanistan  Afghanistan  Afganistan/Afghanistan  ...      13.70
1      Egypt     Ägypten           Misr           ...       4.70
2  Åland Islands  Ålandinseln           Åland           ...       0.00
3      Albania   Albanien           Shqipëria        ...       6.70
4      Algeria   Algerien   Al-Jaza'ir/Algérie        ...       4.30

   Life expectancy  Url
0      51.30  https://www.laenderdaten.info/Asien/Afghanista...
1      72.70  https://www.laenderdaten.info/Afrika/Aegypten/...
2       0.00  https://www.laenderdaten.info/Europa/Aland/ind...
3      78.30  https://www.laenderdaten.info/Europa/Albanien/...
4      76.80  https://www.laenderdaten.info/Afrika/Algerien/...

[5 rows x 17 columns]
```

First, we'll select just the country and continent variables from `worlddata` and rename the column to 'Country'

```
worlddata = worlddata[['Country (en)', 'Continent']]
worlddata = worlddata.rename(columns={'Country (en)': 'Country'})
worlddata.head()
```

```

   Country Continent
0  Afghanistan   Asia
1      Egypt   Africa
2  Åland Islands  Europe
3      Albania  Europe
4      Algeria  Africa
```

We want to merge our new dataframe, `worlddata`, with `realwage_f`.

The pandas merge function allows dataframes to be joined together by rows.

Our dataframes will be merged using country names, requiring us to use the transpose of `realwage_f` so that rows correspond to country names in both dataframes

```
realwage_f.transpose().head()
```

```

Time      2006-01-01  2007-01-01  2008-01-01  ...  2014-01-01  2015-01-01  \
Country
Australia      12.06      12.46      12.24  ...      12.67      12.83
Belgium         9.70       9.82       9.87  ...      10.01       9.95
```

(continues on next page)

(continued from previous page)

Brazil	0.87	0.92	0.96	...	1.21	1.21
Canada	6.89	6.96	7.24	...	8.22	8.35
Chile	1.42	1.45	1.44	...	1.76	1.81
Time	2016-01-01					
Country						
Australia	12.98					
Belgium	9.76					
Brazil	1.24					
Canada	8.48					
Chile	1.91					
[5 rows x 11 columns]						

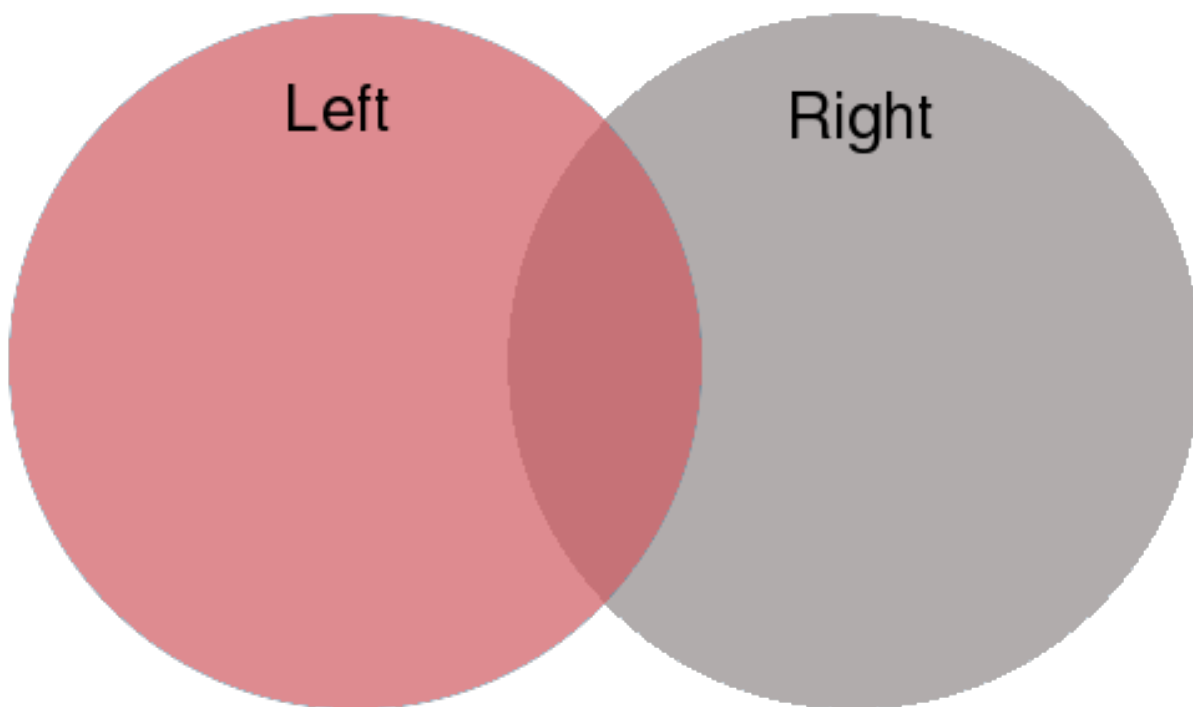
We can use either left, right, inner, or outer join to merge our datasets:

- left join includes only countries from the left dataset
- right join includes only countries from the right dataset
- outer join includes countries that are in either the left and right datasets
- inner join includes only countries common to both the left and right datasets

By default, `merge` will use an inner join.

Here we will pass `how='left'` to keep all countries in `realwage_f`, but discard countries in `worlddata` that do not have a corresponding data entry `realwage_f`.

This is illustrated by the red shading in the following diagram



We will also need to specify where the country name is located in each dataframe, which will be the `key` that is used to merge the dataframes 'on'.

Our 'left' dataframe (`realwage_f.transpose()`) contains countries in the index, so we set `left_index=True`.

Our 'right' dataframe (`worlddata`) contains countries in the 'Country' column, so we set `right_on='Country'`

```
merged = pd.merge(realwage_f.transpose(), worlddata,
                  how='left', left_index=True, right_on='Country')
merged.head()
```

```

      2006-01-01 00:00:00  2007-01-01 00:00:00  2008-01-01 00:00:00  ...  \
17.00                12.06                12.46                12.24  ...
23.00                9.70                9.82                9.87  ...
32.00                0.87                0.92                0.96  ...
100.00              6.89                6.96                7.24  ...
38.00                1.42                1.45                1.44  ...

      2016-01-01 00:00:00  Country  Continent
17.00                12.98  Australia  Australia
23.00                9.76   Belgium    Europe
32.00                1.24   Brazil  South America
100.00              8.48   Canada  North America
38.00                1.91    Chile  South America

[5 rows x 13 columns]
```

Countries that appeared in `realwage_f` but not in `worlddata` will have NaN in the Continent column.

To check whether this has occurred, we can use `.isnull()` on the continent column and filter the merged dataframe

```
merged[merged['Continent'].isnull()]
```

```

      2006-01-01 00:00:00  2007-01-01 00:00:00  2008-01-01 00:00:00  ...  \
nan                3.42                3.74                3.87  ...
nan                0.23                0.45                0.39  ...
nan                1.50                1.64                1.71  ...

      2016-01-01 00:00:00  Country  Continent
nan                5.28    Korea      NaN
nan                0.55  Russian Federation  NaN
nan                2.08  Slovak Republic    NaN

[3 rows x 13 columns]
```

We have three missing values!

One option to deal with NaN values is to create a dictionary containing these countries and their respective continents.

`.map()` will match countries in `merged['Country']` with their continent from the dictionary.

Notice how countries not in our dictionary are mapped with NaN

```
missing_continents = {'Korea': 'Asia',
                     'Russian Federation': 'Europe',
                     'Slovak Republic': 'Europe'}

merged['Country'].map(missing_continents)
```

```

17.00      NaN
23.00      NaN
32.00      NaN
100.00     NaN
38.00      NaN
108.00     NaN
41.00      NaN
225.00     NaN
53.00      NaN
58.00      NaN
45.00      NaN
68.00      NaN
233.00     NaN
86.00      NaN
88.00      NaN
91.00      NaN
nan        Asia
117.00     NaN
122.00     NaN
123.00     NaN
138.00     NaN
153.00     NaN
151.00     NaN
174.00     NaN
175.00     NaN
nan        Europe
nan        Europe
198.00     NaN
200.00     NaN
227.00     NaN
241.00     NaN
240.00     NaN
Name: Country, dtype: object

```

We don't want to overwrite the entire series with this mapping.

`.fillna()` only fills in NaN values in `merged['Continent']` with the mapping, while leaving other values in the column unchanged

```

merged['Continent'] = merged['Continent'].fillna(merged['Country'].map(missing_
    ↪continents))

# Check for whether continents were correctly mapped

merged[merged['Country'] == 'Korea']

```

```

      2006-01-01 00:00:00  2007-01-01 00:00:00  2008-01-01 00:00:00  ...  \
nan                3.42                3.74                3.87  ...

      2016-01-01 00:00:00  Country  Continent
nan                5.28    Korea    Asia

[1 rows x 13 columns]

```

We will also combine the Americas into a single continent - this will make our visualization nicer later on.

To do this, we will use `.replace()` and loop through a list of the continent values we want to replace

```
replace = ['Central America', 'North America', 'South America']

for country in replace:
    merged['Continent'].replace(to_replace=country,
                               value='America',
                               inplace=True)
```

Now that we have all the data we want in a single DataFrame, we will reshape it back into panel form with a Multi-Index.

We should also ensure to sort the index using `.sort_index()` so that we can efficiently filter our dataframe later on.

By default, levels will be sorted top-down

```
merged = merged.set_index(['Continent', 'Country']).sort_index()
merged.head()
```

```
Continent Country      2006-01-01  2007-01-01  2008-01-01  ...  2014-01-01  \
America   Brazil         0.87         0.92         0.96  ...         1.21
          Canada         6.89         6.96         7.24  ...         8.22
          Chile          1.42         1.45         1.44  ...         1.76
          Colombia        1.01         1.02         1.01  ...         1.13
          Costa Rica        nan          nan          nan  ...         2.41

Continent Country      2015-01-01  2016-01-01
America   Brazil         1.21         1.24
          Canada         8.35         8.48
          Chile          1.81         1.91
          Colombia        1.13         1.12
          Costa Rica        2.56         2.63

[5 rows x 11 columns]
```

While merging, we lost our `DatetimeIndex`, as we merged columns that were not in datetime format

```
merged.columns
```

```
Index([2006-01-01 00:00:00, 2007-01-01 00:00:00, 2008-01-01 00:00:00,
       2009-01-01 00:00:00, 2010-01-01 00:00:00, 2011-01-01 00:00:00,
       2012-01-01 00:00:00, 2013-01-01 00:00:00, 2014-01-01 00:00:00,
       2015-01-01 00:00:00, 2016-01-01 00:00:00],
      dtype='object')
```

Now that we have set the merged columns as the index, we can recreate a `DatetimeIndex` using `.to_datetime()`

```
merged.columns = pd.to_datetime(merged.columns)
merged.columns = merged.columns.rename('Time')
merged.columns
```

```
DatetimeIndex(['2006-01-01', '2007-01-01', '2008-01-01', '2009-01-01',
               '2010-01-01', '2011-01-01', '2012-01-01', '2013-01-01',
               '2014-01-01', '2015-01-01', '2016-01-01'],
              dtype='datetime64[ns]', name='Time', freq=None)
```

The `DatetimeIndex` tends to work more smoothly in the row axis, so we will go ahead and transpose `merged`

```
merged = merged.transpose()
merged.head()
```

```
Continent  America      ...  Europe
Country    Brazil Canada Chile ... Slovenia Spain United Kingdom
Time
2006-01-01    0.87    6.89    1.42 ...    3.92    3.99            9.81
2007-01-01    0.92    6.96    1.45 ...    3.88    4.10            10.07
2008-01-01    0.96    7.24    1.44 ...    3.96    4.14            10.04
2009-01-01    1.03    7.67    1.52 ...    4.08    4.32            10.15
2010-01-01    1.08    7.94    1.56 ...    4.81    4.30            9.96

[5 rows x 32 columns]
```

62.4 Grouping and Summarizing Data

Grouping and summarizing data can be particularly useful for understanding large panel datasets.

A simple way to summarize data is to call an [aggregation method](#) on the dataframe, such as `.mean()` or `.max()`.

For example, we can calculate the average real minimum wage for each country over the period 2006 to 2016 (the default is to aggregate over rows)

```
merged.mean().head(10)
```

```
Continent  Country
America    Brazil      1.09
           Canada      7.82
           Chile       1.62
           Colombia     1.07
           Costa Rica   2.53
           Mexico       0.53
           United States 7.15
Asia        Israel      5.95
           Japan        6.18
           Korea        4.22
dtype: float64
```

Using this series, we can plot the average real minimum wage over the past decade for each country in our data set

```
%matplotlib inline
import matplotlib.pyplot as plt
plt.rcParams["figure.figsize"] = (11, 5) #set default figure size
import matplotlib
matplotlib.style.use('seaborn')

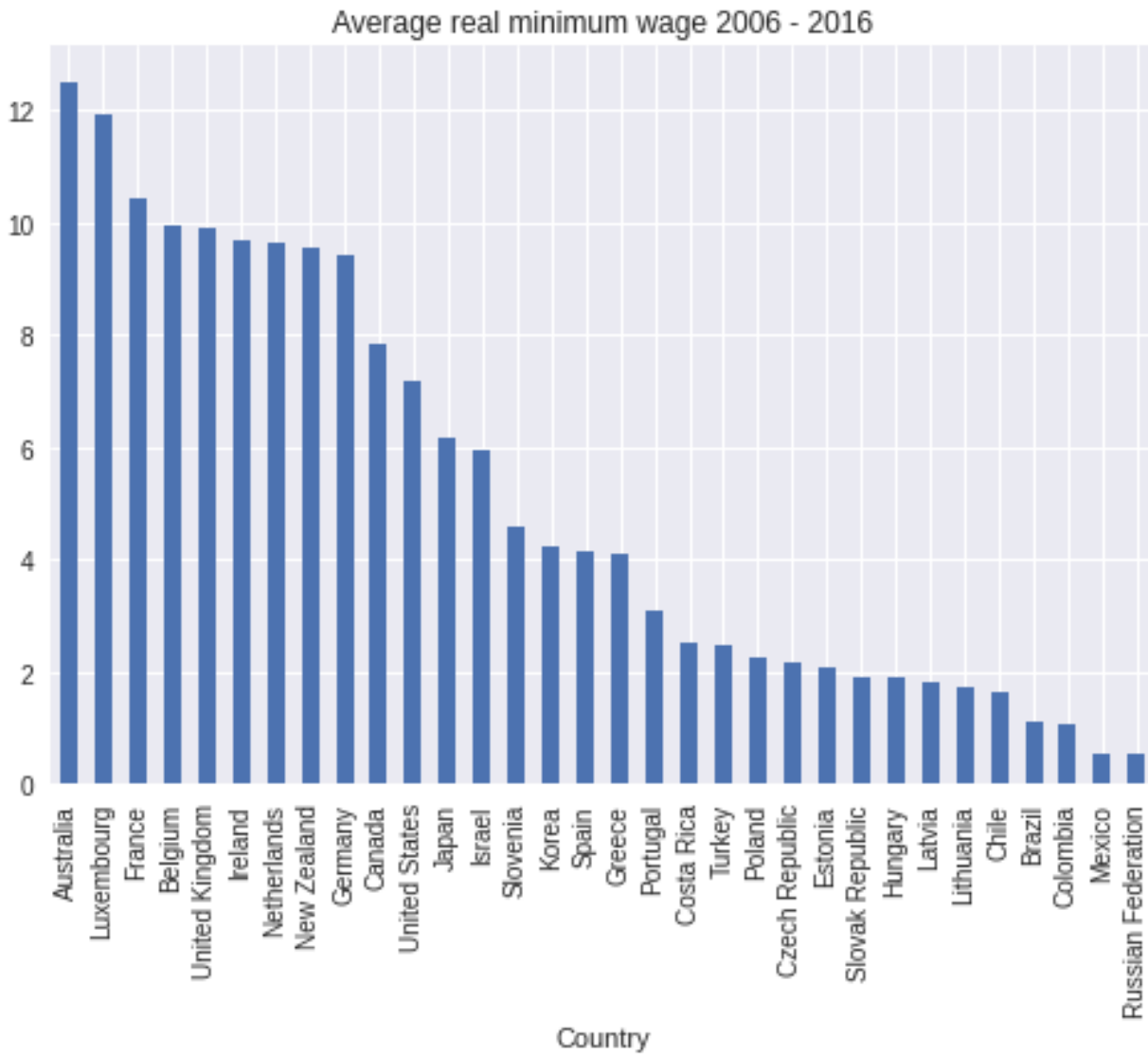
merged.mean().sort_values(ascending=False).plot(kind='bar', title="Average real_
↪ minimum wage 2006 - 2016")

#Set country labels
country_labels = merged.mean().sort_values(ascending=False).index.get_level_values(
↪ 'Country').tolist()
plt.xticks(range(0, len(country_labels)), country_labels)
plt.xlabel('Country')
```

(continues on next page)

(continued from previous page)

