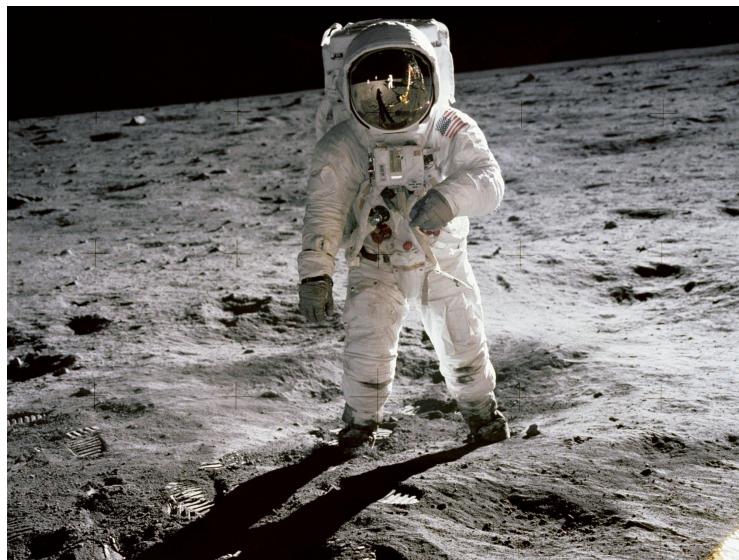


Labs for Foundations of Applied Mathematics

Volume 3
Modeling with Uncertainty and Data

Jeffrey Humpherys & Tyler J. Jarvis, managing editors



List of Contributors

B. Barker <i>Brigham Young University</i>	R. Dorff <i>Brigham Young University</i>
E. Evans <i>Brigham Young University</i>	B. Ehler <i>Brigham Young University</i>
R. Evans <i>Brigham Young University</i>	M. Fabiano <i>Brigham Young University</i>
J. Grout <i>Drake University</i>	K. Finlinson <i>Brigham Young University</i>
J. Humpherys <i>Brigham Young University</i>	J. Fisher <i>Brigham Young University</i>
T. Jarvis <i>Brigham Young University</i>	R. Flores <i>Brigham Young University</i>
J. Whitehead <i>Brigham Young University</i>	R. Fowers <i>Brigham Young University</i>
J. Adams <i>Brigham Young University</i>	A. Frandsen <i>Brigham Young University</i>
K. Baldwin <i>Brigham Young University</i>	R. Fuhriman <i>Brigham Young University</i>
J. Bejarano <i>Brigham Young University</i>	T. Gledhill <i>Brigham Young University</i>
A. Berry <i>Brigham Young University</i>	S. Giddens <i>Brigham Young University</i>
Z. Boyd <i>Brigham Young University</i>	C. Gigena <i>Brigham Young University</i>
M. Brown <i>Brigham Young University</i>	M. Graham <i>Brigham Young University</i>
A. Carr <i>Brigham Young University</i>	F. Glines <i>Brigham Young University</i>
C. Carter <i>Brigham Young University</i>	C. Glover <i>Brigham Young University</i>
T. Christensen <i>Brigham Young University</i>	M. Goodwin <i>Brigham Young University</i>
M. Cook <i>Brigham Young University</i>	R. Grout <i>Brigham Young University</i>

- D. Grundvig
Brigham Young University
- S. Halverson
Brigham Young University
- E. Hannesson
Brigham Young University
- K. Harmer
Brigham Young University
- J. Henderson
Brigham Young University
- J. Hendricks
Brigham Young University
- A. Henriksen
Brigham Young University
- I. Henriksen
Brigham Young University
- C. Hettinger
Brigham Young University
- S. Horst
Brigham Young University
- R. Howell
Brigham Young University
- E. Ibarra-Campos
Brigham Young University
- J. Larsen
Brigham Young University
- K. Jacobson
Brigham Young University
- R. Jenkins
Brigham Young University
- J. Leete
Brigham Young University
- Q. Leishman
Brigham Young University
- J. Lytle
Brigham Young University
- E. Manner
Brigham Young University
- M. Matsushita
Brigham Young University
- R. McMurray
Brigham Young University
- S. McQuarrie
Brigham Young University
- D. Miller
Brigham Young University
- J. Morrise
Brigham Young University
- M. Morrise
Brigham Young University
- A. Morrow
Brigham Young University
- R. Murray
Brigham Young University
- J. Nelson
Brigham Young University
- C. Noorda
Brigham Young University
- A. Oldroyd
Brigham Young University
- A. Oveson
Brigham Young University
- E. Parkinson
Brigham Young University
- M. Probst
Brigham Young University
- M. Proudfoot
Brigham Young University
- D. Reber
Brigham Young University
- H. Ringer
Brigham Young University
- C. Robertson
Brigham Young University
- M. Russell
Brigham Young University
- R. Sandberg
Brigham Young University
- C. Sawyer
Brigham Young University
- D. Smith
Brigham Young University
- J. Smith
Brigham Young University
- P. Smith
Brigham Young University
- M. Stauffer
Brigham Young University

E. Steadman

Brigham Young University

J. Stewart

Brigham Young University

S. Suggs

Brigham Young University

A. Tate

Brigham Young University

T. Thompson

Brigham Young University

M. Victors

Brigham Young University

E. Walker

Brigham Young University

J. Webb

Brigham Young University

R. Webb

Brigham Young University

J. West

Brigham Young University

R. Wonnacott

Brigham Young University

A. Zaitzeff

Brigham Young University

This project is funded in part by the National Science Foundation, grant no. TUES Phase II
DUE-1323785.

Preface

This lab manual is designed to accompany the textbook *Foundations of Applied Mathematics Volume 3: Modeling with Uncertainty and Data* by Humpherys and Jarvis. The labs present various aspects of important machine learning algorithms. The reader should be familiar with Python [VD10] and its NumPy [Oli06, ADH⁺01, Oli07] and Matplotlib [Hun07] packages before attempting these labs. See the Python Essentials manual for introductions to these topics.

©This work is licensed under the Creative Commons Attribution 3.0 United States License. You may copy, distribute, and display this copyrighted work only if you give credit to Dr. J. Humpherys. All derivative works must include an attribution to Dr. J. Humpherys as the owner of this work as well as the web address to

<https://github.com/Foundations-of-Applied-Mathematics/Labs>
as the original source of this work.

To view a copy of the Creative Commons Attribution 3.0 License, visit

<http://creativecommons.org/licenses/by/3.0/us/>
or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.



Contents

Preface	v
I Labs	1
1 Pandas 1: Introduction	3
2 Pandas 2: Plotting	27
3 Pandas 3: Grouping	41
4 Information Theory	53
5 GeoPandas	61
6 Data Cleaning	73
7 LSI and SkLearn	87
8 Intro to Parallel Computing	99
9 Linear Regression	113
10 Logistic Regression	119
11 Naive Bayes	127
12 Random Forests	135
13 Apache Spark	143
14 Web Crawling	161
15 Web Scraping	169
16 K-Means Clustering	181
17 Data Augmentation and Generation	193

18	Metropolis Algorithm	203
19	Gibbs Sampling and LDA	213
20	Gaussian Mixture Models	227
21	Discrete Hidden Markov Models	235
22	Speech Recognition using CDHMMs	247
23	Kalman Filter	255
24	ARMA Models	265
25	Non-negative Matrix Factorization Recommender	283
II	Appendices	289
A	Getting Started	291
B	Installing and Managing Python	299
C	NumPy Visual Guide	305
D	Introduction to Scikit-Learn	309
	Bibliography	325

Part I

Labs

1

Pandas 1: Introduction

Lab Objective: Though NumPy and SciPy are powerful tools for numerical computing, they lack some of the high-level functionality necessary for many data science applications. Python's pandas library, built on NumPy, is designed specifically for data management and analysis. In this lab we introduce pandas data structures, syntax, and explore its capabilities for quickly analyzing and presenting data.

Pandas Basics

Pandas is a python library used primarily to analyze data. It combines functionality of NumPy, MatPlotLib, and SQL to create an easy to understand library that allows for the manipulation of data in various ways. In this lab we focus on the use of Pandas to analyze and manipulate data in ways similar to NumPy and SQL.

Pandas Data Structures

Series

The first pandas data structure is a **Series**. A **Series** is a one-dimensional array that can hold any datatype, similar to a **ndarray**. However, a **Series** has an **index** that gives a label to each entry. An **index** generally is used to label the data.

Typically a **Series** contains information about one feature of the data. For example, the data in a **Series** might show a class's grades on a test and the **Index** would indicate each student in the class. To initialize a **Series**, the first parameter is the data and the second is the index.

```
>>> import pandas as pd  
>>>  
# Initialize Series of student grades  
>>> math = pd.Series(np.random.randint(0,100,4), ['Mark', 'Barbara',  
...     'Eleanor', 'David'])  
>>> english = pd.Series(np.random.randint(0,100,5), ['Mark', 'Barbara',  
...     'David', 'Greg', 'Lauren'])
```

DataFrame

The second key pandas data structure is a `DataFrame`. A `DataFrame` is a collection of multiple `Series`. It can be thought of as a 2-dimensional array, where each row is a separate datapoint and each column is a feature of the data. The rows are labeled with an `index` (as in a `Series`) and the columns are labeled in the attribute `columns`.

There are many different ways to initialize a `DataFrame`. One way to initialize a `DataFrame` is by passing in a dictionary as the data of the `DataFrame`. The keys of the dictionary will become the labels in `columns` and the values are the `Series` associated with the label.

```
# Create a DataFrame of student grades
>>> grades = pd.DataFrame({"Math": math, "English": english})
>>> grades
      Math   English
Barbara    52.0     73.0
David      10.0     39.0
Eleanor    35.0      NaN
Greg        NaN     26.0
Lauren     NaN     99.0
Mark       81.0     68.0
```

Notice that `pd.DataFrame` automatically lines up data from both `Series` that have the same index. If the data only appears in one of the `Series`, the corresponding entry for the other `Series` is `NaN`.

We can also initialize a `DataFrame` with a NumPy array. With this method, the data is passed in as a 2-dimensional NumPy array, while the column labels and the index are passed in as parameters. The first column label goes with the first column of the array, the second with the second, and so forth. The index works similarly.

```
>>> import numpy as np
# Initialize DataFrame with NumPy array. This is identical to the grades ←
# DataFrame above.
>>> data = np.array([[52.0, 73.0], [10.0, 39.0], [35.0, np.nan],
...      [np.nan, 26.0], [np.nan, 99.0], [81.0, 68.0]])
>>> grades = pd.DataFrame(data, columns = ['Math', 'English'], index =
...      ['Barbara', 'David', 'Eleanor', 'Greg', 'Lauren', 'Mark'])

# View the columns
>>> grades.columns
Index(['Math', 'English'], dtype='object')

# View the Index
>>> grades.index
Index(['Barbara', 'David', 'Eleanor', 'Greg', 'Lauren', 'Mark'], dtype='object')
```

A `DataFrame` can also be viewed as a NumPy array using the attribute `values`.

```
# View the DataFrame as a NumPy array
```

```
>>> grades.values
array([[ 52.,  73.],
       [ 10.,  39.],
       [ 35.,   nan],
       [  nan,  26.],
       [  nan,  99.],
       [ 81.,  68.]])
```

Data I/O

The pandas library has functions that make importing and exporting data simple. The functions allow for a variety of file formats to be imported and exported, including CSV, Excel, HDF5, SQL, JSON, HTML, and pickle files.