```
plt.show()
```



Passing in `axis=1` to `.mean()` will aggregate over columns (giving the average minimum wage for all countries over time)

```
merged.mean(axis=1).head()
```

```
Time
2006-01-01    4.69
2007-01-01    4.84
2008-01-01    4.90
2009-01-01    5.08
2010-01-01    5.11
dtype: float64
```

We can plot this time series as a line graph

```
merged.mean(axis=1).plot()
plt.title('Average real minimum wage 2006 - 2016')
plt.ylabel('2015 USD')
plt.xlabel('Year')
plt.show()
```



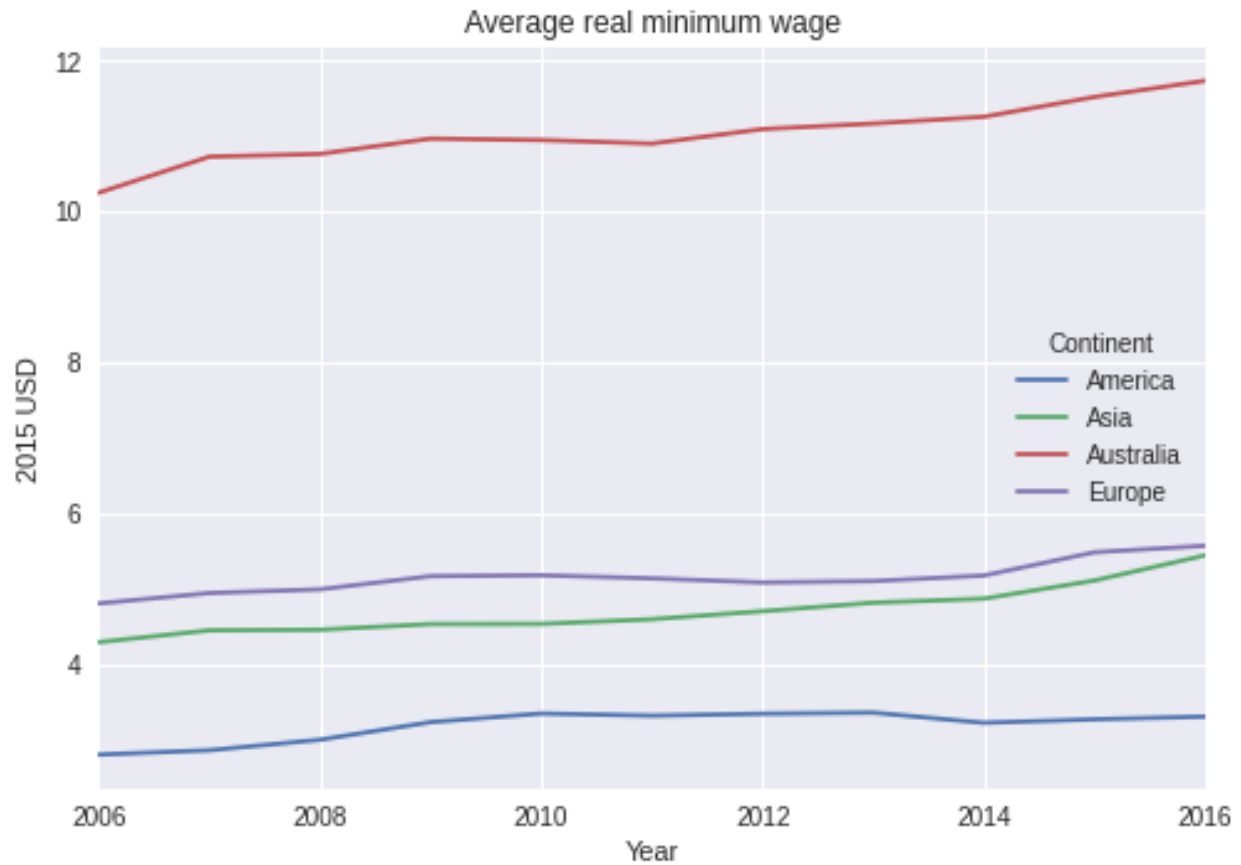
We can also specify a level of the MultiIndex (in the column axis) to aggregate over

```
merged.mean(level='Continent', axis=1).head()
```

Continent	America	Asia	Australia	Europe
Time				
2006-01-01	2.80	4.29	10.25	4.80
2007-01-01	2.85	4.44	10.73	4.94
2008-01-01	2.99	4.45	10.76	4.99
2009-01-01	3.23	4.53	10.97	5.16
2010-01-01	3.34	4.53	10.95	5.17

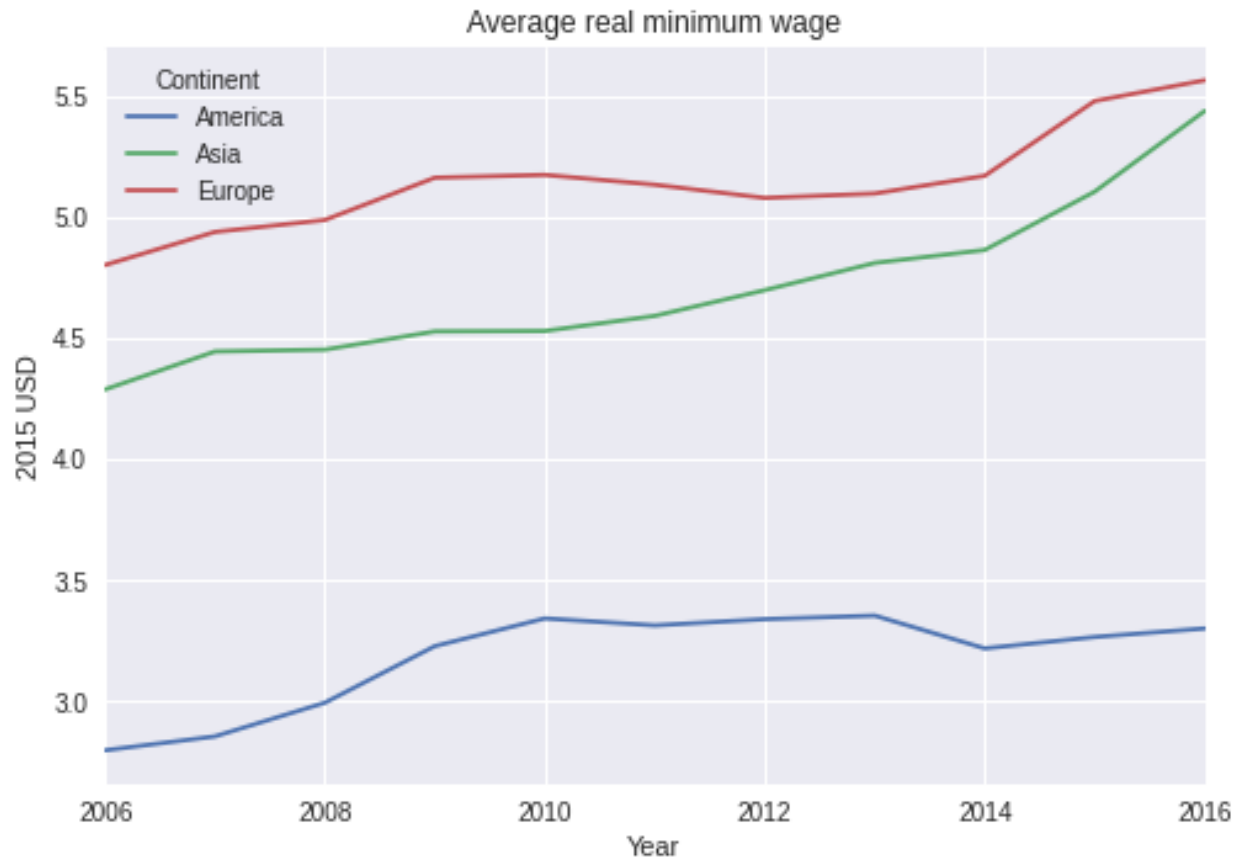
We can plot the average minimum wages in each continent as a time series

```
merged.mean(level='Continent', axis=1).plot()
plt.title('Average real minimum wage')
plt.ylabel('2015 USD')
plt.xlabel('Year')
plt.show()
```



We will drop Australia as a continent for plotting purposes

```
merged = merged.drop('Australia', level='Continent', axis=1)
merged.mean(level='Continent', axis=1).plot()
plt.title('Average real minimum wage')
plt.ylabel('2015 USD')
plt.xlabel('Year')
plt.show()
```

`.describe()` is useful for quickly retrieving a number of common summary statistics

```
merged.stack().describe()
```

Continent	America	Asia	Europe
count	69.00	44.00	200.00
mean	3.19	4.70	5.15
std	3.02	1.56	3.82
min	0.52	2.22	0.23
25%	1.03	3.37	2.02
50%	1.44	5.48	3.54
75%	6.96	5.95	9.70
max	8.48	6.65	12.39

This is a simplified way to use `groupby`.

Using `groupby` generally follows a ‘split-apply-combine’ process:

- split: data is grouped based on one or more keys
- apply: a function is called on each group independently
- combine: the results of the function calls are combined into a new data structure

The `groupby` method achieves the first step of this process, creating a new `DataFrameGroupBy` object with data split into groups.

Let’s split `merged` by continent again, this time using the `groupby` function, and name the resulting object `grouped`

```
grouped = merged.groupby(level='Continent', axis=1)
grouped
```

```
<pandas.core.groupby.generic.DataFrameGroupBy object at 0x7fc17950dfa0>
```

Calling an aggregation method on the object applies the function to each group, the results of which are combined in a new data structure.

For example, we can return the number of countries in our dataset for each continent using `.size()`.

In this case, our new data structure is a `Series`

```
grouped.size()
```

```
Continent
America      7
Asia         4
Europe      19
dtype: int64
```

Calling `.get_group()` to return just the countries in a single group, we can create a kernel density estimate of the distribution of real minimum wages in 2016 for each continent.

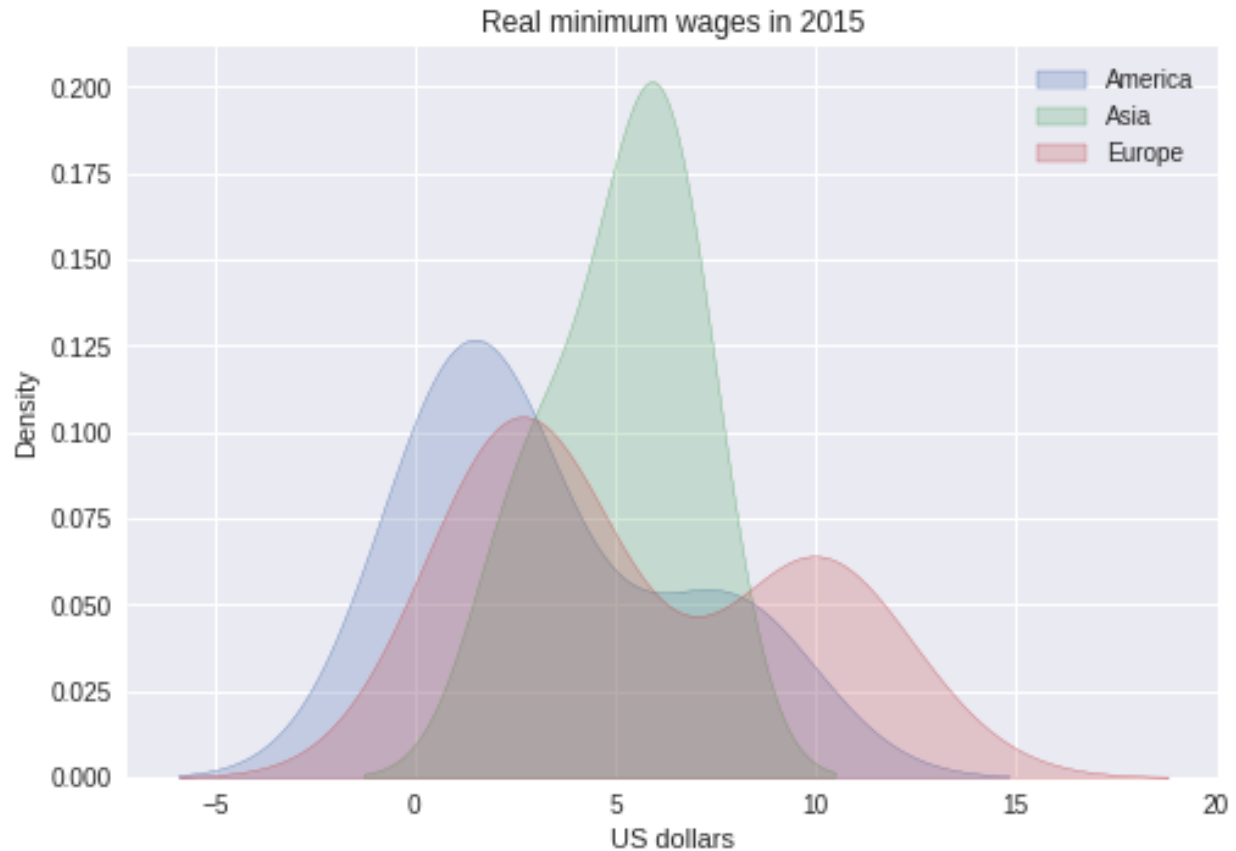
`grouped.groups.keys()` will return the keys from the `groupby` object

```
import seaborn as sns

continents = grouped.groups.keys()

for continent in continents:
    sns.kdeplot(grouped.get_group(continent)['2015'].unstack(), label=continent,
               shade=True)

plt.title('Real minimum wages in 2015')
plt.xlabel('US dollars')
plt.legend()
plt.show()
```



62.5 Final Remarks

This lecture has provided an introduction to some of pandas' more advanced features, including multiindices, merging, grouping and plotting.

Other tools that may be useful in panel data analysis include `xarray`, a python package that extends pandas to N-dimensional data structures.

62.6 Exercises

62.6.1 Exercise 1

In these exercises, you'll work with a dataset of employment rates in Europe by age and sex from [Eurostat](#).

The dataset can be accessed with the following link:

```
url3 = 'https://raw.githubusercontent.com/QuantEcon/lecture-python/master/source/_static/lecture_specific/pandas_panel/employ.csv'
```

Reading in the CSV file returns a panel dataset in long format. Use `.pivot_table()` to construct a wide format dataframe with a `MultiIndex` in the columns.

Start off by exploring the dataframe and the variables available in the `MultiIndex` levels.

Write a program that quickly returns all values in the `MultiIndex`.

62.6.2 Exercise 2

Filter the above dataframe to only include employment as a percentage of ‘active population’.

Create a grouped boxplot using `seaborn` of employment rates in 2015 by age group and sex.

Hint: `GEO` includes both areas and countries.

62.7 Solutions

62.7.1 Exercise 1

```
employ = pd.read_csv(url3)
employ = employ.pivot_table(values='Value',
                             index=['DATE'],
                             columns=['UNIT', 'AGE', 'SEX', 'INDIC_EM', 'GEO'])
employ.index = pd.to_datetime(employ.index) # ensure that dates are datetime format
employ.head()
```

UNIT	Percentage of total population	...	\
AGE	From 15 to 24 years	...	
SEX	Females	...	
INDIC_EM	Active population	...	
GEO	Austria Belgium Bulgaria	...	
DATE		...	
2007-01-01	56.00	31.60	26.00
2008-01-01	56.20	30.80	26.10
2009-01-01	56.20	29.90	24.80
2010-01-01	54.00	29.80	26.60
2011-01-01	54.80	29.80	24.80

UNIT	Thousand persons	...	\
AGE	From 55 to 64 years	...	
SEX	Total	...	
INDIC_EM	Total employment (resident population concept - LFS)	...	
GEO	Switzerland Turkey	...	
DATE		...	
2007-01-01	nan	1,282.00	
2008-01-01	nan	1,354.00	
2009-01-01	nan	1,449.00	
2010-01-01	640.00	1,583.00	
2011-01-01	661.00	1,760.00	

UNIT		
AGE		
SEX		
INDIC_EM		
GEO	United Kingdom	
DATE		
2007-01-01	4,131.00	
2008-01-01	4,204.00	
2009-01-01	4,193.00	

(continues on next page)

(continued from previous page)

```
2010-01-01      4,186.00
2011-01-01      4,164.00

[5 rows x 1440 columns]
```

This is a large dataset so it is useful to explore the levels and variables available

```
employ.columns.names
```

```
FrozenList(['UNIT', 'AGE', 'SEX', 'INDIC_EM', 'GEO'])
```

Variables within levels can be quickly retrieved with a loop

```
for name in employ.columns.names:
    print(name, employ.columns.get_level_values(name).unique())
```

```
UNIT Index(['Percentage of total population', 'Thousand persons'], dtype='object',
↳name='UNIT')
AGE Index(['From 15 to 24 years', 'From 25 to 54 years', 'From 55 to 64 years'],
↳dtype='object', name='AGE')
SEX Index(['Females', 'Males', 'Total'], dtype='object', name='SEX')
INDIC_EM Index(['Active population', 'Total employment (resident population concept -
↳LFS)'], dtype='object', name='INDIC_EM')
GEO Index(['Austria', 'Belgium', 'Bulgaria', 'Croatia', 'Cyprus', 'Czech Republic',
'Denmark', 'Estonia', 'Euro area (17 countries)',
'Euro area (18 countries)', 'Euro area (19 countries)',
'European Union (15 countries)', 'European Union (27 countries)',
'European Union (28 countries)', 'Finland',
'Former Yugoslav Republic of Macedonia, the', 'France',
'France (metropolitan)',
'Germany (until 1990 former territory of the FRG)', 'Greece', 'Hungary',
'Iceland', 'Ireland', 'Italy', 'Latvia', 'Lithuania', 'Luxembourg',
'Malta', 'Netherlands', 'Norway', 'Poland', 'Portugal', 'Romania',
'Slovakia', 'Slovenia', 'Spain', 'Sweden', 'Switzerland', 'Turkey',
'United Kingdom'],
dtype='object', name='GEO')
```

62.7.2 Exercise 2

To easily filter by country, swap GEO to the top level and sort the MultiIndex

```
employ.columns = employ.columns.swaplevel(0,-1)
employ = employ.sort_index(axis=1)
```

We need to get rid of a few items in GEO which are not countries.

A fast way to get rid of the EU areas is to use a list comprehension to find the level values in GEO that begin with 'Euro'

```
geo_list = employ.columns.get_level_values('GEO').unique().tolist()
countries = [x for x in geo_list if not x.startswith('Euro')]
employ = employ[countries]
employ.columns.get_level_values('GEO').unique()
```

```
Index(['Austria', 'Belgium', 'Bulgaria', 'Croatia', 'Cyprus', 'Czech Republic',
      'Denmark', 'Estonia', 'Finland',
      'Former Yugoslav Republic of Macedonia, the', 'France',
      'France (metropolitan)',
      'Germany (until 1990 former territory of the FRG)', 'Greece', 'Hungary',
      'Iceland', 'Ireland', 'Italy', 'Latvia', 'Lithuania', 'Luxembourg',
      'Malta', 'Netherlands', 'Norway', 'Poland', 'Portugal', 'Romania',
      'Slovakia', 'Slovenia', 'Spain', 'Sweden', 'Switzerland', 'Turkey',
      'United Kingdom'],
      dtype='object', name='GEO')
```

Select only percentage employed in the active population from the dataframe

```
employ_f = employ.xs(('Percentage of total population', 'Active population'),
                    level=('UNIT', 'INDIC_EM'),
                    axis=1)
employ_f.head()
```

```
GEO          Austria      ...      United Kingdom      \
AGE      From 15 to 24 years      ...      From 55 to 64 years
SEX          Females Males Total      ...          Females Males
DATE
2007-01-01      56.00 62.90 59.40      ...          49.90 68.90
2008-01-01      56.20 62.90 59.50      ...          50.20 69.80
2009-01-01      56.20 62.90 59.50      ...          50.60 70.30
2010-01-01      54.00 62.60 58.30      ...          51.10 69.20
2011-01-01      54.80 63.60 59.20      ...          51.30 68.40

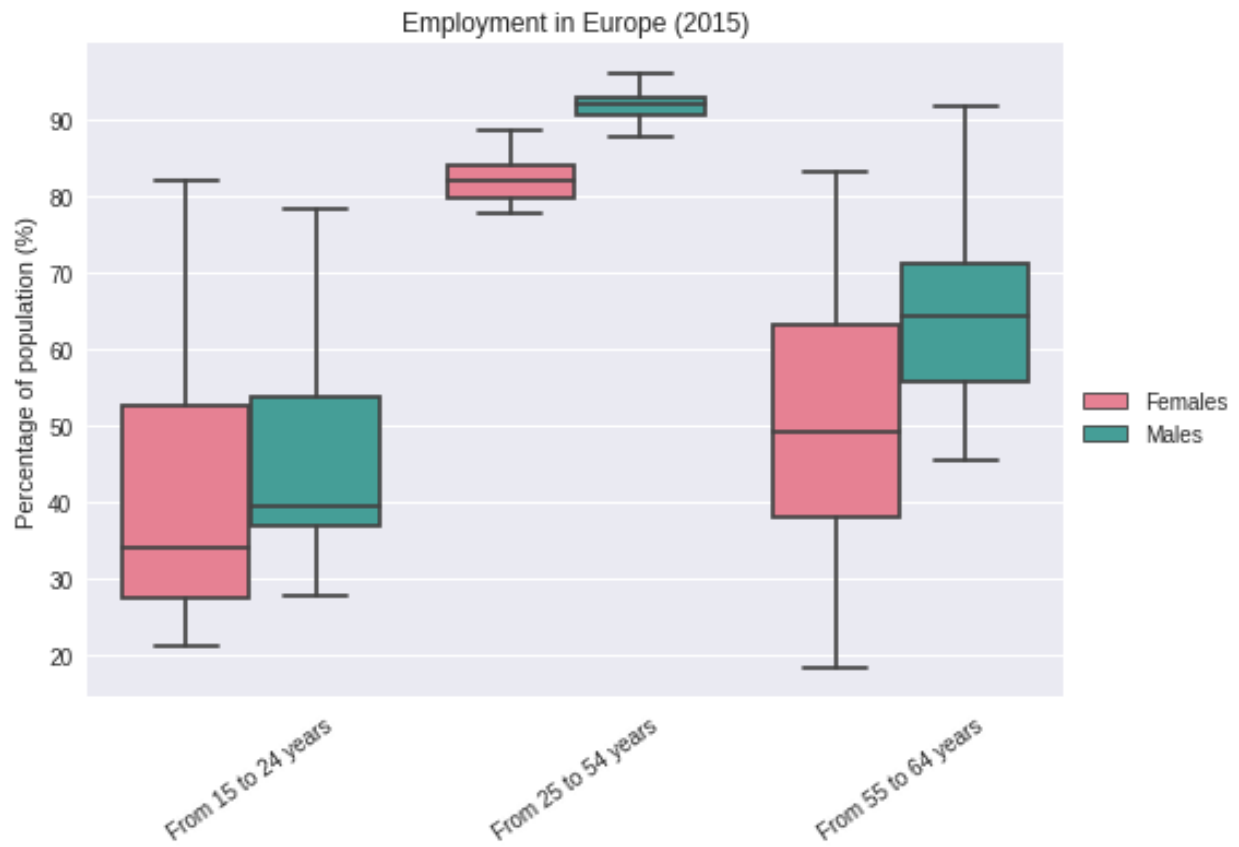
GEO
AGE
SEX      Total
DATE
2007-01-01  59.30
2008-01-01  59.80
2009-01-01  60.30
2010-01-01  60.00
2011-01-01  59.70

[5 rows x 306 columns]
```

Drop the 'Total' value before creating the grouped boxplot

```
employ_f = employ_f.drop('Total', level='SEX', axis=1)
```

```
box = employ_f['2015'].unstack().reset_index()
sns.boxplot(x="AGE", y=0, hue="SEX", data=box, palette="husl", showfliers=False)
plt.xlabel('')
plt.xticks(rotation=35)
plt.ylabel('Percentage of population (%)')
plt.title('Employment in Europe (2015)')
plt.legend(bbox_to_anchor=(1,0.5))
plt.show()
```



LINEAR REGRESSION IN PYTHON

Contents

- *Linear Regression in Python*
 - *Overview*
 - *Simple Linear Regression*
 - *Extending the Linear Regression Model*
 - *Endogeneity*
 - *Summary*
 - *Exercises*
 - *Solutions*

In addition to what's in Anaconda, this lecture will need the following libraries:

```
!pip install linearmodels
```

63.1 Overview

Linear regression is a standard tool for analyzing the relationship between two or more variables.

In this lecture, we'll use the Python package `statsmodels` to estimate, interpret, and visualize linear regression models.

Along the way, we'll discuss a variety of topics, including

- simple and multivariate linear regression
- visualization
- endogeneity and omitted variable bias
- two-stage least squares

As an example, we will replicate results from Acemoglu, Johnson and Robinson's seminal paper [AJR01].

- You can download a copy [here](#).

In the paper, the authors emphasize the importance of institutions in economic development.

The main contribution is the use of settler mortality rates as a source of *exogenous* variation in institutional differences.

Such variation is needed to determine whether it is institutions that give rise to greater economic growth, rather than the other way around.

Let's start with some imports:

```
%matplotlib inline
import matplotlib.pyplot as plt
plt.rcParams["figure.figsize"] = (11, 5) #set default figure size
import numpy as np
import pandas as pd
import statsmodels.api as sm
from statsmodels.iolib.summary2 import summary_col
from linearmodels.iv import IV2SLS
```

63.1.1 Prerequisites

This lecture assumes you are familiar with basic econometrics.

For an introductory text covering these topics, see, for example, [Woo15].

63.2 Simple Linear Regression

[AJR01] wish to determine whether or not differences in institutions can help to explain observed economic outcomes.

How do we measure *institutional differences* and *economic outcomes*?

In this paper,

- economic outcomes are proxied by log GDP per capita in 1995, adjusted for exchange rates.
- institutional differences are proxied by an index of protection against expropriation on average over 1985-95, constructed by the [Political Risk Services Group](#).

These variables and other data used in the paper are available for download on Daron Acemoglu's [webpage](#).

We will use pandas' `.read_stata()` function to read in data contained in the `.dta` files to dataframes

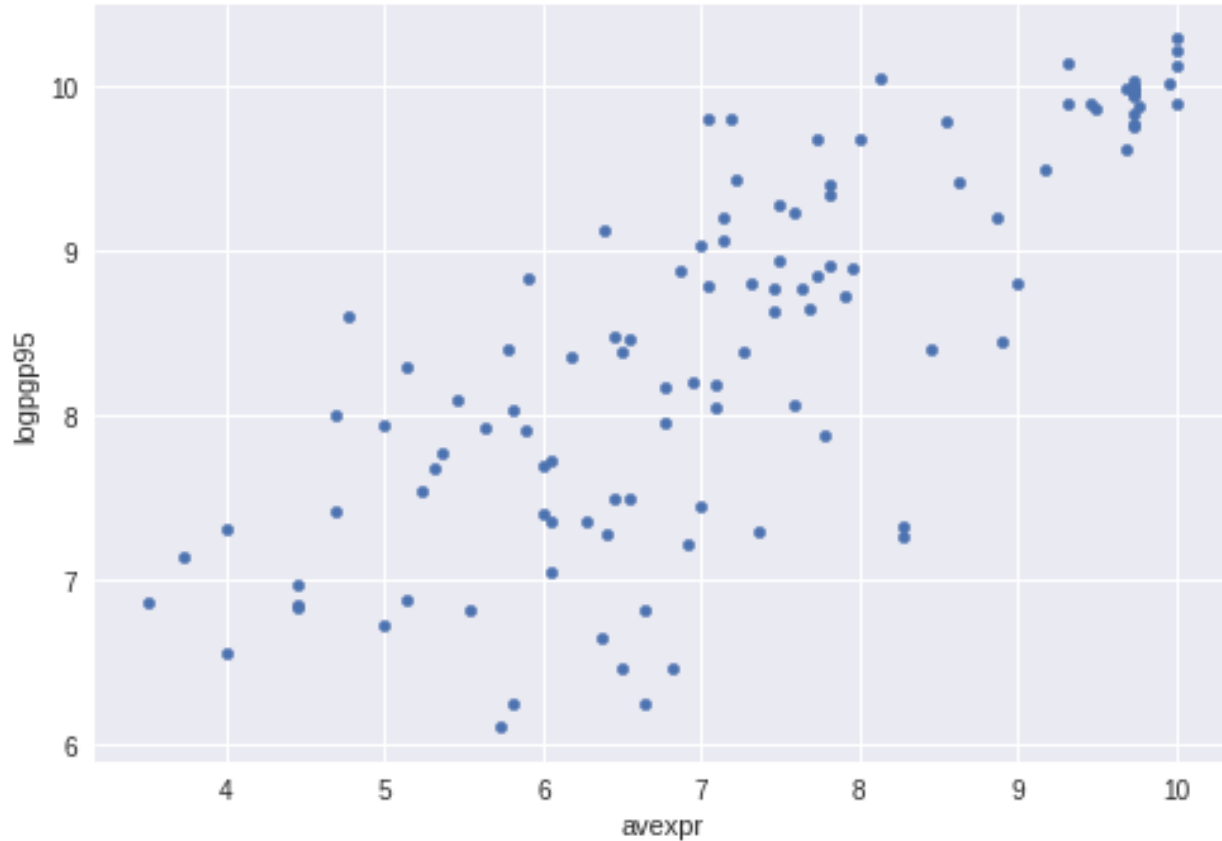
```
df1 = pd.read_stata('https://github.com/QuantEcon/lecture-python/blob/master/source/_
static/lecture_specific/ols/maketable1.dta?raw=true')
df1.head()
```

	shortnam	euro1900	excolony	avexpr	logpgp95	cons1	cons90	democ00a	\
0	AFG	0.000000	1.0	NaN	NaN	1.0	2.0	1.0	
1	AGO	8.000000	1.0	5.363636	7.770645	3.0	3.0	0.0	
2	ARE	0.000000	1.0	7.181818	9.804219	NaN	NaN	NaN	
3	ARG	60.000004	1.0	6.386364	9.133459	1.0	6.0	3.0	
4	ARM	0.000000	0.0	NaN	7.682482	NaN	NaN	NaN	
	cons00a	extmort4	logem4	loghjypl	baseco				
0	1.0	93.699997	4.540098	NaN	NaN				
1	1.0	280.000000	5.634789	-3.411248	1.0				
2	NaN	NaN	NaN	NaN	NaN				
3	3.0	68.900002	4.232656	-0.872274	1.0				
4	NaN	NaN	NaN	NaN	NaN				

Let's use a scatterplot to see whether any obvious relationship exists between GDP per capita and the protection against expropriation index

```
plt.style.use('seaborn')

df1.plot(x='avexpr', y='logpgp95', kind='scatter')
plt.show()
```



The plot shows a fairly strong positive relationship between protection against expropriation and log GDP per capita.

Specifically, if higher protection against expropriation is a measure of institutional quality, then better institutions appear to be positively correlated with better economic outcomes (higher GDP per capita).

Given the plot, choosing a linear model to describe this relationship seems like a reasonable assumption.

We can write our model as

$$\logpgp95_i = \beta_0 + \beta_1 \text{avexpr}_i + u_i$$

where:

- β_0 is the intercept of the linear trend line on the y-axis
- β_1 is the slope of the linear trend line, representing the *marginal effect* of protection against risk on log GDP per capita
- u_i is a random error term (deviations of observations from the linear trend due to factors not included in the model)

Visually, this linear model involves choosing a straight line that best fits the data, as in the following plot (Figure 2 in [AJR01])

```
# Dropping NA's is required to use numpy's polyfit
df1_subset = df1.dropna(subset=['logpgp95', 'avexpr'])

# Use only 'base sample' for plotting purposes
df1_subset = df1_subset[df1_subset['baseco'] == 1]

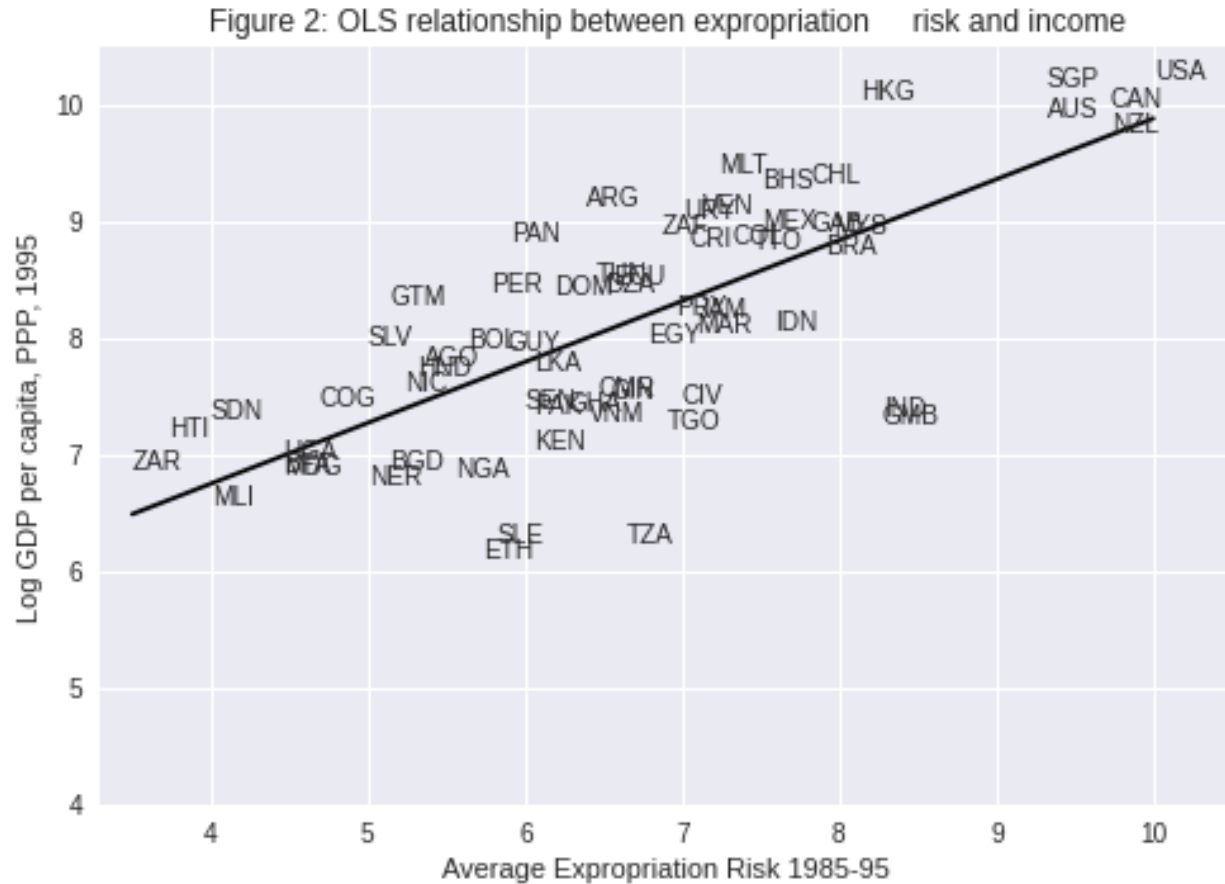
X = df1_subset['avexpr']
y = df1_subset['logpgp95']
labels = df1_subset['shortnam']