Method	Description
<code>to_csv()</code>	Write the index and entries to a CSV file
<code>read_csv()</code>	Read a csv and convert into a DataFrame
<code>to_json()</code>	Convert the object to a JSON string
<code>to_pickle()</code>	Serialize the object and store it in an external file
<code>to_sql()</code>	Write the object data to an open SQL database
<code>read_html()</code>	Read a table in an html page and convert to a DataFrame

Table 1.1: Methods for exporting data in a pandas `Series` or `DataFrame`.

The CSV (comma separated values) format is a simple way of storing tabular data in plain text. Because CSV files are one of the most popular file formats for exchanging data, we will explore the `read_csv()` function in more detail. Some frequently-used keyword arguments include the following:

- **delimiter**: The character that separates data fields. It is often a comma or a whitespace character.
- **header**: The row number (0 indexed) in the CSV file that contains the column names.
- **index_col**: The column (0 indexed) in the CSV file that is the index for the `DataFrame`.
- **skiprows**: If an integer n , skip the first n rows of the file, and then start reading in the data. If a list of integers, skip the specified rows.
- **names**: If the CSV file does not contain the column names, or you wish to use other column names, specify them in a list.

Another particularly useful function is `read_html()`, which is useful when scraping data. It takes in a url or html file and an optional argument `match`, a string or regex, and returns a list of the tables that match the `match` in a DataFrame. While the resulting data will probably need to be cleaned, it is frequently much faster than scraping a website.

Data Manipulation

Accessing Data

In general, the best way to access data in a `Series` or `DataFrame` is through the indexers `loc` and `iloc`. While array slicing can be used, it is more efficient to use these indexers. Accessing `Series` and `DataFrame` objects using these indexing operations is more efficient than slicing because the bracket indexing has to check many cases before it can determine how to slice the data structure. Using `loc` or `iloc` explicitly bypasses these extra checks. The `loc` index selects rows and columns based on their labels, while `iloc` selects them based on their integer position. With these indexers, the first and second arguments refer to the rows and columns, respectively, just as array slicing.

```
# Use loc to select the Math scores of David and Greg
>>> grades.loc[['David', 'Greg'], 'Math']
David    10.0
Greg      NaN
Name: Math, dtype: float64

# Use iloc to select the Math scores of David and Greg
>>> grades.iloc[[1,3], 0]
David    10.0
Greg      NaN
```

To access an entire column of a `DataFrame`, the most efficient method is to use only square brackets and the name of the column, without the indexer. This syntax can also be used to create a new column or reset the values of an entire column.

```
# Create a new History column with array of random values
>>> grades['History'] = np.random.randint(0,100,6)
>>> grades['History']
Barbara     4
David      92
Eleanor    25
Greg       79
Lauren     82
Mark       27
Name: History, dtype: int64

# Reset the column such that everyone has a 100
>>> grades['History'] = 100.0
>>> grades
   Math  English  History
Barbara  52.0     73.0    100.0
David    10.0     39.0    100.0
Eleanor  35.0     55.0    100.0
Greg     26.0     26.0    100.0
Lauren   99.0     68.0    100.0
Mark     81.0     68.0    100.0
```

Datasets can often be very large and thus difficult to visualize. Pandas has various methods to make this easier. The methods `head` and `tail` will show the first or last n data points, respectively, where n defaults to 5. The method `sample` will draw n random entries of the dataset, where n defaults to 1.

```
# Use head to see the first n rows
>>> grades.head(n=2)
      Math   English   History
Barbara    52.0      73.0     100.0
David      10.0      39.0     100.0

# Use sample to sample a random entry
>>> grades.sample()
      Math   English   History
Lauren     NaN      99.0     100.0
```

It may also be useful to re-order the columns or rows or sort according to a given column.

```
# Re-order columns
>>> grades.reindex(columns=['English', 'Math', 'History'])
      English   Math   History
Barbara    73.0    52.0     100.0
David      39.0    10.0     100.0
Eleanor     NaN    35.0     100.0
Greg        26.0    NaN      100.0
Lauren     99.0    NaN      100.0
Mark        68.0    81.0     100.0

# Sort descending according to Math grades
>>> grades.sort_values('Math', ascending=False)
      Math   English   History
Mark      81.0     68.0     100.0
Barbara   52.0     73.0     100.0
Eleanor   35.0     NaN      100.0
David      10.0     39.0     100.0
Greg        NaN     26.0     100.0
Lauren     NaN     99.0     100.0
```

Other methods used for manipulating `DataFrame` and `Series` panda structures can be found in Table 1.2.

Method	Description
<code>append()</code>	Concatenate two or more <code>Series</code> .
<code>drop()</code>	Remove the entries with the specified label or labels
<code>drop_duplicates()</code>	Remove duplicate values
<code>dropna()</code>	Drop null entries
<code>fillna()</code>	Replace null entries with a specified value or strategy
<code>reindex()</code>	Replace the index
<code>sample()</code>	Draw a random entry
<code>shift()</code>	Shift the index
<code>unique()</code>	Return unique values

Table 1.2: Methods for managing or modifying data in a pandas `Series` or `DataFrame`.

Problem 1. The file `budget.csv` contains the budget of a college student over the course of 4 years. Write a function that performs the following operations in this order:

1. Read in `budget.csv` as a `DataFrame` with the index as column 0. Hint: Use `index_col=0` to set the first column as the index when reading in the csv.
2. Reindex the columns such that amount spent on groceries is the first column and all other columns maintain the same ordering.
3. Sort the `DataFrame` in descending order by how much money was spent on `Groceries`.
4. Reset all values in the '`Rent`' column to `800.0`.
5. Reset all values in the first 5 data points to `0.0`.

Return the values of the updated `DataFrame` as a NumPy array.

Basic Data Manipulation

Because the primary pandas data structures are based off of `ndarray`, most NumPy functions work with pandas structures. For example, basic vector operations work as would be expected:

```
# Sum history and english grades of all students
>>> grades['English'] + grades['History']
Barbara    173.0
David      139.0
Eleanor     NaN
Greg       126.0
Lauren     199.0
Mark       168.0
dtype: float64

# Double all Math grades
>>> grades['Math']*2
Barbara    104.0
David      20.0
```

```
Eleanor      70.0
Greg        NaN
Lauren      NaN
Mark       162.0
Name: Math, dtype: float64
```

In addition to arithmetic, `Series` has a variety of other methods similar to NumPy arrays. A collection of these methods is found in Table 1.3.

Method	Returns
<code>abs()</code>	Object with absolute values taken (of numerical data)
<code>idxmax()</code>	The index label of the maximum value
<code>idxmin()</code>	The index label of the minimum value
<code>count()</code>	The number of non-null entries
<code>cumprod()</code>	The cumulative product over an axis
<code>cumsum()</code>	The cumulative sum over an axis
<code>max()</code>	The maximum of the entries
<code>mean()</code>	The average of the entries
<code>median()</code>	The median of the entries
<code>min()</code>	The minimum of the entries
<code>mode()</code>	The most common element(s)
<code>prod()</code>	The product of the elements
<code>sum()</code>	The sum of the elements
<code>var()</code>	The variance of the elements

Table 1.3: Numerical methods of the `Series` and `DataFrame` pandas classes.

Basic Statistical Functions

The pandas library allows us to easily calculate basic summary statistics of our data, which can be useful when we want a quick description of the data. The `describe()` function outputs several such summary statistics for each column in a `DataFrame`:

```
# Use describe to better understand the data
>>> grades.describe()
    Math   English   History
count    4.000000    5.000000     6.0
mean    44.500000   61.000000   100.0
std     29.827281   28.92231     0.0
min     10.000000   26.000000   100.0
25%    28.750000   39.000000   100.0
50%    43.500000   68.000000   100.0
75%    59.250000   73.000000   100.0
max    81.000000   99.000000   100.0
```

Functions for calculating means and variances, the covariance and correlation matrices, and other basic statistics are also available.

```
# Find the average grade for each student
```

```
>>> grades.mean(axis=1)
Barbara    75.000000
David      49.666667
Eleanor   67.500000
Greg       63.000000
Lauren    99.500000
Mark      83.000000
dtype: float64

# Give correlation matrix between subjects
>>> grades.corr()
          Math  English  History
Math      1.00000  0.84996     NaN
English   0.84996  1.00000     NaN
History    NaN      NaN      NaN
```

The method `rank()` can be used to rank the values in a data set, either within each entry or with each column. This function defaults ranking in ascending order: the least will be ranked 1 and the greatest will be ranked the highest number.

```
# Rank each student's performance in their classes in descending order
# (best to worst)
# The method keyword specifies what rank to use when ties occur.
>>> grades.rank(axis=1,method='max',ascending=False)
          Math  English  History
Barbara   3.0      2.0      1.0
David     3.0      2.0      1.0
Eleanor   2.0      NaN      1.0
Greg      NaN      2.0      1.0
Lauren    NaN      2.0      1.0
Mark      2.0      3.0      1.0
```

These methods can be very effective in interpreting data. For example, the `rank()` example above shows use that Barbara does best in History, then English, and then Math.

Dealing with Missing Data

Missing data is a ubiquitous problem in data science. Fortunately, pandas is particularly well-suited to handling missing or anomalous data. As we have already seen, the pandas default for a missing value is `NaN`. In basic arithmetic operations, if one of the operands is `NaN`, then the output is also `NaN`. If we are not interested in the missing values, we can simply drop them from the data altogether, or we can fill them with some other value, such as the mean. `NaN` might also mean something specific, such as some default value, which should inform what to do with `NaN` values.

```
# Grades with all NaN values dropped
>>> grades.dropna()
          Math  English  History
Barbara   52.0     73.0    100.0
David     10.0     39.0    100.0
```

```

Mark      81.0    68.0    100.0

# fill missing data with 50.0
>>> grades.fillna(50.0)
      Math   English   History
Barbara  52.0     73.0    100.0
David    10.0     39.0    100.0
Eleanor  35.0     50.0    100.0
Greg     50.0     26.0    100.0
Lauren   50.0     99.0    100.0
Mark     81.0    68.0    100.0

```

When dealing with missing data, make sure you are aware of the behavior of the pandas functions you are using. For example, `sum()` and `mean()` ignore NaN values in the computation.

Achtung!

Always consider missing data carefully when analyzing a dataset. It may not always be helpful to drop the data or fill it in with a random number. Consider filling the data with the mean of surrounding data or the mean of the feature in question. Overall, the choice for how to fill missing data should make sense with the dataset.

Problem 2. Write a function which uses `budget.csv` to answer the questions "Which category affects living expenses the most? Which affects other expenses the most?" Perform the following manipulations:

1. Fill all NaN values with 0.0.
2. Create two new columns, '`Living Expenses`' and '`Other`'. Set the value of '`Living Expenses`' to be the sum of the columns '`Rent`', '`Groceries`', '`Gas`' and '`Utilities`'. Set the value of '`Other`' to be the sum of the columns '`Dining Out`', '`Out With Friends`' and '`Netflix`'.
3. Identify which column, other than '`Living Expenses`', correlates most with '`Living Expenses`' and which column, other than '`Other`', correlates most with '`Other`'. This can indicate which columns in the budget affect the overarching categories the most.

Return the names of each of those columns as a tuple. The first should be of the column corresponding to '`Living Expenses`' and the second to '`Other`'.

Complex Operations in Pandas

Often times, the data that we have is not exactly the data we want to analyze. In cases like this we use more complex data manipulation tools to access only the data that we need.