# Replace markers with country labels
fig, ax = plt.subplots()
ax.scatter(X, y, marker='')

for i, label in enumerate(labels):
    ax.annotate(label, (X.iloc[i], y.iloc[i]))

# Fit a linear trend line
ax.plot(np.unique(X),
        np.poly1d(np.polyfit(X, y, 1))(np.unique(X)),
        color='black')

ax.set_xlim([3.3, 10.5])
ax.set_ylim([4, 10.5])
ax.set_xlabel('Average Expropriation Risk 1985-95')
ax.set_ylabel('Log GDP per capita, PPP, 1995')
ax.set_title('Figure 2: OLS relationship between expropriation \
             risk and income')
plt.show()
```



The most common technique to estimate the parameters (β 's) of the linear model is Ordinary Least Squares (OLS).

As the name implies, an OLS model is solved by finding the parameters that minimize *the sum of squared residuals*, i.e.

$$\min_{\hat{\beta}} \sum_{i=1}^N \hat{u}_i^2$$

where \hat{u}_i is the difference between the observation and the predicted value of the dependent variable.

To estimate the constant term β_0 , we need to add a column of 1's to our dataset (consider the equation if β_0 was replaced with $\beta_0 x_i$ and $x_i = 1$)

```
df1['const'] = 1
```

Now we can construct our model in statsmodels using the OLS function.

We will use pandas dataframes with statsmodels, however standard arrays can also be used as arguments

```
reg1 = sm.OLS(endog=df1['logppp95'], exog=df1[['const', 'avexpr']], \
              missing='drop')
type(reg1)
```

```
statsmodels.regression.linear_model.OLS
```

So far we have simply constructed our model.

We need to use `.fit()` to obtain parameter estimates $\hat{\beta}_0$ and $\hat{\beta}_1$

```
results = reg1.fit()
type(results)
```

```
statsmodels.regression.linear_model.RegressionResultsWrapper
```

We now have the fitted regression model stored in `results`.

To view the OLS regression results, we can call the `.summary()` method.

Note that an observation was mistakenly dropped from the results in the original paper (see the note located in [maketable2.do](#) from Acemoglu's webpage), and thus the coefficients differ slightly.

```
print(results.summary())
```

```

=====
                        OLS Regression Results
=====
Dep. Variable:          logpgp95      R-squared:                0.611
Model:                  OLS           Adj. R-squared:           0.608
Method:                 Least Squares  F-statistic:              171.4
Date:                  Thu, 07 Oct 2021  Prob (F-statistic):       4.16e-24
Time:                  21:22:23        Log-Likelihood:           -119.71
No. Observations:      111            AIC:                     243.4
Df Residuals:          109            BIC:                     248.8
Df Model:              1
Covariance Type:       nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	4.6261	0.301	15.391	0.000	4.030	5.222
avexpr	0.5319	0.041	13.093	0.000	0.451	0.612

```

=====
Omnibus:                 9.251    Durbin-Watson:              1.689
Prob(Omnibus):           0.010    Jarque-Bera (JB):         9.170
Skew:                   -0.680    Prob(JB):                 0.0102
Kurtosis:                3.362    Cond. No.:                33.2
=====
Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly
    specified.

```

From our results, we see that

- The intercept $\hat{\beta}_0 = 4.63$.
- The slope $\hat{\beta}_1 = 0.53$.
- The positive $\hat{\beta}_1$ parameter estimate implies that institutional quality has a positive effect on economic outcomes, as we saw in the figure.
- The p-value of 0.000 for $\hat{\beta}_1$ implies that the effect of institutions on GDP is statistically significant (using $p < 0.05$ as a rejection rule).
- The R-squared value of 0.611 indicates that around 61% of variation in log GDP per capita is explained by protection against expropriation.

Using our parameter estimates, we can now write our estimated relationship as

$$\widehat{\logpgp95}_i = 4.63 + 0.53 \text{ avexpr}_i$$

This equation describes the line that best fits our data, as shown in Figure 2.

We can use this equation to predict the level of log GDP per capita for a value of the index of expropriation protection.

For example, for a country with an index value of 7.07 (the average for the dataset), we find that their predicted level of log GDP per capita in 1995 is 8.38.

```
mean_expr = np.mean(df1_subset['avexpr'])
mean_expr
```

```
6.515625
```

```
predicted_logpdp95 = 4.63 + 0.53 * 7.07
predicted_logpdp95
```

```
8.3771
```

An easier (and more accurate) way to obtain this result is to use `.predict()` and set `constant = 1` and `avexpri = mean_expr`

```
results.predict(exog=[1, mean_expr])
```

```
array([8.09156367])
```

We can obtain an array of predicted $\logpgp95_i$ for every value of $avexpr_i$ in our dataset by calling `.predict()` on our results.

Plotting the predicted values against $avexpr_i$ shows that the predicted values lie along the linear line that we fitted above.

The observed values of $\logpgp95_i$ are also plotted for comparison purposes

```
# Drop missing observations from whole sample

df1_plot = df1.dropna(subset=['logpgp95', 'avexpr'])

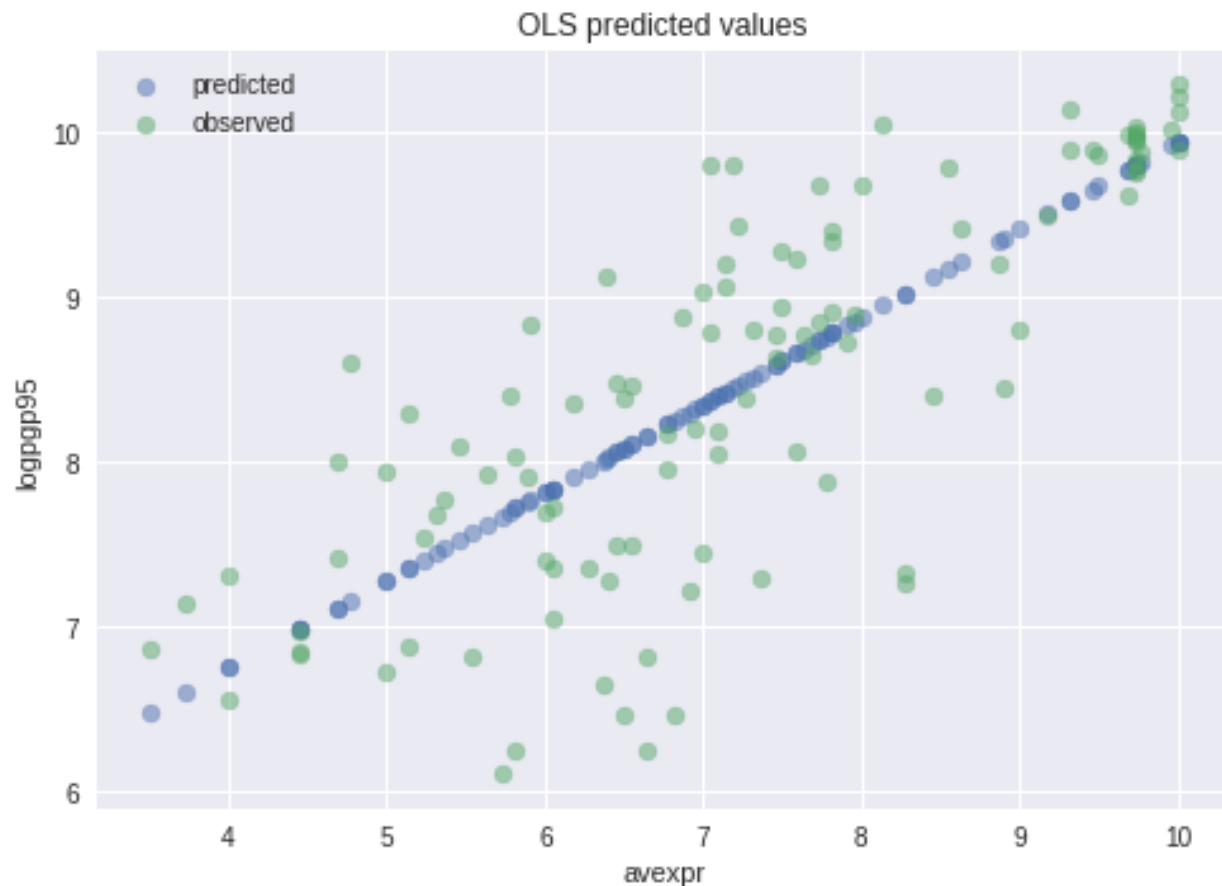
# Plot predicted values

fig, ax = plt.subplots()
ax.scatter(df1_plot['avexpr'], results.predict(), alpha=0.5,
          label='predicted')

# Plot observed values

ax.scatter(df1_plot['avexpr'], df1_plot['logpgp95'], alpha=0.5,
          label='observed')

ax.legend()
ax.set_title('OLS predicted values')
ax.set_xlabel('avexpr')
ax.set_ylabel('logpgp95')
plt.show()
```



63.3 Extending the Linear Regression Model

So far we have only accounted for institutions affecting economic performance - almost certainly there are numerous other factors affecting GDP that are not included in our model.

Leaving out variables that affect $\logpgp95_i$ will result in **omitted variable bias**, yielding biased and inconsistent parameter estimates.

We can extend our bivariate regression model to a **multivariate regression model** by adding in other factors that may affect $\logpgp95_i$.

[AJR01] consider other factors such as:

- the effect of climate on economic outcomes; latitude is used to proxy this
- differences that affect both economic performance and institutions, eg. cultural, historical, etc.; controlled for with the use of continent dummies

Let's estimate some of the extended models considered in the paper (Table 2) using data from `maketable2.dta`

```
df2 = pd.read_stata('https://github.com/QuantEcon/lecture-python/blob/master/source/_
↳static/lecture_specific/ols/maketable2.dta?raw=true')

# Add constant term to dataset
df2['const'] = 1
```

(continues on next page)

(continued from previous page)

```
# Create lists of variables to be used in each regression
X1 = ['const', 'avexpr']
X2 = ['const', 'avexpr', 'lat_abst']
X3 = ['const', 'avexpr', 'lat_abst', 'asia', 'africa', 'other']

# Estimate an OLS regression for each set of variables
reg1 = sm.OLS(df2['logpgp95'], df2[X1], missing='drop').fit()
reg2 = sm.OLS(df2['logpgp95'], df2[X2], missing='drop').fit()
reg3 = sm.OLS(df2['logpgp95'], df2[X3], missing='drop').fit()
```

Now that we have fitted our model, we will use `summary_col` to display the results in a single table (model numbers correspond to those in the paper)

```
info_dict={'R-squared' : lambda x: f"{x.rsquared:.2f}",
          'No. observations' : lambda x: f"{int(x.nobs):d}"}

results_table = summary_col(results=[reg1, reg2, reg3],
                             float_format='%0.2f',
                             stars = True,
                             model_names=['Model 1',
                                           'Model 3',
                                           'Model 4'],
                             info_dict=info_dict,
                             regressor_order=['const',
                                              'avexpr',
                                              'lat_abst',
                                              'asia',
                                              'africa'])

results_table.add_title('Table 2 - OLS Regressions')

print(results_table)
```

```
Table 2 - OLS Regressions
=====
Model 1 Model 3 Model 4
-----
const      4.63*** 4.87*** 5.85***
           (0.30) (0.33) (0.34)
avexpr     0.53*** 0.46*** 0.39***
           (0.04) (0.06) (0.05)
lat_abst           0.87* 0.33
                (0.49) (0.45)
asia                        -0.15
                        (0.15)
africa                       -0.92***
                        (0.17)
other                        0.30
                        (0.37)
R-squared      0.61    0.62    0.72
R-squared Adj. 0.61    0.62    0.70
R-squared      0.61    0.62    0.72
No. observations 111    111    111
=====
Standard errors in parentheses.
```

(continues on next page)

(continued from previous page)

```
* p<.1, ** p<.05, ***p<.01
```

63.4 Endogeneity

As [AJR01] discuss, the OLS models likely suffer from **endogeneity** issues, resulting in biased and inconsistent model estimates.

Namely, there is likely a two-way relationship between institutions and economic outcomes:

- richer countries may be able to afford or prefer better institutions
- variables that affect income may also be correlated with institutional differences
- the construction of the index may be biased; analysts may be biased towards seeing countries with higher income having better institutions

To deal with endogeneity, we can use **two-stage least squares (2SLS) regression**, which is an extension of OLS regression.

This method requires replacing the endogenous variable $avexpr_i$ with a variable that is:

1. correlated with $avexpr_i$
2. not correlated with the error term (ie. it should not directly affect the dependent variable, otherwise it would be correlated with u_i due to omitted variable bias)

The new set of regressors is called an **instrument**, which aims to remove endogeneity in our proxy of institutional differences.

The main contribution of [AJR01] is the use of settler mortality rates to instrument for institutional differences.

They hypothesize that higher mortality rates of colonizers led to the establishment of institutions that were more extractive in nature (less protection against expropriation), and these institutions still persist today.

Using a scatterplot (Figure 3 in [AJR01]), we can see protection against expropriation is negatively correlated with settler mortality rates, coinciding with the authors' hypothesis and satisfying the first condition of a valid instrument.

```
# Dropping NA's is required to use numpy's polyfit
df1_subset2 = df1.dropna(subset=['logem4', 'avexpr'])

X = df1_subset2['logem4']
y = df1_subset2['avexpr']
labels = df1_subset2['shortnam']

# Replace markers with country labels
fig, ax = plt.subplots()
ax.scatter(X, y, marker='')

for i, label in enumerate(labels):
    ax.annotate(label, (X.iloc[i], y.iloc[i]))

# Fit a linear trend line
ax.plot(np.unique(X),
        np.poly1d(np.polyfit(X, y, 1))(np.unique(X)),
        color='black')

ax.set_xlim([1.8, 8.4])
ax.set_ylim([3.3, 10.4])
```

(continues on next page)

(continued from previous page)

```

ax.set_xlabel('Log of Settler Mortality')
ax.set_ylabel('Average Expropriation Risk 1985-95')
ax.set_title('Figure 3: First-stage relationship between settler mortality \
and expropriation risk')
plt.show()

```



The second condition may not be satisfied if settler mortality rates in the 17th to 19th centuries have a direct effect on current GDP (in addition to their indirect effect through institutions).

For example, settler mortality rates may be related to the current disease environment in a country, which could affect current economic performance.

[AJR01] argue this is unlikely because:

- The majority of settler deaths were due to malaria and yellow fever and had a limited effect on local people.
- The disease burden on local people in Africa or India, for example, did not appear to be higher than average, supported by relatively high population densities in these areas before colonization.

As we appear to have a valid instrument, we can use 2SLS regression to obtain consistent and unbiased parameter estimates.

First stage

The first stage involves regressing the endogenous variable ($avexpr_i$) on the instrument.

The instrument is the set of all exogenous variables in our model (and not just the variable we have replaced).

Using model 1 as an example, our instrument is simply a constant and settler mortality rates $logem4_i$.