For the examples below, we will use the following data:

```
>>> name = ['Mylan', 'Regan', 'Justin', 'Jess', 'Jason', 'Remi', 'Matt',
...           'Alexander', 'JeanMarie']
>>> sex = ['M', 'F', 'M', 'F', 'M', 'F', 'M', 'M', 'F']
>>> age = [20, 21, 18, 22, 19, 20, 20, 19, 20]
>>> rank = ['Sp', 'Se', 'Fr', 'Se', 'Sp', 'J', 'J', 'J', 'Se']
>>> ID = range(9)
>>> aid = ['y', 'n', 'n', 'y', 'n', 'n', 'n', 'y', 'n']
>>> GPA = [3.8, 3.5, 3.0, 3.9, 2.8, 2.9, 3.8, 3.4, 3.7]
>>> mathID = [0, 1, 5, 6, 3]
>>> mathGd = [4.0, 3.0, 3.5, 3.0, 4.0]
>>> major = ['y', 'n', 'y', 'n', 'n']
>>> studentInfo = pd.DataFrame({'ID': ID, 'Name': name, 'Sex': sex, 'Age': age,
...                               'Class': rank})
>>> otherInfo = pd.DataFrame({'ID': ID, 'GPA': GPA, 'Financial_Aid': aid})
>>> mathInfo = pd.DataFrame({'ID': mathID, 'Grade': mathGd, 'Math_Major':
...                           major})
```

Before querying our data, it is helpful to know some of its basic properties, such as number of columns, number of rows, and the datatypes of the columns. This can be done by simply calling the `info()` method on the desired `DataFrame`:

```
>>> mathInfo.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 5 entries, 0 to 4
Data columns (total 3 columns):
Grade      5 non-null float64
ID         5 non-null int64
Math_Major 5 non-null object
dtypes: float64(1), int64(1), object(1)
```

Masks

Sometimes, we only want to access data from a single column. For example if we want to only access the ID of the students in the `studentInfo DataFrame`, then we would use the following syntax.

```
# Get the ID column from studentInfo
>>> studentInfo.ID # or studentInfo['ID']
   ID
0   0
1   1
2   2
3   3
4   4
5   5
6   6
7   7
8   8
```

If we want to access multiple columns at once we can use a list of column names.

```
# Get the ID and Age columns.
>>> studentInfo[['ID', 'Age']]
   ID  Age
0    0   20
1    1   21
2    2   18
3    3   22
4    4   19
5    5   20
6    6   20
7    7   19
8    8   29
```

Now we can access the specific columns that we want. However, some of these columns may still contain data points that we don't want to consider. In this case we can build a mask. Each mask that we build will return a pandas `Series` object with a `bool` value at each index indicating if the condition is satisfied.

```
# Create a mask for all student receiving financial aid.
>>> mask = otherInfo['Financial_Aid'] == 'y'
# Access other info where the mask is true and display the ID and GPA ←
# columns.
>>> otherInfo[mask][['ID', 'GPA']]
   ID  GPA
0    0  3.8
3    3  3.9
7    7  3.4
```

We can also create compound masks with multiple statements. We do this using the same syntax you would use for a compound mask in a normal NumPy array. Useful operators are `&`, the AND operator; `|`, the OR operator; and `~`, the NOT operator.

```
# Get all student names where Class = 'J' OR Class = 'Sp'.
>>> mask = (studentInfo.Class == 'J') | (studentInfo.Class == 'Sp')
>>> studentInfo[mask].Name
0      Mylan
4     Jason
5      Remi
6      Matt
7  Alexander
Name: Name, dtype: object
# This can also be accomplished with the following command:
# studentInfo[studentInfo['Class'].isin(['J','Sp'])]['Name']
```

Problem 3. Read in the file `crime_data.csv` as a pandas object. The file contains data on types of crimes in the U.S. from 1960 to 2016. Set the index as the column '`Year`'. Answer the following questions using the pandas methods learned in this lab. The answer of each question should be saved as indicated. Return the answers to all three questions as a tuple (i.e. `(answer_1, answer_2, answer_3)`).

1. Identify the three crimes that have a mean yearly number of occurrences over 1,500,000. Of these three crimes, which two are very correlated? Which of these two crimes has a greater maximum value? Save the title of this column as a variable to return as the answer.
2. Examine the data from 2000 and later. Sort this data (in ascending order) according to number of murders. Find the years where aggravated assault is greater than 850,000. Save the indices (the years) of the masked and reordered `DataFrame` as a NumPy array to return as the answer.
3. What year had the highest crime rate? In this year, which crime was committed the most? What percentage of the total crime that year was it? Save this value as a float.

Working with Dates and Times

The `datetime` module in the standard library provides a few tools for representing and operating on dates and times. The `datetime.datetime` object represents a time stamp: a specific time of day on a certain day. Its constructor accepts a four-digit year, a month (starting at 1 for January), a day, and, optionally, an hour, minute, second, and microsecond. Each of these arguments must be an integer, with the hour ranging from 0 to 23.

```
>>> from datetime import datetime

# Represent November 18th, 1991, at 2:01 PM.
>>> bday = datetime(1991, 11, 18, 14, 1)
>>> print(bday)
1991-11-18 14:01:00

# Find the number of days between 11/18/1991 and 11/9/2017.
>>> dt = datetime(2017, 11, 9) - bday
>>> dt.days
9487
```

The `datetime.datetime` object has a parser method, `strptime()`, that converts a string into a new `datetime.datetime` object. The parser is flexible so the user must specify the format that the dates are in. For example, if the dates are in the format "`Month/Day//Year::Hour`", specify `format="%m/%d//%Y::%H"` to parse the string appropriately. See Table 1.4 for formatting options.

Pattern	Description
%Y	4-digit year
%y	2-digit year
%m	1- or 2-digit month
%d	1- or 2-digit day
%H	Hour (24-hour)
%I	Hour (12-hour)
%M	2-digit minute
%S	2-digit second

Table 1.4: Formats recognized by `datetime.strptime()`

```
>>> print(datetime.strptime("1991-11-18 / 14:01", "%Y-%m-%d / %H:%M"),
...       datetime.strptime("1/22/1996", "%m/%d/%Y"),
...       datetime.strptime("19-8, 1998", "%d-%m, %Y"), sep='\n')
1991-11-18 14:01:00          # The date formats are now standardized.
1996-01-22 00:00:00          # If no hour/minute/seconds data is given,
1998-08-19 00:00:00          # the default is midnight.
```

Converting Dates to an Index

The `TimeStamp` class is the pandas equivalent to a `datetime.datetime` object. A pandas index composed of `TimeStamp` objects is a `DatetimeIndex`, and a `Series` or `DataFrame` with a `DatetimeIndex` is called a time series. The function `pd.to_datetime()` converts a collection of dates in a parsable format to a `DatetimeIndex`. The format of the dates is inferred if possible, but it can be specified explicitly with the same syntax as `datetime.strptime()`.

```
>>> import pandas as pd

# Convert some dates (as strings) into a DatetimeIndex.
>>> dates = ["2010-1-1", "2010-2-1", "2012-1-1", "2012-1-2"]
>>> pd.to_datetime(dates)
DatetimeIndex(['2010-01-01', '2010-02-01', '2012-01-01', '2012-01-02'],
               dtype='datetime64[ns]', freq=None)

# Create a time series, specifying the format for the DatetimeIndex.
>>> dates = ["1/1, 2010", "1/2, 2010", "1/1, 2012", "1/2, 2012"]
>>> date_index = pd.to_datetime(dates, format="%m/%d, %Y")
>>> pd.Series([x**2 for x in range(4)], index=date_index)
2010-01-01    0
2010-01-02    1
2012-01-01    4
2012-01-02    9
dtype: int64
```

Problem 4. The file `DJIA.csv` contains daily closing values of the Dow Jones Industrial Average from 2006–2016. Read the data into a `Series` or `DataFrame` with a `DatetimeIndex` as the index. Drop any rows without numerical values, cast the "`VALUE`" column to floats, then return the updated `DataFrame`.

Hint: You can change the column type the same way you'd change a numpy array type.

Generating Time-based Indices

Some time series datasets come without explicit labels but have instructions for deriving timestamps. For example, a list of bank account balances might have records from the beginning of every month, or heart rate readings could be recorded by an app every 10 minutes. Use `pd.date_range()` to generate a `DatetimeIndex` where the timestamps are equally spaced. The function is analogous to `np.arange()` and has the following parameters:

Parameter	Description
<code>start</code>	Starting date
<code>end</code>	End date
<code>periods</code>	Number of dates to include
<code>freq</code>	Amount of time between consecutive dates
<code>normalize</code>	Normalizes the start and end times to midnight

Table 1.5: Parameters for `pd.date_range()`.

Exactly three of the parameters `start`, `end`, `periods`, and `freq` must be specified to generate a range of dates. The `freq` parameter accepts a variety of string representations, referred to as offset aliases. See Table 1.6 for a sampling of some of the options. For a complete list of the options, see https://pandas.pydata.org/pandas-docs/stable/user_guide/timeseries.html#timeseries-offset-aliases1.

Parameter	Description
<code>"D"</code>	calendar daily (default)
<code>"B"</code>	business daily (every business day)
<code>"H"</code>	hourly
<code>"T"</code>	minutely
<code>"S"</code>	secondly
<code>"MS"</code>	first day of the month (Month Start)
<code>"BMS"</code>	first business day of the month (Business Month Start)
<code>"W-MON"</code>	every Monday (Week-Monday)
<code>"WOM-3FRI"</code>	every 3rd Friday of the month (Week of the Month - 3rd Friday)

Table 1.6: Options for the `freq` parameter to `pd.date_range()`.

```
# Create a DatetimeIndex for 5 consecutive days starting on September 28, 2016.
>>> pd.date_range(start='9/28/2016 16:00', periods=5)
DatetimeIndex(['2016-09-28 16:00:00', '2016-09-29 16:00:00',
```

```

'2016-09-30 16:00:00', '2016-10-01 16:00:00',
'2016-10-02 16:00:00'],
dtype='datetime64[ns]', freq='D')

# Create a DatetimeIndex with the first weekday of every other month in 2016.
>>> pd.date_range(start='1/1/2016', end='1/1/2017', freq="2BMS")
DatetimeIndex(['2016-01-01', '2016-03-01', '2016-05-02', '2016-07-01',
               '2016-09-01', '2016-11-01'],
              dtype='datetime64[ns]', freq='2BMS')

# Create a DatetimeIndex for 10 minute intervals between 4:00 PM and 4:30 PM on ←
# September 9, 2016.
>>> pd.date_range(start='9/28/2016 16:00',
                  end='9/28/2016 16:30', freq="10T")
DatetimeIndex(['2016-09-28 16:00:00', '2016-09-28 16:10:00',
               '2016-09-28 16:20:00', '2016-09-28 16:30:00'],
              dtype='datetime64[ns]', freq='10T')

# Create a DatetimeIndex for 2 hour 30 minute intervals between 4:30 PM and ←
# 2:30 AM on September 29, 2016.
>>> pd.date_range(start='9/28/2016 16:30', periods=5, freq="2h30min")
DatetimeIndex(['2016-09-28 16:30:00', '2016-09-28 19:00:00',
               '2016-09-28 21:30:00', '2016-09-29 00:00:00',
               '2016-09-29 02:30:00'],
              dtype='datetime64[ns]', freq='150T')

```

Problem 5. The file `paychecks.csv` contains values of an hourly employee's last 93 paychecks. Paychecks are given every other Friday, starting on March 14, 2008, and the employee started working on March 13, 2008.

Read in the data, using `pd.date_range()` to generate the `DatetimeIndex`. Set this as the new index of the `DataFrame` and return the `DataFrame`.

Elementary Time Series Analysis

Shifting

`DataFrame` and `Series` objects have a `shift()` method that allows you to move data up or down relative to the index. When dealing with time series data, we can also shift the `DatetimeIndex` relative to a time offset.

```

>>> df = pd.DataFrame(dict(VALUE=np.random.rand(5)),
                      index=pd.date_range("2016-10-7", periods=5, freq='D'))
>>> df
          VALUE
2016-10-07  0.127895

```

```

2016-10-08  0.811226
2016-10-09  0.656711
2016-10-10  0.351431
2016-10-11  0.608767

>>> df.shift(1)
      VALUE
2016-10-07      NaN
2016-10-08  0.127895
2016-10-09  0.811226
2016-10-10  0.656711
2016-10-11  0.351431

>>> df.shift(-2)
      VALUE
2016-10-07  0.656711
2016-10-08  0.351431
2016-10-09  0.608767
2016-10-10      NaN
2016-10-11      NaN

>>> df.shift(14, freq="D")
      VALUE
2016-10-21  0.127895
2016-10-22  0.811226
2016-10-23  0.656711
2016-10-24  0.351431
2016-10-25  0.608767

```

Shifting data makes it easy to gather statistics about changes from one timestamp or period to the next.

```

# Find the changes from one period/timestamp to the next
>>> df - df.shift(1)           # Equivalent to df.diff().
      VALUE
2016-10-07      NaN
2016-10-08  0.683331
2016-10-09 -0.154516
2016-10-10 -0.305279
2016-10-11  0.257336

```

Problem 6. Compute the following information about the DJIA dataset from Problem 4 that has a DateTimeIndex.

- The single day with the largest gain.
- The single day with the largest loss.

Return the DateTimeIndex of the day with the largest gain and the day with the largest loss.

(Hint: Call your function from Problem 4 to get the DataFrame already cleaned and with DatetimeIndex).

More information on how to use `datetime` with Pandas is in the additional material section. This includes working with `Periods` and more analysis with time series.

Additional Material

SQL Operations in pandas

DataFrames are tabular data structures bearing an obvious resemblance to a typical relational database table. SQL is the standard for working with relational databases; however, pandas can accomplish many of the same tasks as SQL. The SQL-like functionality of pandas is one of its biggest advantages, eliminating the need to switch between programming languages for different tasks. Within pandas, we can handle both the querying and data analysis.

For the examples below, we will use the following data:

```
>>> name = ['Mylan', 'Regan', 'Justin', 'Jess', 'Jason', 'Remi', 'Matt',
...          'Alexander', 'JeanMarie']
>>> sex = ['M', 'F', 'M', 'F', 'M', 'F', 'M', 'M', 'F']
>>> age = [20, 21, 18, 22, 19, 20, 20, 19, 20]
>>> rank = ['Sp', 'Se', 'Fr', 'Se', 'Sp', 'J', 'J', 'J', 'Se']
>>> ID = range(9)
>>> aid = ['y', 'n', 'n', 'y', 'n', 'n', 'n', 'y', 'n']
>>> GPA = [3.8, 3.5, 3.0, 3.9, 2.8, 2.9, 3.8, 3.4, 3.7]
>>> mathID = [0, 1, 5, 6, 3]
>>> mathGd = [4.0, 3.0, 3.5, 3.0, 4.0]
>>> major = ['y', 'n', 'y', 'n', 'n']
>>> studentInfo = pd.DataFrame({'ID': ID, 'Name': name, 'Sex': sex, 'Age': age,
...                                'Class': rank})
>>> otherInfo = pd.DataFrame({'ID': ID, 'GPA': GPA, 'Financial_Aid': aid})
>>> mathInfo = pd.DataFrame({'ID': mathID, 'Grade': mathGd, 'Math_Major':
...                                major})
```

SQL SELECT statements can be done by column indexing. WHERE statements can be included by adding masks (just like in a NumPy array). The method `isin()` can also provide a useful WHERE statement. This method accepts a list, dictionary, or Series containing possible values of the DataFrame or Series. When called upon, it returns a Series of booleans, indicating whether an entry contained a value in the parameter pass into `isin()`.

```
# SELECT ID, Age FROM studentInfo
>>> studentInfo[['ID', 'Age']]
   ID  Age
0    0   20
1    1   21
2    2   18
3    3   22
4    4   19
5    5   20
6    6   20
7    7   19
8    8   29

# SELECT ID, GPA FROM otherInfo WHERE Financial_Aid = 'y'
>>> mask = otherInfo['Financial_Aid'] == 'y'
>>> otherInfo[mask][['ID', 'GPA']]
```

```

      ID  GPA
0    0  3.8
3    3  3.9
7    7  3.4

# SELECT Name FROM studentInfo WHERE Class = 'J' OR Class = 'Sp'
>>> studentInfo[studentInfo['Class'].isin(['J', 'Sp'])]['Name']
0        Mylan
4       Jason
5       Remi
6       Matt
7  Alexander
Name: Name, dtype: object

```

Next, let's look at JOIN statements. In pandas, this is done with the `merge` function. `merge` takes the two `DataFrame` objects to join as parameters, as well as keyword arguments specifying the column on which to join, along with the type (left, right, inner, outer).

```

# SELECT * FROM studentInfo INNER JOIN mathInfo ON studentInfo.ID = mathInfo.ID
>>> pd.merge(studentInfo, mathInfo, on='ID') # INNER JOIN is the default
   Age Class  ID    Name  Sex  Grade Math_Major
0   20    Sp   0  Mylan   M    4.0        y
1   21    Se   1  Regan   F    3.0        n
2   22    Se   3   Jess   F    4.0        n
3   20     J   5  Remi   F    3.5        y
4   20     J   6   Matt   M    3.0        n
[5 rows x 7 columns]

# SELECT GPA, Grade FROM otherInfo FULL OUTER JOIN mathInfo ON otherInfo.
# ID = mathInfo.ID
>>> pd.merge(otherInfo, mathInfo, on='ID', how='outer')[['GPA', 'Grade']]
   GPA  Grade
0  3.8    4.0
1  3.5    3.0
2  3.0    NaN
3  3.9    4.0
4  2.8    NaN
5  2.9    3.5
6  3.8    3.0
7  3.4    NaN
8  3.7    NaN
[9 rows x 2 columns]

```

More Datetime with Pandas

Periods

A pandas `Timestamp` object represents a precise moment in time on a given day. Some data, however, is recorded over a time interval, and it wouldn't make sense to place an exact timestamp on any of the measurements. For example, a record of the number of steps walked in a day, box office earnings per week, quarterly earnings, and so on. This kind of data is better represented with the pandas `Period` object and the corresponding `PeriodIndex`.

The `Period` class accepts a `value` and a `freq`. The `value` parameter indicates the label for a given `Period`. This label is tied to the `end` of the defined `Period`. The `freq` indicates the length of the `Period` and in some cases can also indicate the offset of the `Period`. The default value for `freq` is "M" for months. The `freq` parameter accepts the majority, but not all, of frequencies listed in Table 1.6.

```
# Creates a period for month of Oct, 2016.
>>> p1 = pd.Period("2016-10")
>>> p1.start_time                      # The start and end times of the period
Timestamp('2016-10-01 00:00:00')      # are recorded as Timestamps.
>>> p1.end_time
Timestamp('2016-10-31 23:59:59.999999999')

# Represent the annual period ending in December that includes 10/03/2016.
>>> p2 = pd.Period("2016-10-03", freq="A-DEC")
>>> p2.start_time
Timestamp('2016-01-01 00:00:00')
> p2.end_time
Timestamp('2016-12-31 23:59:59.999999999')

# Get the weekly period ending on a Saturday that includes 10/03/2016.
>>> print(pd.Period("2016-10-03", freq="W-SAT"))
2016-10-02/2016-10-08
```

Like the `pd.date_range()` method, the `pd.period_range()` method is useful for generating a `PeriodIndex` for unindexed data. The syntax is essentially identical to that of `pd.date_range()`. When using `pd.period_range()`, remember that the `freq` parameter marks the end of the period. After creating a `PeriodIndex`, the `freq` parameter can be changed via the `asfreq()` method.

```
# Represent quarters from 2008 to 2010, with Q4 ending in December.
>>> pd.period_range(start="2008", end="2010-12", freq="Q-DEC")
PeriodIndex(['2008Q1', '2008Q2', '2008Q3', '2008Q4', '2009Q1', '2009Q2',
             '2009Q3', '2009Q4', '2010Q1', '2010Q2', '2010Q3', '2010Q4'],
            dtype='period[Q-DEC]', freq='Q-DEC')

# Get every three months from March 2010 to the start of 2011.
>>> p = pd.period_range("2010-03", "2011", freq="3M")
>>> p
PeriodIndex(['2010-03', '2010-06', '2010-09', '2010-12'],
            dtype='period[3M]', freq='3M')
```

```
# Change frequency to be quarterly.
>>> p.asfreq("Q-DEC")
PeriodIndex(['2010Q2', '2010Q3', '2010Q4', '2011Q1'],
            dtype='period[Q-DEC]', freq='Q-DEC')
```

The bounds of a `PeriodIndex` object can be shifted by adding or subtracting an integer. `PeriodIndex` will be shifted by $n \times \text{freq}$.

```
# Shift index by 1
>>> p -= 1
>>> p
PeriodIndex(['2010Q1', '2010Q2', '2010Q3', '2010Q4'],
            dtype='int64', freq='Q-DEC')
```

If for any reason you need to switch from periods to timestamps, pandas provides a very simple method to do so. The `how` parameter can be `start` or `end` and determines if the timestamp is the beginning or the end of the period. Similarly, you can switch from timestamps to periods.

```
# Convert to timestamp (last day of each quarter)
>>> p = p.to_timestamp(how='end')
>>> p
DatetimeIndex(['2010-03-31', '2010-06-30', '2010-09-30', '2010-12-31'],
               dtype='datetime64[ns]', freq='Q-DEC')

>>> p.to_period("Q-DEC")
PeriodIndex(['2010Q1', '2010Q2', '2010Q3', '2010Q4'],
            dtype='int64', freq='Q-DEC')
```

Operations on Time Series

There are certain operations only available to Series and DataFrames that have a `DatetimeIndex`. A sampling of this functionality is described throughout the remainder of this lab.

Slicing

Slicing is much more flexible in pandas for time series. We can slice by year, by month, or even use traditional slicing syntax to select a range of dates.

```
# Select all rows in a given year
>>> df["2010"]
              0          1
2010-01-01  0.566694  1.093125
2010-02-01 -0.219856  0.852917
2010-03-01  1.511347 -1.324036

# Select all rows in a given month of a given year
>>> df["2012-01"]
              0          1
```

```

2012-01-01  0.212141  0.859555
2012-01-02  1.483123 -0.520873
2012-01-03  1.436843  0.596143

# Select a range of dates using traditional slicing syntax
>>> df["2010-1-2":"2011-12-31"]
          0      1
2010-02-01 -0.219856  0.852917
2010-03-01  1.511347 -1.324036
2011-01-01  0.300766  0.934895

```

Resampling

Some datasets do not have datapoints at a fixed frequency. For example, a dataset of website traffic has datapoints that occur at irregular intervals. In situations like these, resampling can help provide insight on the data.