Therefore, we will estimate the first-stage regression as

$$avexpr_i = \delta_0 + \delta_1 \logem4_i + v_i$$

The data we need to estimate this equation is located in `maketable4.dta` (only complete data, indicated by `baseco = 1`, is used for estimation)

```
# Import and select the data
df4 = pd.read_stata('https://github.com/QuantEcon/lecture-python/blob/master/source/_
↳static/lecture_specific/ols/maketable4.dta?raw=true')
df4 = df4[df4['baseco'] == 1]

# Add a constant variable
df4['const'] = 1

# Fit the first stage regression and print summary
results_fs = sm.OLS(df4['avexpr'],
                    df4[['const', 'logem4']],
                    missing='drop').fit()
print(results_fs.summary())
```

OLS Regression Results						
=====						
Dep. Variable:	avexpr	R-squared:	0.270			
Model:	OLS	Adj. R-squared:	0.258			
Method:	Least Squares	F-statistic:	22.95			
Date:	Thu, 07 Oct 2021	Prob (F-statistic):	1.08e-05			
Time:	21:22:25	Log-Likelihood:	-104.83			
No. Observations:	64	AIC:	213.7			
Df Residuals:	62	BIC:	218.0			
Df Model:	1					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

const	9.3414	0.611	15.296	0.000	8.121	10.562
logem4	-0.6068	0.127	-4.790	0.000	-0.860	-0.354
=====						
Omnibus:	0.035	Durbin-Watson:	2.003			
Prob(Omnibus):	0.983	Jarque-Bera (JB):	0.172			
Skew:	0.045	Prob(JB):	0.918			
Kurtosis:	2.763	Cond. No.	19.4			
=====						
Notes:						
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.						

Second stage

We need to retrieve the predicted values of $avexpr_i$ using `.predict()`.

We then replace the endogenous variable $avexpr_i$ with the predicted values \widehat{avexpr}_i in the original linear model.

Our second stage regression is thus

$$\logpgp95_i = \beta_0 + \beta_1 \widehat{avexpr}_i + u_i$$

```
df4['predicted_avexpr'] = results_fs.predict()

results_ss = sm.OLS(df4['logpgp95'],
                    df4[['const', 'predicted_avexpr']]).fit()
print(results_ss.summary())
```

```

=====
                        OLS Regression Results
=====
Dep. Variable:          logpgp95      R-squared:                0.477
Model:                  OLS          Adj. R-squared:           0.469
Method:                 Least Squares  F-statistic:              56.60
Date:                   Thu, 07 Oct 2021  Prob (F-statistic):      2.66e-10
Time:                   21:22:25      Log-Likelihood:           -72.268
No. Observations:       64           AIC:                     148.5
Df Residuals:           62           BIC:                     152.9
Df Model:               1
Covariance Type:        nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	1.9097	0.823	2.320	0.024	0.264	3.555
predicted_avexpr	0.9443	0.126	7.523	0.000	0.693	1.195

```

=====
Omnibus:                 10.547      Durbin-Watson:           2.137
Prob(Omnibus):            0.005      Jarque-Bera (JB):        11.010
Skew:                    -0.790      Prob(JB):                0.00407
Kurtosis:                 4.277      Cond. No.                58.1
=====

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly
    specified.

```

The second-stage regression results give us an unbiased and consistent estimate of the effect of institutions on economic outcomes.

The result suggests a stronger positive relationship than what the OLS results indicated.

Note that while our parameter estimates are correct, our standard errors are not and for this reason, computing 2SLS ‘manually’ (in stages with OLS) is not recommended.

We can correctly estimate a 2SLS regression in one step using the `linearmodels` package, an extension of `statsmodels`

Note that when using IV2SLS, the exogenous and instrument variables are split up in the function arguments (whereas before the instrument included exogenous variables)

```
iv = IV2SLS(dependent=df4['logpgp95'],
            exog=df4['const'],
            endog=df4['avexpr'],
            instruments=df4['logem4']).fit(cov_type='unadjusted')

print(iv.summary)
```

```

=====
                        IV-2SLS Estimation Summary
=====
Dep. Variable:          logpgp95      R-squared:                0.1870
Estimator:              IV-2SLS      Adj. R-squared:           0.1739
No. Observations:       64           F-statistic:              37.568

```

(continues on next page)

(continued from previous page)

```

Date:          Thu, Oct 07 2021    P-value (F-stat)      0.0000
Time:          21:22:25           Distribution:        chi2(1)
Cov. Estimator:      unadjusted

```

Parameter Estimates

```

=====
Parameter  Std. Err.    T-stat    P-value    Lower CI    Upper CI
-----
const      1.9097      1.0106    1.8897    0.0588     -0.0710    3.8903
avexpr     0.9443      0.1541    6.1293    0.0000      0.6423    1.2462
=====

```

```

Endogenous: avexpr
Instruments: logem4
Unadjusted Covariance (Homoskedastic)
Debiased: False

```

Given that we now have consistent and unbiased estimates, we can infer from the model we have estimated that institutional differences (stemming from institutions set up during colonization) can help to explain differences in income levels across countries today.

[AJR01] use a marginal effect of 0.94 to calculate that the difference in the index between Chile and Nigeria (ie. institutional quality) implies up to a 7-fold difference in income, emphasizing the significance of institutions in economic development.

63.5 Summary

We have demonstrated basic OLS and 2SLS regression in `statsmodels` and `linearmodels`.

If you are familiar with R, you may want to use the [formula interface](#) to `statsmodels`, or consider using `r2py` to call R from within Python.

63.6 Exercises

63.6.1 Exercise 1

In the lecture, we think the original model suffers from endogeneity bias due to the likely effect income has on institutional development.

Although endogeneity is often best identified by thinking about the data and model, we can formally test for endogeneity using the **Hausman test**.

We want to test for correlation between the endogenous variable, $avexpr_i$, and the errors, u_i

$$\begin{aligned}
 H_0 : Cov(avexpr_i, u_i) &= 0 \quad (\text{no endogeneity}) \\
 H_1 : Cov(avexpr_i, u_i) &\neq 0 \quad (\text{endogeneity})
 \end{aligned}$$

This test is running in two stages.

First, we regress $avexpr_i$ on the instrument, $logem4_i$

$$avexpr_i = \pi_0 + \pi_1 logem4_i + v_i$$

Second, we retrieve the residuals \hat{v}_i and include them in the original equation

$$\logpgp95_i = \beta_0 + \beta_1 \text{avexpr}_i + \alpha \hat{v}_i + u_i$$

If α is statistically significant (with a p-value < 0.05), then we reject the null hypothesis and conclude that avexpr_i is endogenous.

Using the above information, estimate a Hausman test and interpret your results.

63.6.2 Exercise 2

The OLS parameter β can also be estimated using matrix algebra and `numpy` (you may need to review the [numpy](#) lecture to complete this exercise).

The linear equation we want to estimate is (written in matrix form)

$$y = X\beta + u$$

To solve for the unknown parameter β , we want to minimize the sum of squared residuals

$$\min_{\hat{\beta}} \hat{u}'\hat{u}$$

Rearranging the first equation and substituting into the second equation, we can write

$$\min_{\hat{\beta}} (Y - X\hat{\beta})'(Y - X\hat{\beta})$$

Solving this optimization problem gives the solution for the $\hat{\beta}$ coefficients

$$\hat{\beta} = (X'X)^{-1}X'y$$

Using the above information, compute $\hat{\beta}$ from model 1 using `numpy` - your results should be the same as those in the `statsmodels` output from earlier in the lecture.

63.7 Solutions

63.7.1 Exercise 1

```
# Load in data
df4 = pd.read_stata('https://github.com/QuantEcon/lecture-python/blob/master/source/_
↳static/lecture_specific/ols/maketable4.dta?raw=true')

# Add a constant term
df4['const'] = 1

# Estimate the first stage regression
reg1 = sm.OLS(endog=df4['avexpr'],
               exog=df4[['const', 'logem4']],
               missing='drop').fit()

# Retrieve the residuals
df4['resid'] = reg1.resid
```

(continues on next page)

(continued from previous page)

```
# Estimate the second stage residuals
reg2 = sm.OLS(endog=df4['logpgp95'],
              exog=df4[['const', 'avexpr', 'resid']],
              missing='drop').fit()

print(reg2.summary())
```

```

=====
                        OLS Regression Results
=====
Dep. Variable:          logpgp95      R-squared:                0.689
Model:                  OLS          Adj. R-squared:            0.679
Method:                 Least Squares  F-statistic:              74.05
Date:                  Thu, 07 Oct 2021  Prob (F-statistic):      1.07e-17
Time:                  21:22:25       Log-Likelihood:           -62.031
No. Observations:      70            AIC:                     130.1
Df Residuals:          67            BIC:                     136.8
Df Model:              2
Covariance Type:       nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	2.4782	0.547	4.530	0.000	1.386	3.570
avexpr	0.8564	0.082	10.406	0.000	0.692	1.021
resid	-0.4951	0.099	-5.017	0.000	-0.692	-0.298

```

=====
Omnibus:                 17.597      Durbin-Watson:           2.086
Prob(Omnibus):           0.000      Jarque-Bera (JB):        23.194
Skew:                   -1.054      Prob(JB):                9.19e-06
Kurtosis:                4.873      Cond. No.                53.8
=====

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly
    specified.

```

The output shows that the coefficient on the residuals is statistically significant, indicating $avexpr_i$ is endogenous.

63.7.2 Exercise 2

```
# Load in data
df1 = pd.read_stata('https://github.com/QuantEcon/lecture-python/blob/master/source/_
    static/lecture_specific/ols/maketable1.dta?raw=true')
df1 = df1.dropna(subset=['logpgp95', 'avexpr'])

# Add a constant term
df1['const'] = 1

# Define the X and y variables
y = np.asarray(df1['logpgp95'])
X = np.asarray(df1[['const', 'avexpr']])

# Compute  $\beta_{\text{hat}}$ 
 $\beta_{\text{hat}}$  = np.linalg.solve(X.T @ X, X.T @ y)
```

(continues on next page)

(continued from previous page)

```
# Print out the results from the 2 x 1 vector  $\beta_{\text{hat}}$ 
print(f' $\beta_0 = \{\beta_{\text{hat}}[0]:.2\}$ ')
print(f' $\beta_1 = \{\beta_{\text{hat}}[1]:.2\}$ ')
```

```
 $\beta_0 = 4.6$ 
 $\beta_1 = 0.53$ 
```

It is also possible to use `np.linalg.inv(X.T @ X) @ X.T @ y` to solve for β , however `.solve()` is preferred as it involves fewer computations.

MAXIMUM LIKELIHOOD ESTIMATION

Contents

- *Maximum Likelihood Estimation*
 - *Overview*
 - *Set Up and Assumptions*
 - *Conditional Distributions*
 - *Maximum Likelihood Estimation*
 - *MLE with Numerical Methods*
 - *Maximum Likelihood Estimation with `statsmodels`*
 - *Summary*
 - *Exercises*
 - *Solutions*

64.1 Overview

In a *previous lecture*, we estimated the relationship between dependent and explanatory variables using linear regression.

But what if a linear relationship is not an appropriate assumption for our model?

One widely used alternative is maximum likelihood estimation, which involves specifying a class of distributions, indexed by unknown parameters, and then using the data to pin down these parameter values.

The benefit relative to linear regression is that it allows more flexibility in the probabilistic relationships between variables.

Here we illustrate maximum likelihood by replicating Daniel Treisman's (2016) paper, [Russia's Billionaires](#), which connects the number of billionaires in a country to its economic characteristics.

The paper concludes that Russia has a higher number of billionaires than economic factors such as market size and tax rate predict.

We'll require the following imports:

```
%matplotlib inline
import matplotlib.pyplot as plt
plt.rcParams["figure.figsize"] = (11, 5)  #set default figure size
import numpy as np
```

(continues on next page)

(continued from previous page)

```
from numpy import exp
from scipy.special import factorial
import pandas as pd
from mpl_toolkits.mplot3d import Axes3D
import statsmodels.api as sm
from statsmodels.api import Poisson
from scipy import stats
from scipy.stats import norm
from statsmodels.iolib.summary2 import summary_col
```

64.1.1 Prerequisites

We assume familiarity with basic probability and multivariate calculus.

64.2 Set Up and Assumptions

Let's consider the steps we need to go through in maximum likelihood estimation and how they pertain to this study.

64.2.1 Flow of Ideas

The first step with maximum likelihood estimation is to choose the probability distribution believed to be generating the data.

More precisely, we need to make an assumption as to which *parametric class* of distributions is generating the data.

- e.g., the class of all normal distributions, or the class of all gamma distributions.

Each such class is a family of distributions indexed by a finite number of parameters.

- e.g., the class of normal distributions is a family of distributions indexed by its mean $\mu \in (-\infty, \infty)$ and standard deviation $\sigma \in (0, \infty)$.

We'll let the data pick out a particular element of the class by pinning down the parameters.

The parameter estimates so produced will be called **maximum likelihood estimates**.

64.2.2 Counting Billionaires

Treisman [Tre16] is interested in estimating the number of billionaires in different countries.

The number of billionaires is integer-valued.

Hence we consider distributions that take values only in the nonnegative integers.

(This is one reason least squares regression is not the best tool for the present problem, since the dependent variable in linear regression is not restricted to integer values)

One integer distribution is the [Poisson distribution](#), the probability mass function (pmf) of which is

$$f(y) = \frac{\mu^y}{y!} e^{-\mu}, \quad y = 0, 1, 2, \dots, \infty$$

We can plot the Poisson distribution over y for different values of μ as follows

```

poisson_pmf = lambda y, mu: mu**y / factorial(y) * exp(-mu)
y_values = range(0, 25)

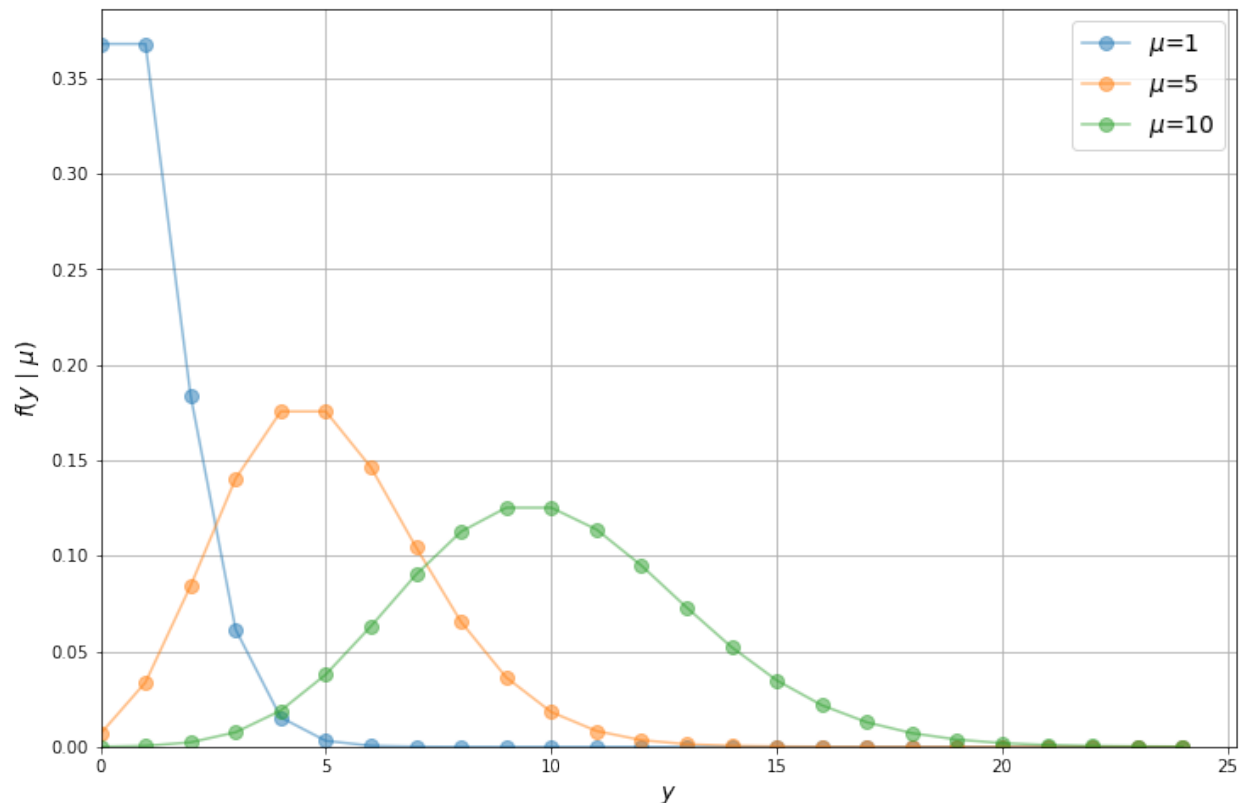
fig, ax = plt.subplots(figsize=(12, 8))

for mu in [1, 5, 10]:
    distribution = []
    for y_i in y_values:
        distribution.append(poisson_pmf(y_i, mu))
    ax.plot(y_values,
            distribution,
            label=f'$\mu$={mu}',
            alpha=0.5,
            marker='o',
            markersize=8)

ax.grid()
ax.set_xlabel('$y$', fontsize=14)
ax.set_ylabel('$f(y | \mu)$', fontsize=14)
ax.axis(xmin=0, ymin=0)
ax.legend(fontsize=14)

plt.show()

```



Notice that the Poisson distribution begins to resemble a normal distribution as the mean of y increases.

Let's have a look at the distribution of the data we'll be working with in this lecture.

Treisman's main source of data is *Forbes'* annual rankings of billionaires and their estimated net worth.