The two main forms of resampling are downsampling, aggregating data into fewer intervals, and upsampling, adding more intervals.

To downsample, use the `resample()` method of the `Series` or `DataFrame`. This method is similar to `groupby()` in that it groups different entries together. Then aggregation produces a new data set. The first parameter to `resample()` is an offset string from Table 1.6: "`D`" for daily, "`H`" for hourly, and so on.

```

>>> import numpy as np

# Get random data for every day from 2000 to 2010.
>>> dates = pd.date_range(start="2000-1-1", end='2009-12-31', freq='D')
>>> df = pd.Series(np.random(len(dates)), index=dates)
>>> df
2000-01-01    0.559
2000-01-02    0.874
2000-01-03    0.774
...
2009-12-29    0.837
2009-12-30    0.472
2009-12-31    0.211
Freq: D, Length: 3653, dtype: float64

# Group the data by year.
>>> years = df.resample("A")           # 'A' for 'annual'.
>>> years.agg(len)                 # Number of entries per year.
2000-12-31    366.0
2001-12-31    365.0
2002-12-31    365.0
...
2007-12-31    365.0
2008-12-31    366.0
2009-12-31    365.0

```

```

Freq: A-DEC, dtype: float64

>>> years.mean()                               # Average entry by year.
2000-12-31    0.491
2001-12-31    0.514
2002-12-31    0.484
...
2007-12-31    0.508
2008-12-31    0.521
2009-12-31    0.523
Freq: A-DEC, dtype: float64

# Group the data by month.
>>> months = df.resample("M")
>>> len(months.mean())                         # 12 months x 10 years = 120 months.
120

```

Elementary Time Series Analysis

Rolling Functions and Exponentially-Weighted Moving Functions

Many time series are inherently noisy. To analyze general trends in data, we use rolling functions and exponentially-weighted moving (EWM) functions. Rolling functions, or moving window functions, perform a calculation on a window of data. There are a few rolling functions that come standard with pandas.

Rolling Functions (Moving Window Functions)

One of the most commonly used rolling functions is the rolling average, which takes the average value over a window of data.

```

# Generate a time series using random walk from a uniform distribution.
N = 10000
bias = 0.01
s = np.zeros(N)
s[1:] = np.random.uniform(low=-1, high=1, size=N-1) + bias
s = pd.Series(s.cumsum(),
              index=pd.date_range("2015-10-20", freq='H', periods=N))

# Plot the original data together with a rolling average.
ax1 = plt.subplot(121)
s.plot(color="gray", lw=.3, ax=ax1)
s.rolling(window=200).mean().plot(color='r', lw=1, ax=ax1)
ax1.legend(["Actual", "Rolling"], loc="lower right")
ax1.set_title("Rolling Average")

```

The function call `s.rolling(window=200)` creates a `pd.core.rolling.Window` object that can be aggregated with a function like `mean()`, `std()`, `var()`, `min()`, `max()`, and so on.

Exponentially-Weighted Moving (EWM) Functions

Whereas a moving window function gives equal weight to the whole window, an exponentially-weighted moving function gives more weight to the most recent data points.

In the case of a exponentially-weighted moving average (EWMA), each data point is calculated as follows.

$$z_i = \alpha \bar{x}_i + (1 - \alpha) z_{i-1},$$

where z_i is the value of the EWMA at time i , \bar{x}_i is the average for the i -th window, and α is the decay factor that controls the importance of previous data points. Notice that $\alpha = 1$ reduces to the rolling average.

More commonly, the decay is expressed as a function of the window size. In fact, the `span` for an EWMA is nearly analogous to `window size` for a rolling average.

Notice the syntax for EWM functions is very similar to that of rolling functions.

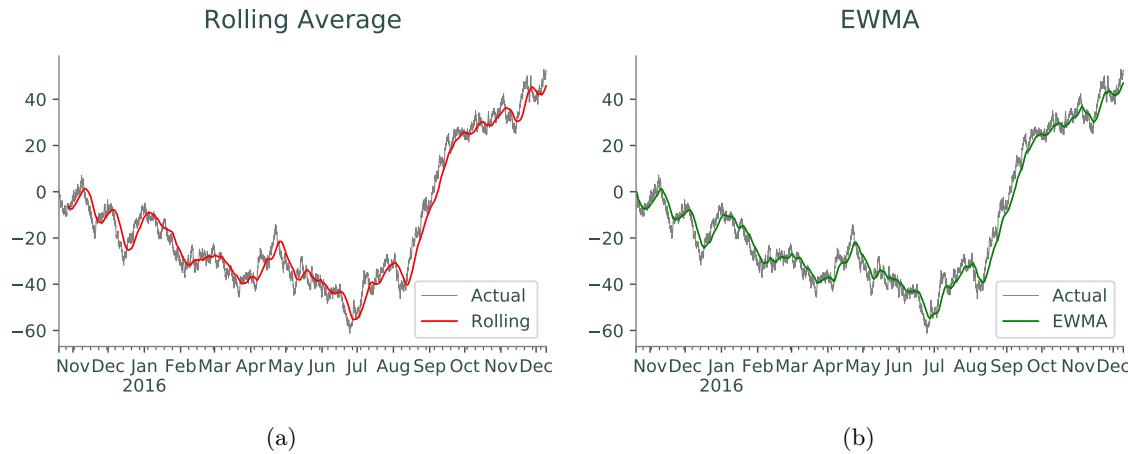


Figure 1.1: Rolling average and EWMA.

```
ax2 = plt.subplot(122)
s.plot(color="gray", lw=.3, ax=ax2)
s.ewm(span=200).mean().plot(color='g', lw=1, ax=ax2)
ax2.legend(["Actual", "EWMA"], loc="lower right")
ax2.set_title("EWMA")
```

2

Pandas 2: Plotting

Lab Objective: Clear, insightful visualizations are a crucial part of data analysis. To facilitate quick data visualization, pandas includes several tools that wrap around matplotlib. These tools make it easy to compare different parts of a data set, explore the data as a whole, and spot patterns and correlations in the data.

Overview of Plotting Tools

The main tool for visualization in pandas is the `plot()` method for `Series` and `DataFrames`. The method has a keyword argument `kind` that specifies the type of plot to draw. The valid options for `kind` are detailed below.

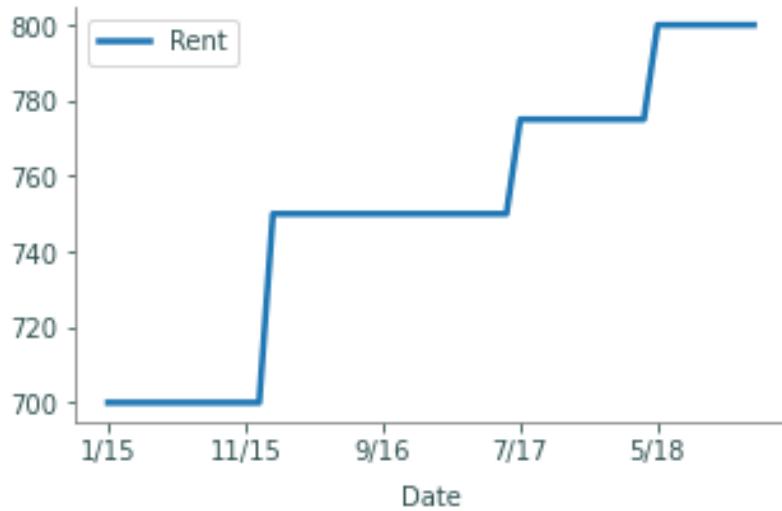
Plot Type	plot() ID	Uses and Advantages
Line plot	<code>"line"</code>	Show trends ordered in data; easy to compare multiple data sets
Scatter plot	<code>"scatter"</code>	Compare exactly two data sets, independent of ordering
Bar plot	<code>"bar", "barh"</code>	Compare categorical or sequential data
Histogram	<code>"hist"</code>	Show frequencies of one set of values, independent of ordering
Box plot	<code>"box"</code>	Display min, median, max, and quartiles; compare data distributions
Hexbin plot	<code>"hexbin"</code>	2D histogram; reveal density of cluttered scatter plots

Table 2.1: Types of plots in pandas. The plot ID is the value of the keyword argument `kind`. That is, `df.plot(kind="scatter")` creates a scatter plot. The default `kind` is `"line"`.

The `plot()` method calls `plt.plot()`, `plt.hist()`, `plt.scatter()`, and other matplotlib plotting functions, but it also assigns axis labels, tick marks, legends, and a few other things based on the index and the data. Most calls to `plot()` specify the kind of plot and which `Series` to use as the x and y axes. By default, the `index` of the `Series` or `DataFrame` is used for the x axis.

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

>>> budget = pd.read_csv("budget.csv", index_col="Date")
>>> budget.plot(y="Rent") # Plot rent against the index (date).
```



In this case, the call to the `plot()` method is essentially equivalent to the following code.

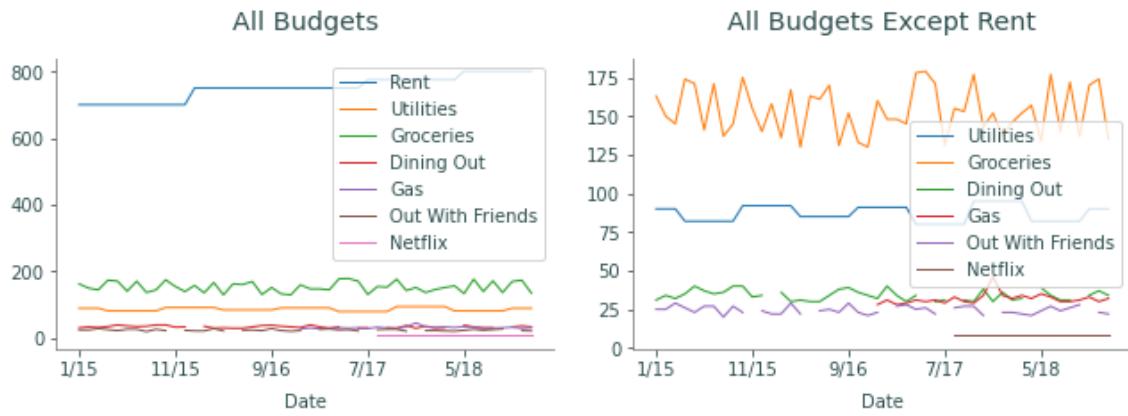
```
>>> plt.plot(budget.index, budget['Rent'], label='Rent')
>>> plt.xlabel(budget.index.name)
>>> plt.xlim(min(budget.index), max(budget.index))
>>> plt.legend(loc='best')
```

The `plot()` method also takes in many keyword arguments for matplotlib plotting and annotation functions. For example, setting `legend=False` disables the legend, providing a value for `title` sets the figure title, `grid=True` turns a grid on, and so on. For more customizations, see <https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.plot.html>.

Visualizing an Entire Data Set

A good way to start analyzing an unfamiliar data set is to visualize as much of the data as possible to determine which parts are most important or interesting. For example, since the columns in a `DataFrame` share the same index, the columns can all be graphed together using the index as the *x*-axis. By default, the `plot()` method attempts to plot **every Series** (column) in a `DataFrame`. This is especially useful with sequential data, like the budget data set.