The dataset `mle/fp.dta` can be downloaded from [here](#) or its [AER page](#).

```
pd.options.display.max_columns = 10

# Load in data and view
df = pd.read_stata('https://github.com/QuantEcon/lecture-python/blob/master/source/_
↳static/lecture_specific/mle/fp.dta?raw=true')
df.head()
```

```

      country  ccode  year  cyear  numbil  ...  topint08  rintr  \
0  United States    2.0  1990.0  21990.0    NaN  ...  39.799999  4.988405
1  United States    2.0  1991.0  21991.0    NaN  ...  39.799999  4.988405
2  United States    2.0  1992.0  21992.0    NaN  ...  39.799999  4.988405
3  United States    2.0  1993.0  21993.0    NaN  ...  39.799999  4.988405
4  United States    2.0  1994.0  21994.0    NaN  ...  39.799999  4.988405

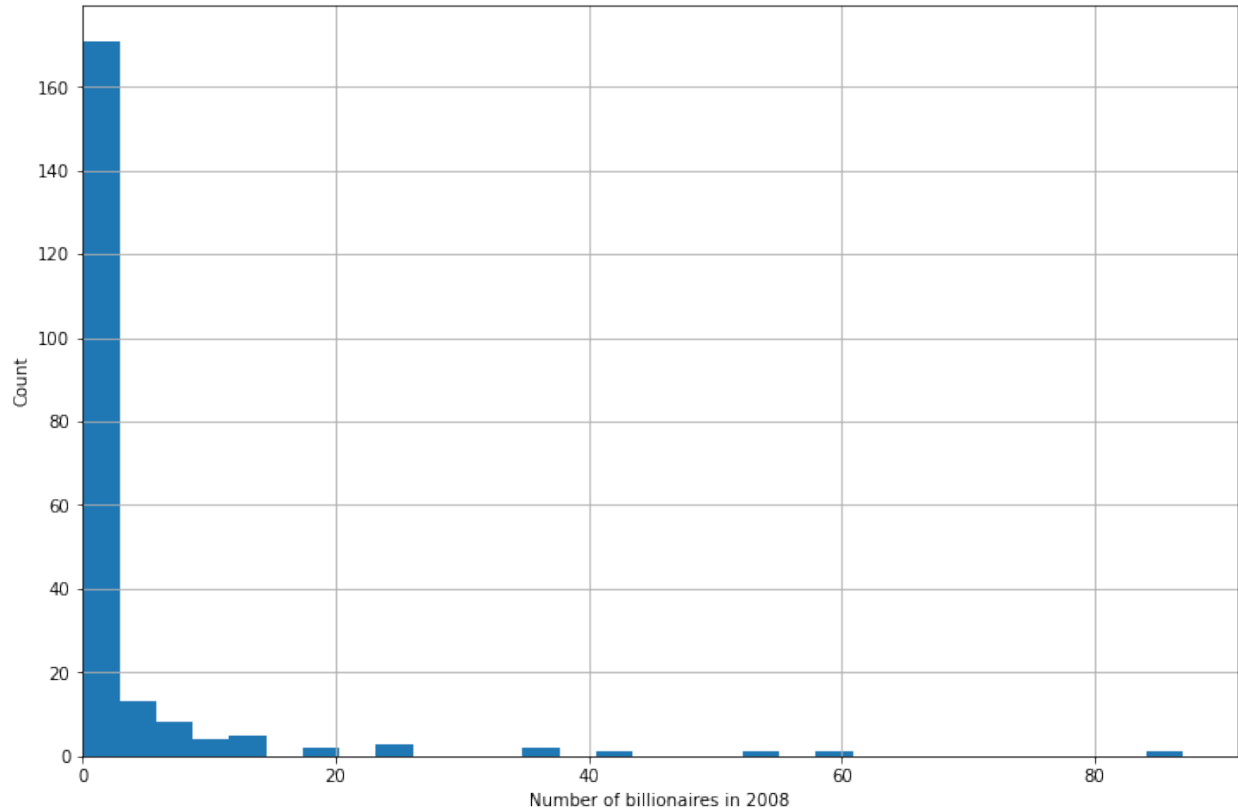
      noyrs  roflaw  nrrents
0    20.0    1.61    NaN
1    20.0    1.61    NaN
2    20.0    1.61    NaN
3    20.0    1.61    NaN
4    20.0    1.61    NaN

[5 rows x 36 columns]
```

Using a histogram, we can view the distribution of the number of billionaires per country, `numbil0`, in 2008 (the United States is dropped for plotting purposes)

```
numbil0_2008 = df[(df['year'] == 2008) & (
    df['country'] != 'United States')].loc[:, 'numbil0']

plt.subplots(figsize=(12, 8))
plt.hist(numbil0_2008, bins=30)
plt.xlim(left=0)
plt.grid()
plt.xlabel('Number of billionaires in 2008')
plt.ylabel('Count')
plt.show()
```



From the histogram, it appears that the Poisson assumption is not unreasonable (albeit with a very low μ and some outliers).

64.3 Conditional Distributions

In Treisman's paper, the dependent variable — the number of billionaires y_i in country i — is modeled as a function of GDP per capita, population size, and years membership in GATT and WTO.

Hence, the distribution of y_i needs to be conditioned on the vector of explanatory variables \mathbf{x}_i .

The standard formulation — the so-called *poisson regression* model — is as follows:

$$f(y_i | \mathbf{x}_i) = \frac{\mu_i^{y_i}}{y_i!} e^{-\mu_i}; \quad y_i = 0, 1, 2, \dots, \infty. \quad (1)$$

$$\text{where } \mu_i = \exp(\mathbf{x}_i' \boldsymbol{\beta}) = \exp(\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik})$$

To illustrate the idea that the distribution of y_i depends on \mathbf{x}_i let's run a simple simulation.

We use our `poisson_pmf` function from above and arbitrary values for $\boldsymbol{\beta}$ and \mathbf{x}_i

```
y_values = range(0, 20)

# Define a parameter vector with estimates
beta = np.array([0.26, 0.18, 0.25, -0.1, -0.22])

# Create some observations X
datasets = [np.array([0, 1, 1, 1, 2]),
```

(continues on next page)

(continued from previous page)

```

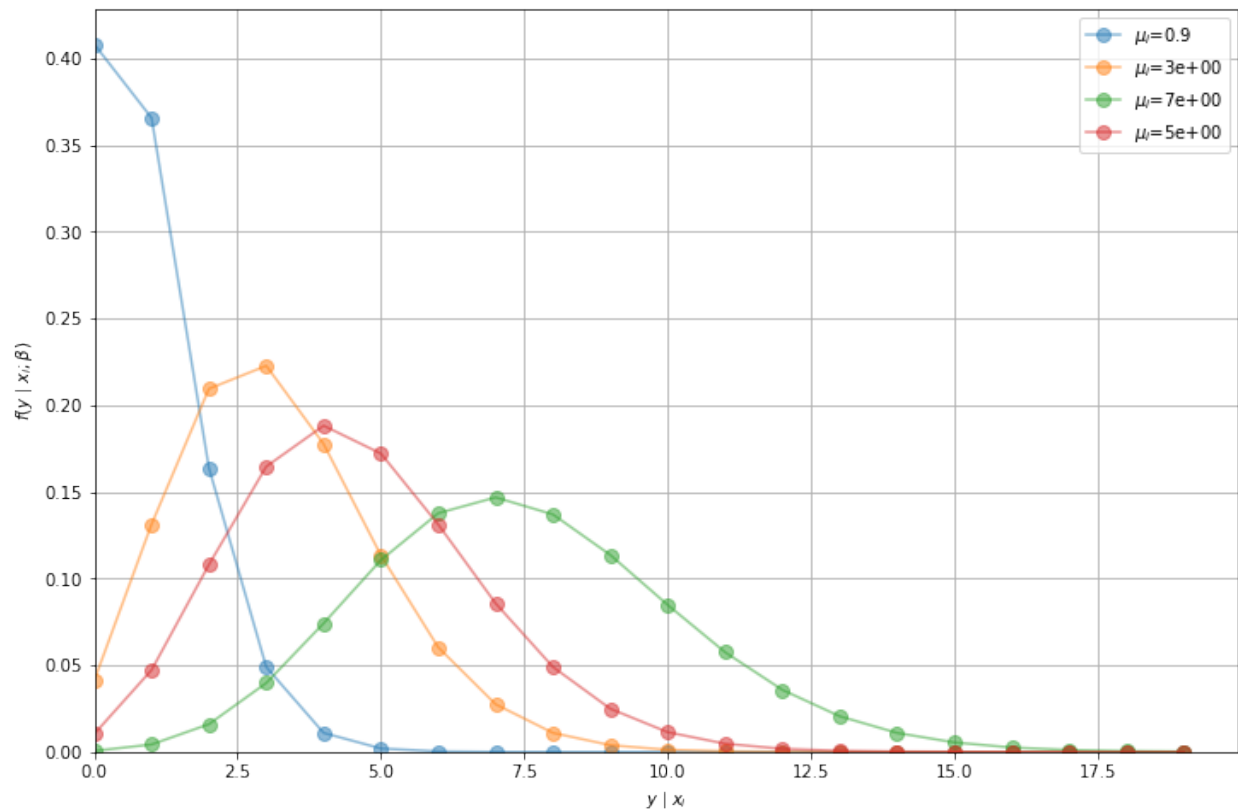
np.array([2, 3, 2, 4, 0]),
np.array([3, 4, 5, 3, 2]),
np.array([6, 5, 4, 4, 7])]

fig, ax = plt.subplots(figsize=(12, 8))

for X in datasets:
    μ = exp(X @ β)
    distribution = []
    for y_i in y_values:
        distribution.append(poisson_pmf(y_i, μ))
    ax.plot(y_values,
            distribution,
            label=f'$\mu_i$={μ:.1}',
            marker='o',
            markersize=8,
            alpha=0.5)

ax.grid()
ax.legend()
ax.set_xlabel('$y \mid x_i$')
ax.set_ylabel(r'$f(y \mid x_i; \beta)$')
ax.axis(xmin=0, ymin=0)
plt.show()

```



We can see that the distribution of y_i is conditional on \mathbf{x}_i (μ_i is no longer constant).

64.4 Maximum Likelihood Estimation

In our model for number of billionaires, the conditional distribution contains 4 ($k = 4$) parameters that we need to estimate.

We will label our entire parameter vector as β where

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$$

To estimate the model using MLE, we want to maximize the likelihood that our estimate $\hat{\beta}$ is the true parameter β .

Intuitively, we want to find the $\hat{\beta}$ that best fits our data.

First, we need to construct the likelihood function $\mathcal{L}(\beta)$, which is similar to a joint probability density function.

Assume we have some data $y_i = \{y_1, y_2\}$ and $y_i \sim f(y_i)$.

If y_1 and y_2 are independent, the joint pmf of these data is $f(y_1, y_2) = f(y_1) \cdot f(y_2)$.

If y_i follows a Poisson distribution with $\lambda = 7$, we can visualize the joint pmf like so

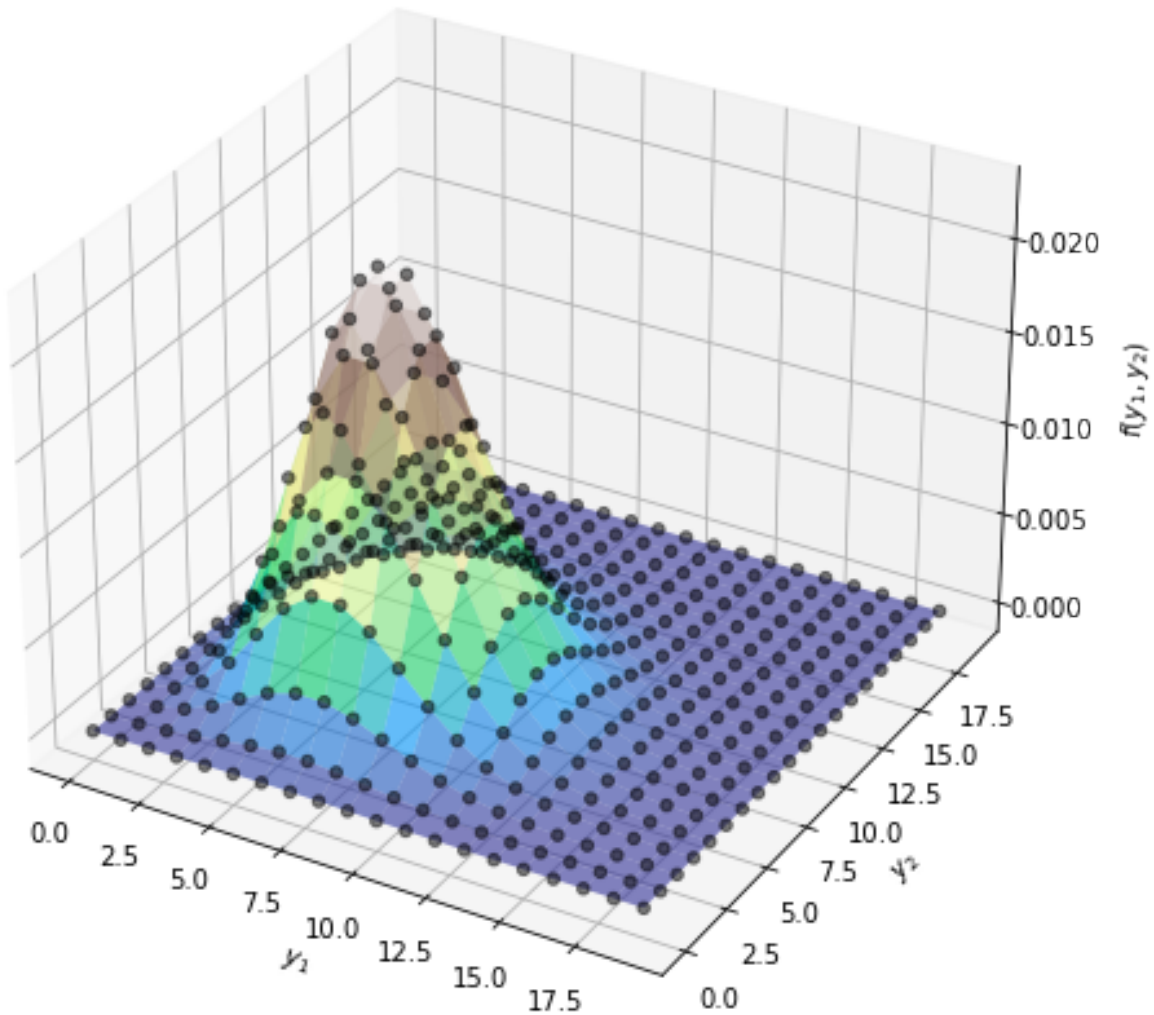
```
def plot_joint_poisson(mu=7, y_n=20):
    yi_values = np.arange(0, y_n, 1)

    # Create coordinate points of X and Y
    X, Y = np.meshgrid(yi_values, yi_values)

    # Multiply distributions together
    Z = poisson_pmf(X, mu) * poisson_pmf(Y, mu)

    fig = plt.figure(figsize=(12, 8))
    ax = fig.add_subplot(111, projection='3d')
    ax.plot_surface(X, Y, Z.T, cmap='terrain', alpha=0.6)
    ax.scatter(X, Y, Z.T, color='black', alpha=0.5, linewidths=1)
    ax.set_xlabel('$y_1$', ylabel='$y_2$')
    ax.set_zlabel('$f(y_1, y_2)$', labelpad=10)
    plt.show()

plot_joint_poisson(mu=7, y_n=20)
```

Similarly, the joint pmf of our data (which is distributed as a conditional Poisson distribution) can be written as

$$f(y_1, y_2, \dots, y_n \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n; \beta) = \prod_{i=1}^n \frac{\mu_i^{y_i}}{y_i!} e^{-\mu_i}$$

y_i is conditional on both the values of \mathbf{x}_i and the parameters β .

The likelihood function is the same as the joint pmf, but treats the parameter β as a random variable and takes the observations (y_i, \mathbf{x}_i) as given

$$\begin{aligned} \mathcal{L}(\beta \mid y_1, y_2, \dots, y_n; \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) &= \prod_{i=1}^n \frac{\mu_i^{y_i}}{y_i!} e^{-\mu_i} \\ &= f(y_1, y_2, \dots, y_n \mid \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n; \beta) \end{aligned}$$

Now that we have our likelihood function, we want to find the $\hat{\beta}$ that yields the maximum likelihood value

$$\max_{\beta} \mathcal{L}(\beta)$$

In doing so it is generally easier to maximize the log-likelihood (consider differentiating $f(x) = x \exp(x)$ vs. $f(x) = \log(x) + x$).

Given that taking a logarithm is a monotone increasing transformation, a maximizer of the likelihood function will also be a maximizer of the log-likelihood function.