```
# Plot all columns together against the index.
>>> budget.plot(title="All Budgets", linewidth=1)
>>> budget.drop(["Rent"], axis=1).plot(linewidth=1, title="All Budgets Except ←
    Rent")
```



- (a) All columns of the budget data set on the same figure, using the index as the *x*-axis.
(b) All columns of the budget data set except "Living Expenses" and "Rent".

Figure 2.1

While plotting every `Series` at once can give an overview of all the data, the resulting plot is often difficult for the reader to understand. For example, the budget data set has 9 columns, so the resulting figure, Figure 2.1a, is fairly cluttered.

One way to declutter a visualization is to examine less data. For example, the columns '`Living Expenses`' and '`Rent`' have values that are much larger than the other columns. Dropping these columns gives a better overview of the remaining data, as shown in Figure 2.1b.

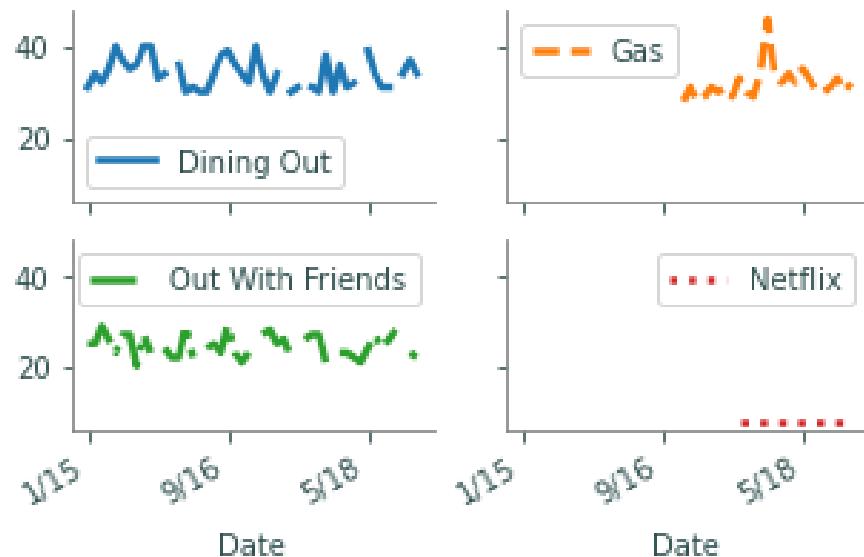
Achtung!

Often plotting all data at once is unwise because columns have **different units of measure**. Be careful not to plot parts of a data set together if those parts do not have the same units or are otherwise incomparable.

Another way to declutter a plot is to use subplots. To quickly plot several columns in separate subplots, use `subplots=True` and specify a shape tuple as the `layout` for the plots. Subplots automatically share the same *x*-axis. Set `sharey=True` to force them to share the same *y*-axis as well.

```
>>> budget.plot(y=['Dining Out', 'Gas', 'Out With Friends', 'Netflix'],
...     subplots=True, layout=(2,2), sharey=True,
...     style=['-', '--', '-.', ':'], title="Plots of Dollars Spent for Different ↵
...     Budgets")
```

Plots of Dollars Spent for Different Budgets



As mentioned previously, the `plot()` method can be used to plot different kinds of plots. One possible kind of plot is a histogram. Since plots made by the `plot()` method share an *x*-axis by default, histograms turn out poorly whenever there are columns with very different data ranges or when more than one column is plotted at once.

```
# Plot three histograms together.
>>> budget.plot(kind='hist',y=['Gas','Dining Out','Out With Friends'],
...     alpha=.7,bins=10,title="Frequency of Amount (in dollars) Spent")

# Plot three histograms, stacking one on top of the other.
>>> budget.plot(kind='hist',y=['Gas','Dining Out','Out With Friends'],
...     bins=10,stacked=True,title="Frequency of Amount (in dollars) Spent")
```

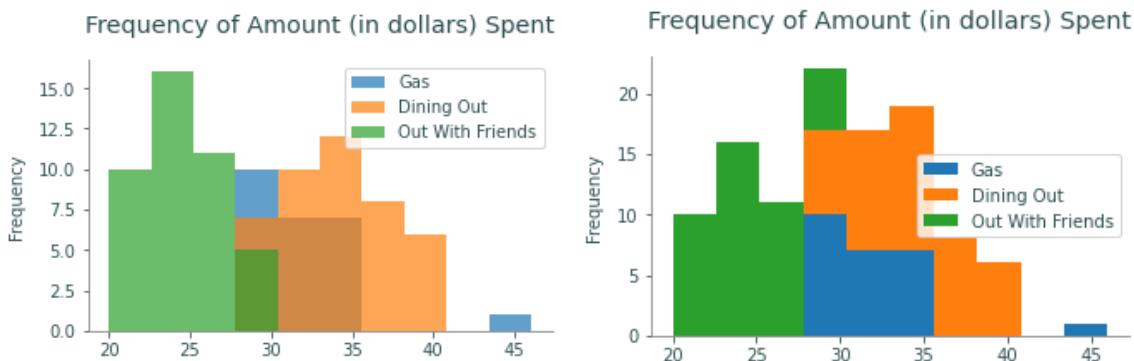


Figure 2.2: Two examples of histograms that are difficult to understand because multiple columns are plotted.

Thus, histograms are good for examining the distribution of a **single** column in a data set. For histograms, use the `hist()` method of the `DataFrame` instead of the `plot()` method. Specify the number of bins with the `bins` parameter. Choose a number of bins that accurately represents the data; the wrong number of bins can create a misleading or uninformative visualization.

```
>>> budget[["Dining Out", "Gas"]].hist(grid=False, bins=10)
```

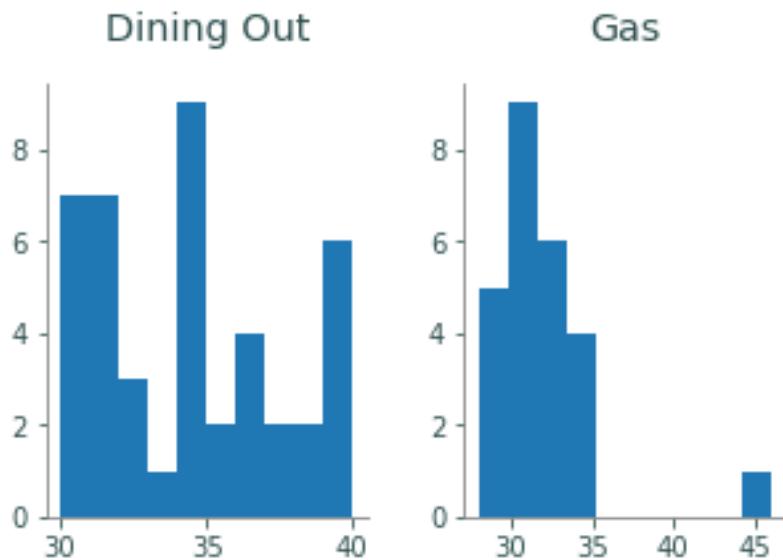


Figure 2.3: Histograms of "Dining Out" and "Gas".

Problem 1. Create 3 visualizations for the data in `crime_data.csv`. Make one of the visualizations a histogram. The visualizations should be well labeled and easy to understand.

Patterns and Correlations

After visualizing the entire data set initially, a good next step is to closely compare related parts of the data. This can be done with different types of visualizations. For example, Figure 2.1b suggests that the "Dining Out" and "Out With Friends" columns are roughly on the same scale. Since this data is sequential (indexed by time), start by plotting these two columns against the index. Next, create a scatter plot of one of the columns versus the other to investigate correlations that are independent of the index. Unlike other types of plots, using `kind="scatter"` requires both `x` and `y` columns as arguments.

```
# Plot 'Dining Out' and 'Out With Friends' as lines against the index.
>>> budget.plot(y=["Dining Out", "Out With Friends"], title="Amount Spent on ←
    Dining Out and Out with Friends per Day")
```

```
# Make a scatter plot of 'Dining Out' against 'Out With Friends'
>>> budget.plot(kind="scatter", x="Dining Out", y="Out With Friends",
...     alpha=.8,xlim=(0,max(budget['Dining Out'])+1),
...     ylim=(0,max(budget['Out With Friends'])+1))
```

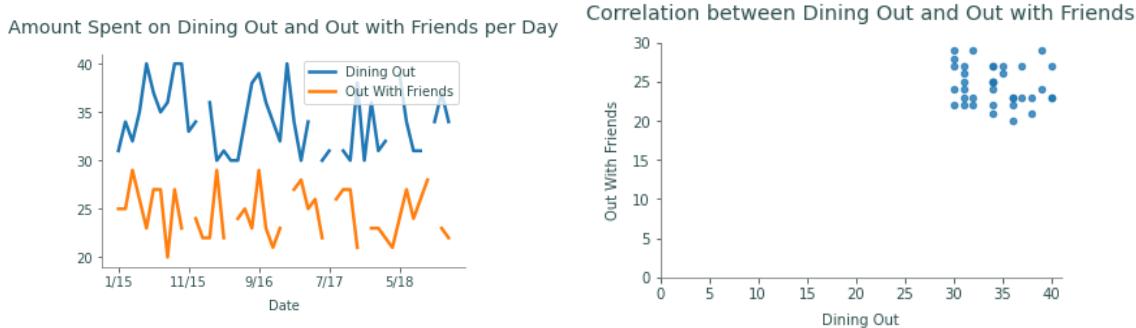


Figure 2.4: Correlations between "Dining Out" and "Out With Friends".

The first plot shows us that more money is spent on dining out than being out with friends overall. However, both categories stay in the same range for most of the data. This is confirmed in the scatter plot by the block in the upper right corner, indicating the common range spent on dining out and being out with friends.

Achtung!

When analyzing data, especially while searching for patterns and correlations, **always** ask yourself if the data makes sense and is trustworthy. What lurking variables could have influenced the data measurements as they were being gathered?

The crime data set from Problem 1 is somewhat suspect in this regard. The murder rate is likely accurate, since murder is conspicuous and highly reported, but what about the rape rate? Are the number of rapes increasing, or is the percentage of rapes being reported increasing? It's probably both! Be careful about drawing conclusions for sensitive or questionable data.

Another useful visualization used to understand correlations in a data set is a scatter matrix. The function `pd.plotting.scatter_matrix()` produces a table of plots where each column is plotted against each other column in separate scatter plots. The plots on the diagonal, instead of plotting a column against itself, displays a histogram of that column. This provides a very quick method for an initial analysis of the correlation between different columns.

```
>>> pd.plotting.scatter_matrix(budget[['Living Expenses', 'Other']])
```

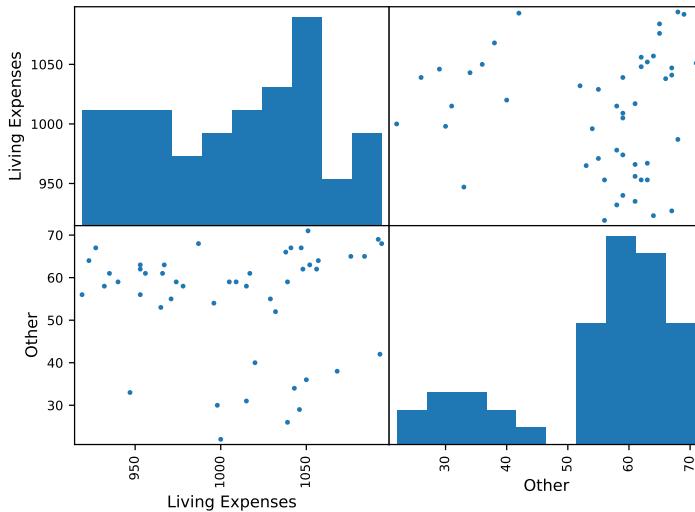


Figure 2.5: Scatter matrix comparing "Living Expenses" and "Other".

Bar Graphs

Different types of graphs help to identify different patterns. Note that the data set `budget` gives monthly expenses. It may be beneficial to look at one specific month. Bar graphs are a good way to compare small portions of the data set.

As a general rule, horizontal bar charts (`kind="barh"`) are better than the default vertical bar charts (`kind="bar"`) because most humans can detect horizontal differences more easily than vertical differences. If the labels are too long to fit on a normal figure, use `plt.tight_layout()` to adjust the plot boundaries to fit the labels in.