In our case the log-likelihood is

$$\begin{aligned}\log \mathcal{L}(\beta) &= \log \left(f(y_1; \beta) \cdot f(y_2; \beta) \cdot \dots \cdot f(y_n; \beta) \right) \\ &= \sum_{i=1}^n \log f(y_i; \beta) \\ &= \sum_{i=1}^n \log \left(\frac{\mu_i^{y_i}}{y_i!} e^{-\mu_i} \right) \\ &= \sum_{i=1}^n y_i \log \mu_i - \sum_{i=1}^n \mu_i - \sum_{i=1}^n \log y_i!\end{aligned}$$

The MLE of the Poisson to the Poisson for $\hat{\beta}$ can be obtained by solving

$$\max_{\beta} \left(\sum_{i=1}^n y_i \log \mu_i - \sum_{i=1}^n \mu_i - \sum_{i=1}^n \log y_i! \right)$$

However, no analytical solution exists to the above problem – to find the MLE we need to use numerical methods.

64.5 MLE with Numerical Methods

Many distributions do not have nice, analytical solutions and therefore require numerical methods to solve for parameter estimates.

One such numerical method is the Newton-Raphson algorithm.

Our goal is to find the maximum likelihood estimate $\hat{\beta}$.

At $\hat{\beta}$, the first derivative of the log-likelihood function will be equal to 0.

Let's illustrate this by supposing

$$\log \mathcal{L}(\beta) = -(\beta - 10)^2 - 10$$

```
β = np.linspace(1, 20)
logL = -(β - 10) ** 2 - 10
dlogL = -2 * β + 20

fig, (ax1, ax2) = plt.subplots(2, sharex=True, figsize=(12, 8))

ax1.plot(β, logL, lw=2)
ax2.plot(β, dlogL, lw=2)

ax1.set_ylabel(r'$\log \mathcal{L}(\beta)$',
               rotation=0,
               labelpad=35,
               fontsize=15)
ax2.set_ylabel(r'$\frac{d \log \mathcal{L}(\beta)}{d \beta}$',
               rotation=0,
```

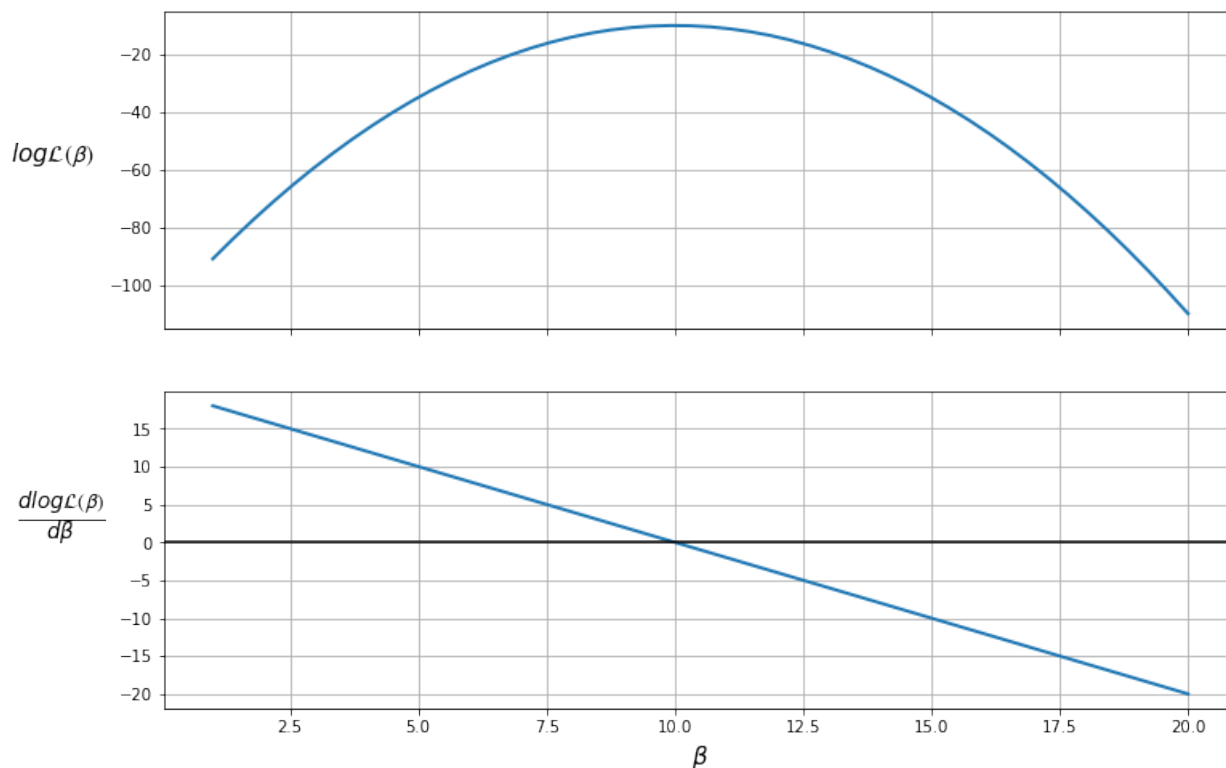
(continues on next page)

(continued from previous page)

```

        labelpad=35,
        fontsize=19)
ax2.set_xlabel(r'$\beta$', fontsize=15)
ax1.grid(), ax2.grid()
plt.axhline(c='black')
plt.show()

```



The plot shows that the maximum likelihood value (the top plot) occurs when $\frac{d \log \mathcal{L}(\beta)}{d \beta} = 0$ (the bottom plot).

Therefore, the likelihood is maximized when $\beta = 10$.

We can also ensure that this value is a *maximum* (as opposed to a minimum) by checking that the second derivative (slope of the bottom plot) is negative.

The Newton-Raphson algorithm finds a point where the first derivative is 0.

To use the algorithm, we take an initial guess at the maximum value, β_0 (the OLS parameter estimates might be a reasonable guess), then

1. Use the updating rule to iterate the algorithm

$$\beta_{(k+1)} = \beta_{(k)} - H^{-1}(\beta_{(k)})G(\beta_{(k)})$$

where:

$$G(\beta_{(k)}) = \frac{d \log \mathcal{L}(\beta_{(k)})}{d \beta_{(k)}}$$

$$H(\beta_{(k)}) = \frac{d^2 \log \mathcal{L}(\beta_{(k)})}{d \beta_{(k)} d \beta'_{(k)}}$$

2. Check whether $\beta_{(k+1)} - \beta_{(k)} < tol$

- If true, then stop iterating and set $\hat{\beta} = \beta_{(k+1)}$
- If false, then update $\beta_{(k+1)}$

As can be seen from the updating equation, $\beta_{(k+1)} = \beta_{(k)}$ only when $G(\beta_{(k)}) = 0$ ie. where the first derivative is equal to 0.

(In practice, we stop iterating when the difference is below a small tolerance threshold)

Let's have a go at implementing the Newton-Raphson algorithm.

First, we'll create a class called `PoissonRegression` so we can easily recompute the values of the log likelihood, gradient and Hessian for every iteration

```
class PoissonRegression:

    def __init__(self, y, X, β):
        self.X = X
        self.n, self.k = X.shape
        # Reshape y as a n_by_1 column vector
        self.y = y.reshape(self.n,1)
        # Reshape β as a k_by_1 column vector
        self.β = β.reshape(self.k,1)

    def μ(self):
        return np.exp(self.X @ self.β)

    def logL(self):
        y = self.y
        μ = self.μ()
        return np.sum(y * np.log(μ) - μ - np.log(factorial(y)))

    def G(self):
        y = self.y
        μ = self.μ()
        return X.T @ (y - μ)

    def H(self):
        X = self.X
        μ = self.μ()
        return -(X.T @ (μ * X))
```

Our function `newton_raphson` will take a `PoissonRegression` object that has an initial guess of the parameter vector β_0 .

The algorithm will update the parameter vector according to the updating rule, and recalculate the gradient and Hessian matrices at the new parameter estimates.

Iteration will end when either:

- The difference between the parameter and the updated parameter is below a tolerance level.
- The maximum number of iterations has been achieved (meaning convergence is not achieved).

So we can get an idea of what's going on while the algorithm is running, an option `display=True` is added to print out values at each iteration.

```
def newton_raphson(model, tol=1e-3, max_iter=1000, display=True):

    i = 0
    error = 100 # Initial error value
```

(continues on next page)

(continued from previous page)

```

# Print header of output
if display:
    header = f'{"Iteration_k":<13}{ "Log-likelihood":<16}{ "θ":<60}'
    print(header)
    print("-" * len(header))

# While loop runs while any value in error is greater
# than the tolerance until max iterations are reached
while np.any(error > tol) and i < max_iter:
    H, G = model.H(), model.G()
    β_new = model.β - (np.linalg.inv(H) @ G)
    error = β_new - model.β
    model.β = β_new

    # Print iterations
    if display:
        β_list = [f'{t:.3}' for t in list(model.β.flatten())]
        update = f'{i:<13}{model.logL():<16.8}{β_list}'
        print(update)

    i += 1

print(f'Number of iterations: {i}')
print(f'β_hat = {model.β.flatten()}')

# Return a flat array for β (instead of a k_by_1 column vector)
return model.β.flatten()

```

Let's try out our algorithm with a small dataset of 5 observations and 3 variables in **X**.

```

X = np.array([[1, 2, 5],
              [1, 1, 3],
              [1, 4, 2],
              [1, 5, 2],
              [1, 3, 1]])

y = np.array([1, 0, 1, 1, 0])

# Take a guess at initial βs
init_β = np.array([0.1, 0.1, 0.1])

# Create an object with Poisson model values
poi = PoissonRegression(y, X, β=init_β)

# Use newton_raphson to find the MLE
β_hat = newton_raphson(poi, display=True)

```

Iteration_k	Log-likelihood	θ
0	-4.3447622	['-1.49', '0.265', '0.244']
1	-3.5742413	['-3.38', '0.528', '0.474']
2	-3.3999526	['-5.06', '0.782', '0.702']
3	-3.3788646	['-5.92', '0.909', '0.82']
4	-3.3783559	['-6.07', '0.933', '0.843']

(continues on next page)

(continued from previous page)

```

5          -3.3783555          ['-6.08', '0.933', '0.843']
Number of iterations: 6
β_hat = [-6.07848205  0.93340226  0.84329625]

```

As this was a simple model with few observations, the algorithm achieved convergence in only 6 iterations.

You can see that with each iteration, the log-likelihood value increased.

Remember, our objective was to maximize the log-likelihood function, which the algorithm has worked to achieve.

Also, note that the increase in $\log \mathcal{L}(\beta_{(k)})$ becomes smaller with each iteration.

This is because the gradient is approaching 0 as we reach the maximum, and therefore the numerator in our updating equation is becoming smaller.

The gradient vector should be close to 0 at $\hat{\beta}$

```
poi.G()
```

```

array([[ -3.95169226e-07],
       [-1.00114804e-06],
       [-7.73114559e-07]])

```

The iterative process can be visualized in the following diagram, where the maximum is found at $\beta = 10$

```

logL = lambda x: -(x - 10) ** 2 - 10

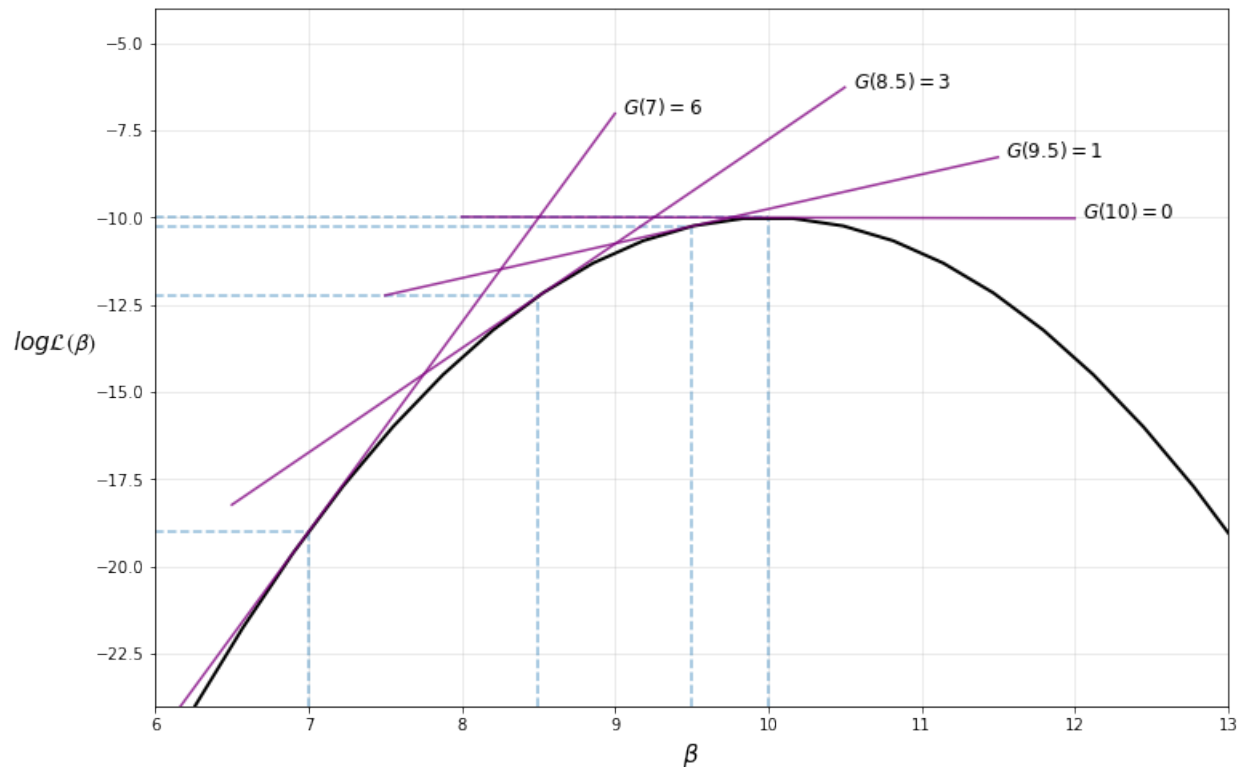
def find_tangent(β, a=0.01):
    y1 = logL(β)
    y2 = logL(β+a)
    x = np.array([[β, 1], [β+a, 1]])
    m, c = np.linalg.lstsq(x, np.array([y1, y2]), rcond=None)[0]
    return m, c

β = np.linspace(2, 18)
fig, ax = plt.subplots(figsize=(12, 8))
ax.plot(β, logL(β), lw=2, c='black')

for β in [7, 8.5, 9.5, 10]:
    β_line = np.linspace(β-2, β+2)
    m, c = find_tangent(β)
    y = m * β_line + c
    ax.plot(β_line, y, '-', c='purple', alpha=0.8)
    ax.text(β+2.05, y[-1], f'$G(\beta) = {abs(m):.0f}$', fontsize=12)
    ax.vlines(β, -24, logL(β), linestyles='--', alpha=0.5)
    ax.hlines(logL(β), 6, β, linestyles='--', alpha=0.5)

ax.set(ylim=(-24, -4), xlim=(6, 13))
ax.set_xlabel(r'$\beta$', fontsize=15)
ax.set_ylabel(r'$\log \mathcal{L}(\beta)$',
              rotation=0,
              labelpad=25,
              fontsize=15)
ax.grid(alpha=0.3)
plt.show()

```



Note that our implementation of the Newton-Raphson algorithm is rather basic — for more robust implementations see, for example, `scipy.optimize`.

64.6 Maximum Likelihood Estimation with `statsmodels`

Now that we know what's going on under the hood, we can apply MLE to an interesting application.

We'll use the Poisson regression model in `statsmodels` to obtain a richer output with standard errors, test values, and more.

`statsmodels` uses the same algorithm as above to find the maximum likelihood estimates.

Before we begin, let's re-estimate our simple model with `statsmodels` to confirm we obtain the same coefficients and log-likelihood value.

```
X = np.array([[1, 2, 5],
              [1, 1, 3],
              [1, 4, 2],
              [1, 5, 2],
              [1, 3, 1]])

y = np.array([1, 0, 1, 1, 0])

stats_poisson = Poisson(y, X).fit()
print(stats_poisson.summary())
```

```
Optimization terminated successfully.
Current function value: 0.675671
Iterations 7
```

(continues on next page)

(continued from previous page)

Poisson Regression Results						
=====						
Dep. Variable:	y	No. Observations:	5			
Model:	Poisson	Df Residuals:	2			
Method:	MLE	Df Model:	2			
Date:	Thu, 07 Oct 2021	Pseudo R-squ.:	0.2546			
Time:	21:25:01	Log-Likelihood:	-3.3784			
converged:	True	LL-Null:	-4.5325			
Covariance Type:	nonrobust	LLR p-value:	0.3153			
=====						
	coef	std err	z	P> z	[0.025	0.975]

const	-6.0785	5.279	-1.151	0.250	-16.425	4.268
x1	0.9334	0.829	1.126	0.260	-0.691	2.558
x2	0.8433	0.798	1.057	0.291	-0.720	2.407
=====						

Now let's replicate results from Daniel Treisman's paper, [Russia's Billionaires](#), mentioned earlier in the lecture.

Treisman starts by estimating equation (1), where:

- y_i is number of billionaires_{*i*}
- x_{i1} is log GDP per capita_{*i*}
- x_{i2} is log population_{*i*}
- x_{i3} is years in GATT_{*i*} – years membership in GATT and WTO (to proxy access to international markets)

The paper only considers the year 2008 for estimation.