```
# Plot all data for the last month in the budget
>>> budget.iloc[-1,:].plot(kind='barh')
>>> plt.tight_layout()

# Plot all data for the last month without 'Rent' and 'Living Expenses'
>>> budget.drop(['Rent','Living Expenses'],axis=1).iloc[-1,:].plot(kind='barh')
>>> plt.tight_layout()
```

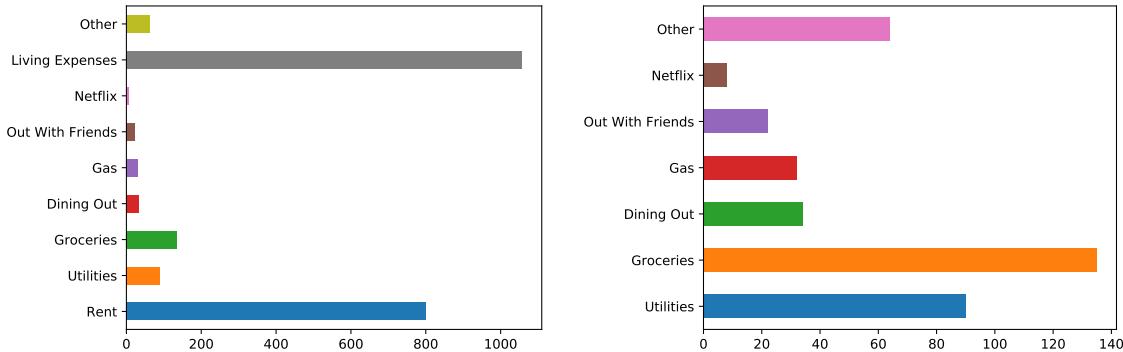


Figure 2.6: Bar graphs showing expenses paid in the last month of budget.

Problem 2. Using the crime data from the previous problem, identify if a trend exists between `Forcible Rape` and the following variables:

1. `Violent`
2. `Burglary`
3. `Aggravated Assault`

Make sure each graph is clearly labelled and readable. Return a tuple of booleans describing whether `Forcible Rape` correlates with each of the other variables.

Distributional Visualizations

While histograms are good at displaying the distributions for one column, a different visualization is needed to show the distribution of an entire set. A box plot, sometimes called a “cat-and-whisker” plot, shows the five number summary: the minimum, first quartile, median, third quartile, and maximum of the data. Box plots are useful for comparing the distributions of relatable data. However, box plots are a basic summary, meaning that they are susceptible to miss important information such as how many points were in each distribution.

```
# Compare the distributions of four columns.
>>> budget.plot(kind="box", y=["Gas", "Dining Out", "Out With Friends", "Other"])

# Compare the distributions of all columns but 'Rent' and 'Living Expenses'.
>>> budget.drop(["Rent", "Living Expenses"], axis=1).plot(kind="box",
...             vert=False)
```

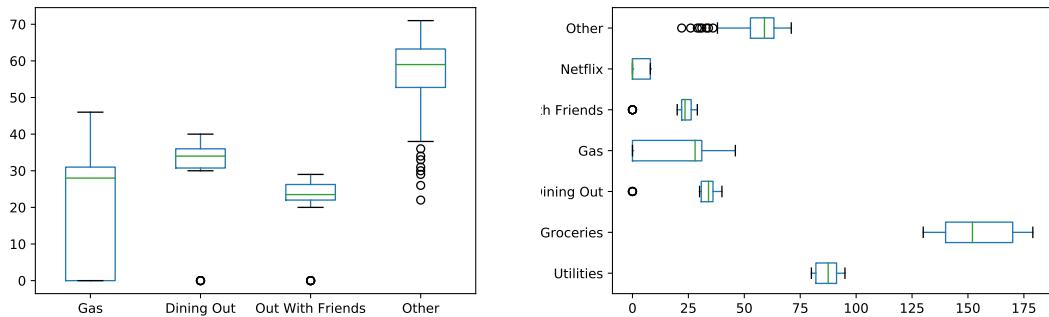


Figure 2.7: Vertical and horizontal box plots of `budget` dataset.

Hexbin Plots

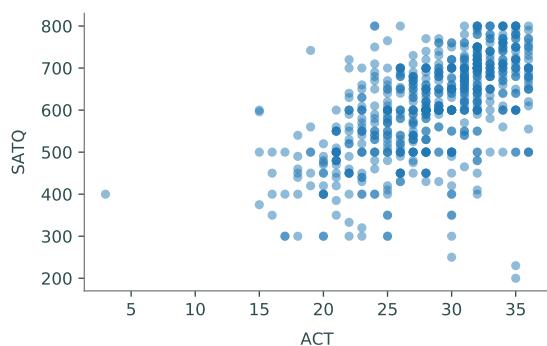
A scatter plot is essentially a plot of samples from the joint distribution of two columns. However, scatter plots can be uninformative for large data sets when the points in a scatter plot are closely clustered. Hexbin plots solve this problem by plotting point density in hexagonal bins—essentially creating a 2-dimensional histogram.

The file `sat_act.csv` contains 700 self reported scores on the SAT Verbal, SAT Quantitative and ACT, collected as part of the Synthetic Aperture Personality Assessment (SAPA) web based personality assessment project. The obvious question with this data set is “how correlated are ACT and SAT scores?” The scatter plot of ACT scores versus SAT Quantitative scores, Figure 2.8a, is highly cluttered, even though the points have some transparency. A hexbin plot of the same data, Figure 2.8b, reveals the **frequency** of points in binned regions.

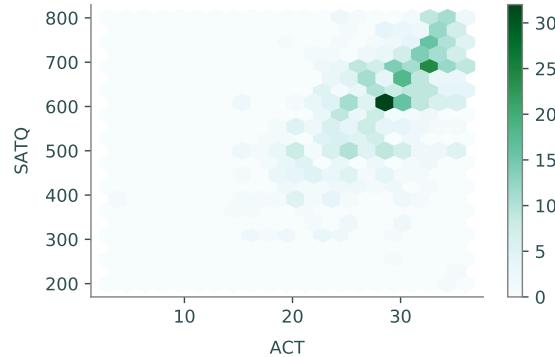
```
>>> satact = pd.read_csv("sat_act.csv", index_col="ID")
>>> list(satact.columns)
['gender', 'education', 'age', 'ACT', 'SATV', 'SATQ']

# Plot the ACT scores against the SAT Quant scores in a regular scatter plot.
>>> satact.plot(kind="scatter", x="ACT", y="SATQ", alpha=.8)

# Plot the densities of the ACT vs. SATQ scores with a hexbin plot.
>>> satact.plot(kind="hexbin", x="ACT", y="SATQ", gridsize=20)
```



(a) ACT vs. SAT Quant scores.



(b) Frequency of ACT vs. SAT Quant scores.

Figure 2.8: Scatter plots and hexbin plot of SAT and ACT scores.

Just as choosing a good number of `bins` is important for a good histogram, choosing a good `gridsize` is crucial for an informative hexbin plot. A large `gridsize` creates many small bins and a small `gridsize` creates fewer, larger bins.

Note

Since hexbins are based on frequencies, they are prone to being misleading if the dataset is not understood well. For example, when plotting information that deals with geographic position, increases in frequency may be results in higher populations rather than the actual information being plotted.

See <http://pandas.pydata.org/pandas-docs/stable/visualization.html> for more types of plots available in Pandas and further examples.

Problem 3. Use `crime_data.csv` to display the following distributions.

1. The distributions of `Burglary`, `Violent`, and `Vehicle Theft`,
2. The distributions of `Vehicle Thefts` against the values of `Robbery`.

As usual, all plots should be labeled and easy to read.

Hint: To get the x-axis label to display, you might need to set the `sharex` parameter of `plot()` to False.

Principles of Good Data Visualization

Data visualization is a powerful tool for analysis and communication. When writing a paper or report, the author must make many decisions about how to use graphics effectively to convey useful information to the reader. Here we will go over a simple process for making deliberate, effective, and efficient design decisions.

Attention to Detail

Consider the plot in Figure 2.9. It is a scatter plot of positively correlated data of some kind, with `temp`—likely temperature—on the *x* axis and `cons` on the *y* axis. However, the picture is not really communicating anything about the dataset. It has not specified the units for the *x* or the *y* axis, nor does it tell what `cons` is. There is no title, and the source of the data is unknown.

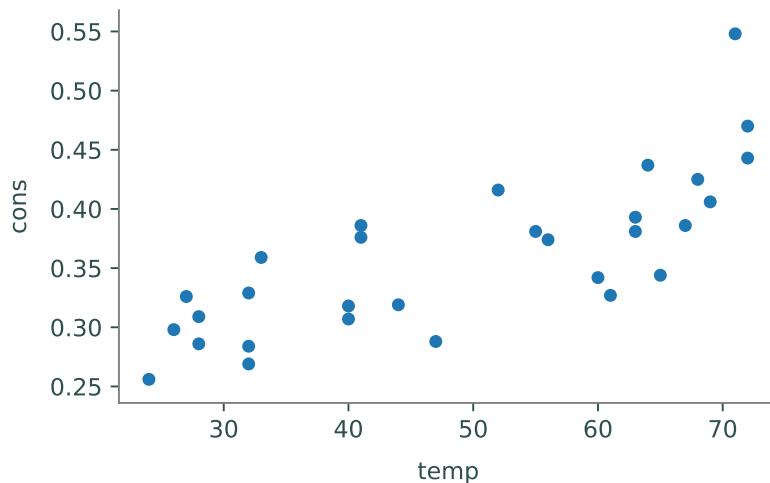


Figure 2.9: Non-specific data.

Labels and Citations

In a homework or lab setting, we sometimes (mistakenly) think that it is acceptable to leave off appropriate labels, legends, titles, and sourcing. In a published report or presentation, this kind of carelessness is confusing at best and, when the source is not included, even plagiaristic. Data needs to be explained in a useful manner that includes all of the vital information.

Consider again Figure 2.9. This figure comes from the `Icecream` dataset within the `pydataset` package, which we store here in a dataframe and then plot:

```
>>> from pydataset import data
>>> icecream = data("Icecream")
```

```
>>> icecream.plot(kind="scatter", x="temp", y="cons")
```

This code produces the rather substandard plot in Figure 2.9. Examining the source of the dataset can give important details to create better plots. When plotting data, make sure to understand what the variable names represent and where the data was taken from. Use this information to create a more effective plot.

The ice cream data used in Figure 2.9 is better understood with the following information:

1. The dataset details ice cream consumption via 30 four-week periods from March 1951 to July 1953 in the United States.
2. `cons` corresponds to “consumption of ice cream per capita” and is measured in pints.
3. `income` is the family’s weekly income in dollars.
4. `price` is the price of a pint of ice cream.
5. `temp` corresponds to temperature, degrees Fahrenheit.
6. The listed source is: “Hildreth, C. and J. Lu (1960) Demand relations with autocorrelated disturbances, Technical Bulletin No 2765, Michigan State University.”

This information gives important details that can be used in the following code. As seen in previous examples, pandas automatically generates legends when appropriate. Pandas also automatically labels the x and y axes, however our data frame column titles may be insufficient. Appropriate titles for the x and y axes must also list appropriate units. For example, the y axis should specify that the consumption is in units of pints per head, in place of the ambiguous label `cons`.