We will set up our variables for estimation like so (you should have the data assigned to `df` from earlier in the lecture)

```
# Keep only year 2008
df = df[df['year'] == 2008]

# Add a constant
df['const'] = 1

# Variable sets
reg1 = ['const', 'lngdppc', 'lnpop', 'gattwto08']
reg2 = ['const', 'lngdppc', 'lnpop',
        'gattwto08', 'lnmcap08', 'rintr', 'topint08']
reg3 = ['const', 'lngdppc', 'lnpop', 'gattwto08', 'lnmcap08',
        'rintr', 'topint08', 'nrrents', 'roflaw']
```

Then we can use the Poisson function from `statsmodels` to fit the model.

We'll use robust standard errors as in the author's paper

```
# Specify model
poisson_reg = sm.Poisson(df[['numbil0']], df[reg1],
                        missing='drop').fit(cov_type='HC0')
print(poisson_reg.summary())
```

```
Optimization terminated successfully.
      Current function value: 2.226090
      Iterations 9

Poisson Regression Results
```

(continues on next page)

(continued from previous page)

```

=====
Dep. Variable:          numbil0      No. Observations:          197
Model:                  Poisson      Df Residuals:              193
Method:                  MLE          Df Model:                  3
Date:                   Thu, 07 Oct 2021  Pseudo R-squ.:          0.8574
Time:                   21:25:01      Log-Likelihood:          -438.54
converged:              True          LL-Null:                  -3074.7
Covariance Type:        HC0          LLR p-value:              0.000
=====

```

	coef	std err	z	P> z	[0.025	0.975]
const	-29.0495	2.578	-11.268	0.000	-34.103	-23.997
lngdppc	1.0839	0.138	7.834	0.000	0.813	1.355
lnpop	1.1714	0.097	12.024	0.000	0.980	1.362
gattwto08	0.0060	0.007	0.868	0.386	-0.008	0.019

```

=====

```

Success! The algorithm was able to achieve convergence in 9 iterations.

Our output indicates that GDP per capita, population, and years of membership in the General Agreement on Tariffs and Trade (GATT) are positively related to the number of billionaires a country has, as expected.

Let's also estimate the author's more full-featured models and display them in a single table

```

regs = [reg1, reg2, reg3]
reg_names = ['Model 1', 'Model 2', 'Model 3']
info_dict = {'Pseudo R-squared': lambda x: f"{x.prsquared:.2f}",
             'No. observations': lambda x: f"{int(x.nobs):d}"}
regressor_order = ['const',
                   'lngdppc',
                   'lnpop',
                   'gattwto08',
                   'lnmcap08',
                   'rintr',
                   'topint08',
                   'nrrents',
                   'roflaw']

results = []

for reg in regs:
    result = sm.Poisson(df[['numbil0']], df[reg],
                       missing='drop').fit(cov_type='HC0',
                                           maxiter=100, disp=0)

    results.append(result)

results_table = summary_col(results=results,
                             float_format='%0.3f',
                             stars=True,
                             model_names=reg_names,
                             info_dict=info_dict,
                             regressor_order=regressor_order)
results_table.add_title('Table 1 - Explaining the Number of Billionaires \
                        in 2008')
print(results_table)

```

Table 1 - Explaining the Number of Billionaires	in 2008
=====	

(continues on next page)

(continued from previous page)

	Model 1	Model 2	Model 3
-----	-----	-----	-----
const	-29.050*** (2.578)	-19.444*** (4.820)	-20.858*** (4.255)
lngdppc	1.084*** (0.138)	0.717*** (0.244)	0.737*** (0.233)
lnpop	1.171*** (0.097)	0.806*** (0.213)	0.929*** (0.195)
gattwto08	0.006 (0.007)	0.007 (0.006)	0.004 (0.006)
lnmcap08		0.399** (0.172)	0.286* (0.167)
rintr		-0.010 (0.010)	-0.009 (0.010)
topint08		-0.051*** (0.011)	-0.058*** (0.012)
nrrents			-0.005 (0.010)
roflaw			0.203 (0.372)
Pseudo R-squared	0.86	0.90	0.90
No. observations	197	131	131
=====	=====	=====	=====
Standard errors in parentheses.			
* p<.1, ** p<.05, ***p<.01			

The output suggests that the frequency of billionaires is positively correlated with GDP per capita, population size, stock market capitalization, and negatively correlated with top marginal income tax rate.

To analyze our results by country, we can plot the difference between the predicted and actual values, then sort from highest to lowest and plot the first 15

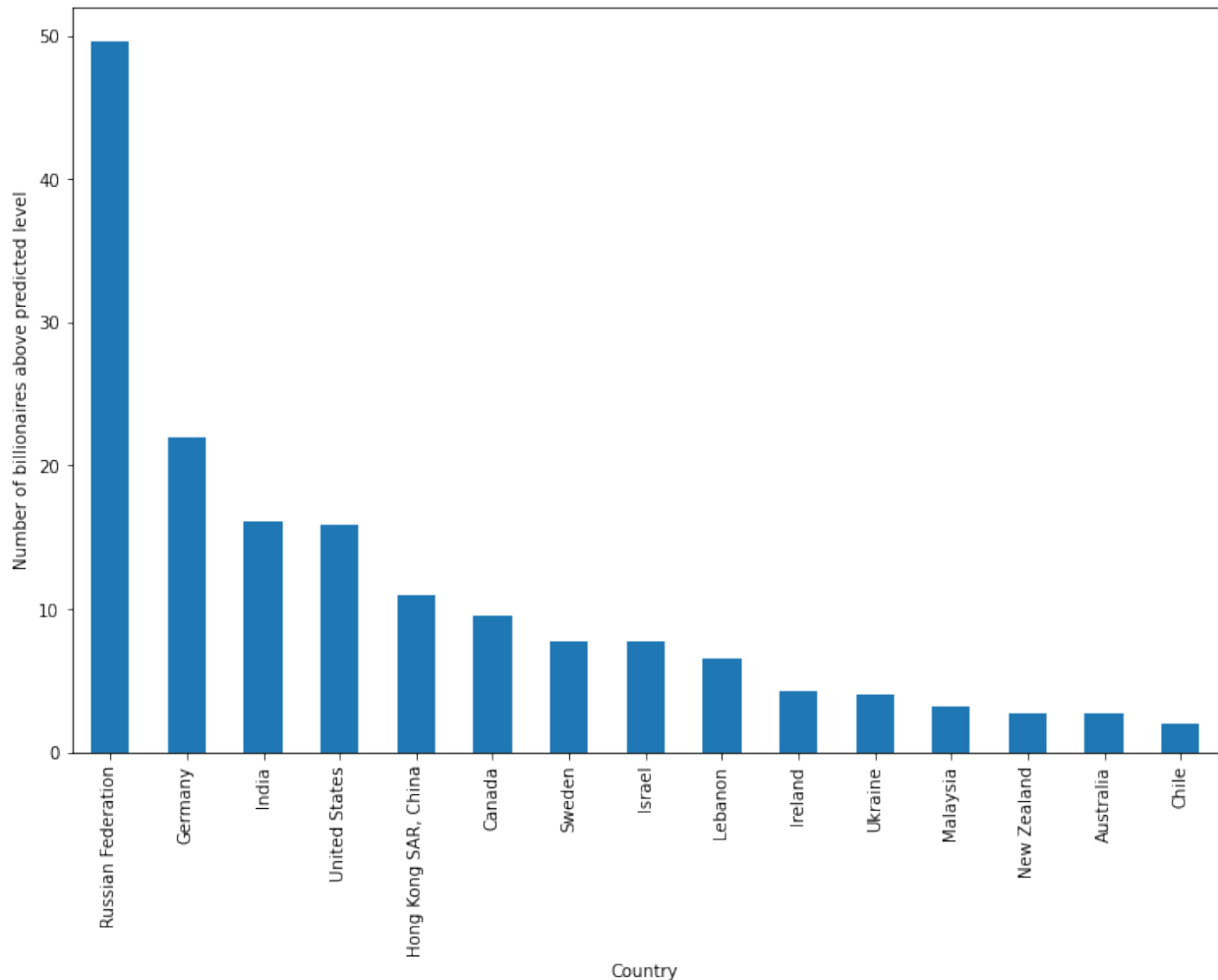
```
data = ['const', 'lngdppc', 'lnpop', 'gattwto08', 'lnmcap08', 'rintr',
        'topint08', 'nrrents', 'roflaw', 'numbil0', 'country']
results_df = df[data].dropna()

# Use last model (model 3)
results_df['prediction'] = results[-1].predict()

# Calculate difference
results_df['difference'] = results_df['numbil0'] - results_df['prediction']

# Sort in descending order
results_df.sort_values('difference', ascending=False, inplace=True)

# Plot the first 15 data points
results_df[:15].plot('country', 'difference', kind='bar',
                    figsize=(12,8), legend=False)
plt.ylabel('Number of billionaires above predicted level')
plt.xlabel('Country')
plt.show()
```



As we can see, Russia has by far the highest number of billionaires in excess of what is predicted by the model (around 50 more than expected).

Treisman uses this empirical result to discuss possible reasons for Russia's excess of billionaires, including the origination of wealth in Russia, the political climate, and the history of privatization in the years after the USSR.

64.7 Summary

In this lecture, we used Maximum Likelihood Estimation to estimate the parameters of a Poisson model.

`statsmodels` contains other built-in likelihood models such as [Probit](#) and [Logit](#).

For further flexibility, `statsmodels` provides a way to specify the distribution manually using the `GenericLikelihoodModel` class - an example notebook can be found [here](#).

64.8 Exercises

64.8.1 Exercise 1

Suppose we wanted to estimate the probability of an event y_i occurring, given some observations.

We could use a probit regression model, where the pmf of y_i is

$$f(y_i; \beta) = \mu_i^{y_i} (1 - \mu_i)^{1-y_i}, \quad y_i = 0, 1$$

where $\mu_i = \Phi(\mathbf{x}_i' \beta)$

Φ represents the *cumulative normal distribution* and constrains the predicted y_i to be between 0 and 1 (as required for a probability).

β is a vector of coefficients.

Following the example in the lecture, write a class to represent the Probit model.

To begin, find the log-likelihood function and derive the gradient and Hessian.

The `scipy` module `stats.norm` contains the functions needed to compute the cmf and pmf of the normal distribution.

64.8.2 Exercise 2

Use the following dataset and initial values of β to estimate the MLE with the Newton-Raphson algorithm developed earlier in the lecture

$$\mathbf{X} = \begin{bmatrix} 1 & 2 & 4 \\ 1 & 1 & 1 \\ 1 & 4 & 3 \\ 1 & 5 & 6 \\ 1 & 3 & 5 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \quad \beta_{(0)} = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}$$

Verify your results with `statsmodels` - you can import the Probit function with the following import statement

```
from statsmodels.discrete.discrete_model import Probit
```

Note that the simple Newton-Raphson algorithm developed in this lecture is very sensitive to initial values, and therefore you may fail to achieve convergence with different starting values.

64.9 Solutions

64.9.1 Exercise 1

The log-likelihood can be written as

$$\log \mathcal{L} = \sum_{i=1}^n [y_i \log \Phi(\mathbf{x}_i' \beta) + (1 - y_i) \log(1 - \Phi(\mathbf{x}_i' \beta))]$$

Using the **fundamental theorem of calculus**, the derivative of a cumulative probability distribution is its marginal distribution

$$\frac{\partial}{\partial s} \Phi(s) = \phi(s)$$

where ϕ is the marginal normal distribution.

The gradient vector of the Probit model is

$$\frac{\partial \log \mathcal{L}}{\partial \beta} = \sum_{i=1}^n \left[y_i \frac{\phi(\mathbf{x}'_i \beta)}{\Phi(\mathbf{x}'_i \beta)} - (1 - y_i) \frac{\phi(\mathbf{x}'_i \beta)}{1 - \Phi(\mathbf{x}'_i \beta)} \right] \mathbf{x}_i$$

The Hessian of the Probit model is

$$\frac{\partial^2 \log \mathcal{L}}{\partial \beta \partial \beta'} = - \sum_{i=1}^n \phi(\mathbf{x}'_i \beta) \left[y_i \frac{\phi(\mathbf{x}'_i \beta) + \mathbf{x}'_i \beta \Phi(\mathbf{x}'_i \beta)}{[\Phi(\mathbf{x}'_i \beta)]^2} + (1 - y_i) \frac{\phi(\mathbf{x}'_i \beta) - \mathbf{x}'_i \beta (1 - \Phi(\mathbf{x}'_i \beta))}{[1 - \Phi(\mathbf{x}'_i \beta)]^2} \right] \mathbf{x}_i \mathbf{x}'_i$$

Using these results, we can write a class for the Probit model as follows

```
class ProbitRegression:

    def __init__(self, y, X, β):
        self.X, self.y, self.β = X, y, β
        self.n, self.k = X.shape

    def μ(self):
        return norm.cdf(self.X @ self.β.T)

    def φ(self):
        return norm.pdf(self.X @ self.β.T)

    def logL(self):
        μ = self.μ()
        return np.sum(y * np.log(μ) + (1 - y) * np.log(1 - μ))

    def G(self):
        μ = self.μ()
        φ = self.φ()
        return np.sum((X.T * y * φ / μ - X.T * (1 - y) * φ / (1 - μ)),
                      axis=1)

    def H(self):
        X = self.X
        β = self.β
        μ = self.μ()
        φ = self.φ()
        a = (φ + (X @ β.T) * μ) / μ**2
        b = (φ - (X @ β.T) * (1 - μ)) / (1 - μ)**2
        return -(φ * (y * a + (1 - y) * b) * X.T) @ X
```

64.9.2 Exercise 2

```
X = np.array([[1, 2, 4],
              [1, 1, 1],
              [1, 4, 3],
              [1, 5, 6],
              [1, 3, 5]])

y = np.array([1, 0, 1, 1, 0])

# Take a guess at initial βs
```

(continues on next page)

(continued from previous page)

```

β = np.array([0.1, 0.1, 0.1])

# Create instance of Probit regression class
prob = ProbitRegression(y, X, β)

# Run Newton-Raphson algorithm
newton_raphson(prob)

```

```

Iteration_k  Log-likelihood  θ
-----
↪----
0           -2.3796884      ['-1.34', '0.775', '-0.157']
1           -2.3687526      ['-1.53', '0.775', '-0.0981']
2           -2.3687294      ['-1.55', '0.778', '-0.0971']
3           -2.3687294      ['-1.55', '0.778', '-0.0971']
Number of iterations: 4
β_hat = [-1.54625858  0.77778952 -0.09709757]

```

```
array([-1.54625858,  0.77778952, -0.09709757])
```

```

# Use statsmodels to verify results

print(Probit(y, X).fit().summary())

```

```

Optimization terminated successfully.
      Current function value: 0.473746
      Iterations 6

                    Probit Regression Results
=====
Dep. Variable:          y      No. Observations:          5
Model:                  Probit    Df Residuals:            2
Method:                  MLE      Df Model:              2
Date:                   Thu, 07 Oct 2021    Pseudo R-squ.:        0.2961
Time:                   21:25:01    Log-Likelihood:        -2.3687
converged:               True      LL-Null:              -3.3651
Covariance Type:         nonrobust    LLR p-value:           0.3692
=====
              coef      std err          z      P>|z|      [0.025      0.975]
-----
const         -1.5463      1.866      -0.829      0.407      -5.204      2.111
x1              0.7778      0.788       0.986      0.324      -0.768      2.323
x2            -0.0971      0.590      -0.165      0.869      -1.254      1.060
=====

```


Part XI

Other

TROUBLESHOOTING

Contents

- *Troubleshooting*
 - *Fixing Your Local Environment*
 - *Reporting an Issue*

This page is for readers experiencing errors when running the code from the lectures.

65.1 Fixing Your Local Environment

The basic assumption of the lectures is that code in a lecture should execute whenever

1. it is executed in a Jupyter notebook and
2. the notebook is running on a machine with the latest version of Anaconda Python.

You have installed Anaconda, haven't you, following the instructions in [this lecture](#)?

Assuming that you have, the most common source of problems for our readers is that their Anaconda distribution is not up to date.

[Here's a useful article](#) on how to update Anaconda.

Another option is to simply remove Anaconda and reinstall.

You also need to keep the external code libraries, such as [QuantEcon.py](#) up to date.

For this task you can either

- use `conda install -y quantecon` on the command line, or
- execute `!conda install -y quantecon` within a Jupyter notebook.

If your local environment is still not working you can do two things.

First, you can use a remote machine instead, by clicking on the Launch Notebook icon available for each lecture



Second, you can report an issue, so we can try to fix your local set up.

We like getting feedback on the lectures so please don't hesitate to get in touch.

65.2 Reporting an Issue

One way to give feedback is to raise an issue through our [issue tracker](#).

Please be as specific as possible. Tell us where the problem is and as much detail about your local set up as you can provide.

Another feedback option is to use our [discourse forum](#).

Finally, you can provide direct feedback to contact@quantecon.org