```
>>> icecream = data("Icecream")
# Set title via the title keyword argument
>>> icecream.plot(kind="scatter", x="temp", y="cons",
...     title="Ice Cream Consumption in the U.S., 1951-1953")
# Override pandas automatic labelling using xlabel and ylabel
>>> plt.xlabel("Temp (Fahrenheit)")
>>> plt.ylabel("Consumption per head (pints)")
```

To add the necessary text to the figure, use either `plt.annotate()` or `plt.text()`. Alternatively, add text immediately below wherever the figure is displayed. The first two parameters of `plt.text` are the x and y coordinates to place the text. The third parameter is the text to write. For instance, using `plt.text(0.5, 0.5, "Hello World")` will center the Hello World string in the axes.

```
>>> plt.text(20, .1, r"Source: Hildreth, C. and J. Lu (1960) \emph{Demand}
...     "relations with autocorrelated disturbances}\nTechnical Bulletin No"
...     "2765, Michigan State University.", fontsize=7)
```

Both of these methods are imperfect but can normally be easily replaced by a caption attached to the figure. Again, we reiterate how important it is that you source any data you use; failing to do so is plagiarism.

Finally, we have a clear and demonstrative graphic in Figure 2.10.

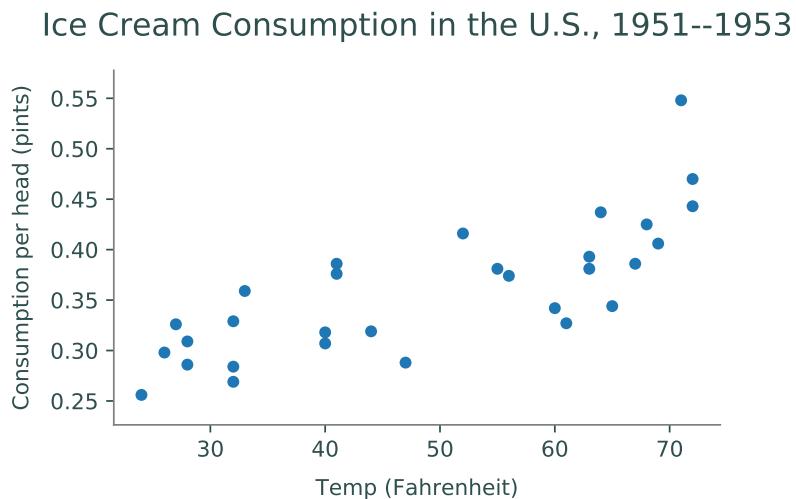


Figure 2.10: Source: Hildreth, C. and J. Lu (1960) Demand relations with autocorrelated disturbances, Technical Bulletin No 2765, Michigan State University.

Achtung!

Visualizing data can inherit many biases of the visualizer and as a result can be intentionally misleading. Examples of this include, but are not limited to, visualizing subsets of data that do not represent the whole of the data and having purposely misconstrued axes. Every data visualizer has the responsibility to avoid including biases in their visualizations to ensure data is being represented informatively and accurately.

Problem 4. The dataset `college.csv` contains information from 1995 on universities in the United States. To access information on variable names, go to <https://cran.r-project.org/web/packages/ISLR/ISLR.pdf>. Create 3 plots that compare variables or universities. These plots should answer questions about the data, e.g. what is the distribution of graduation rates or do schools with lower student to faculty ratios have higher tuition costs. These three plots should be easy to understand and have clear variable names and citations.

3

Pandas 3: Grouping

Lab Objective: Many data sets contain categorical values that naturally sort the data into groups. Analyzing and comparing such groups is an important part of data analysis. In this lab we explore pandas tools for grouping data and presenting tabular data more compactly, primarily through groupby and pivot tables.

Groupby

The file `mammal_sleep.csv`¹ contains data on the sleep cycles of different mammals, classified by order, genus, species, and diet (carnivore, herbivore, omnivore, or insectivore). The "`sleep_total`" column gives the total number of hours that each animal sleeps (on average) every 24 hours. To get an idea of how many animals sleep for how long, we start off with a histogram of the "`sleep_total`" column.

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

# Read in the data and print a few random entries.
>>> msleep = pd.read_csv("mammal_sleep.csv")
>>> msleep.sample(5)
   name   genus   vore      order  sleep_total  sleep_rem  sleep_cycle
51  Jaguar  Panthera  carni  Carnivora       10.4        NaN        NaN
77  Tenrec    Tenrec  omni  Afrosoricida     15.6        2.3        NaN
10   Goat     Capri  herbi  Artiodactyla      5.3        0.6        NaN
80   Genet    Genetta  carni  Carnivora       6.3        1.3        NaN
33  Human      Homo  omni   Primates        8.0        1.9        1.5

# Plot the distribution of the sleep_total variable.
>>> msleep.plot(kind="hist", y="sleep_total", title="Mammalian Sleep Data")
>>> plt.xlabel("Hours")
```

¹Proceedings of the National Academy of Sciences, 104 (3):1051–1056, 2007. Updates from V. M. Savage and G. B. West, with additional variables supplemented by Wikipedia. Available in `pydataset` (with a few more columns) under the key "`msleep`".

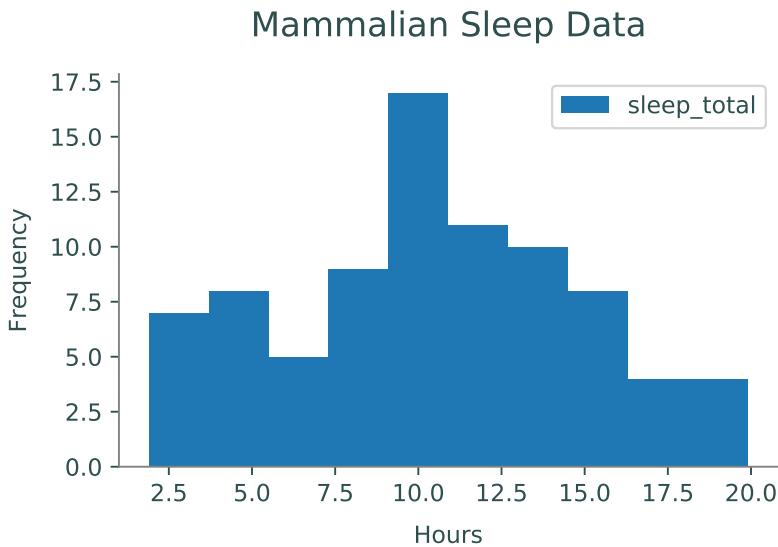


Figure 3.1: "`sleep_total`" frequencies from the mammalian sleep data set.

While this visualization is a good start, it doesn't provide any information about how different kinds of animals have different sleeping habits. How long do carnivores sleep compared to herbivores? Do mammals of the same genus have similar sleep patterns?

A powerful tool for answering these kinds of questions is the `groupby()` method of the pandas `DataFrame` class, which partitions the original `DataFrame` into groups based on the values in one or more columns. The `groupby()` method does **not** return a new `DataFrame`; it returns a pandas `GroupBy` object, an interface for analyzing the original `DataFrame` by groups.

For example, the columns "`genus`", "`vore`", and "`order`" in the mammal sleep data all have a discrete number of categorical values that could be used to group the data. Since the "`vore`" column has only a few unique values, we start by grouping the animals by diet.

```
# List all of the unique values in the 'vore' column.
>>> set(msleep["vore"])
{nan, 'herbi', 'omni', 'carni', 'insecti'}

# Group the data by the 'vore' column.
>>> vores = msleep.groupby("vore")
>>> list(vores.groups)
['carni', 'herbi', 'insecti', 'omni']           # NaN values for vore were dropped.

# Get a single group and sample a few rows. Note vore='carni' in each entry.
>>> vores.get_group("carni").sample(5)
   name    genus     vore     order  sleep_total  sleep_rem  sleep_cycle
80  Genet  Genetta    carni  Carnivora       6.3        1.3        NaN
50   Tiger  Panthera    carni  Carnivora      15.8        NaN        NaN
8     Dog     Canis    carni  Carnivora      10.1        2.9       0.333
0  Cheetah  Acinonyx    carni  Carnivora      12.1        NaN        NaN
82  Red fox    Vulpes    carni  Carnivora       9.8        2.4       0.350
```

As shown above, `groupby()` is useful for filtering a DataFrame by column values; the command `df.groupby(col).get_group(value)` returns the rows of df where the entry of the `col` column is `value`. The real advantage of `groupby()`, however, is how easily it compares groups of data. Standard DataFrame methods like `describe()`, `mean()`, `std()`, `min()`, and `max()` all work on GroupBy objects to produce a new data frame that describes the statistics of each group.

```
# Get averages of the numerical columns for each group.
>>> vores.mean()
           sleep_total   sleep_rem   sleep_cycle
vore
carni          10.379      2.290       0.373
herbi          9.509       1.367       0.418
insecti        14.940      3.525       0.161
omni           10.925      1.956       0.592

# Get more detailed statistics for 'sleep_total' by group.
>>> vores["sleep_total"].describe()
    count     mean      std    min    25%    50%    75%    max
vore
carni     19.0  10.379  4.669   2.7   6.25  10.4  13.000  19.4
herbi     32.0   9.509  4.879   1.9   4.30  10.3  14.225  16.6
insecti     5.0  14.940  5.921   8.4   8.60  18.1  19.700  19.9
omni      20.0  10.925  2.949   8.0   9.10  9.9  10.925  18.0
```

Multiple columns can be used simultaneously for grouping. In this case, the `get_group()` method of the GroupBy object requires a tuple specifying the values for each of the grouping columns.

```
>>> msleep_small = msleep.drop(["sleep_rem", "sleep_cycle"], axis=1)
>>> vores_orders = msleep_small.groupby(["vore", "order"])
>>> vores_orders.get_group(("carni", "Cetacea"))
      name      genus  vore  order  sleep_total
30    Pilot whale  Globicephalus  carni  Cetacea      2.7
59  Common porpoise      Phocoena  carni  Cetacea      5.6
79  Bottle-nosed dolphin      Tursiops  carni  Cetacea      5.2
```

Problem 1. Read in the data `college.csv` containing information on various United States universities in 1995. To access information on variable names, go to <https://cran.r-project.org/web/packages/ISLR/ISLR.pdf>. Use a `groupby` object to group the colleges by private and public universities. Read in the data as a DataFrame object and use `groupby` and `describe` to examine the following columns by group:

1. Student to faculty ratio
2. Percent of students from the top 10% of their high school class
3. Percent of students from the top 25% of their high school class

Determine whether private or public universities have a higher mean for each of these columns. For the type of university with the higher mean, save the values of the describe function on said column as an array using `.values`. Return a tuple with these arrays in the order described above.

For example, if we were comparing whether the average number of professors with PhDs was higher at private or public universities, we would find that public universities have a higher average, and we would return the following array:

```
array([212., 76.83490566, 12.31752531, 33., 71., 78.5, 86., 103.])
```

Visualizing Groups

There are a few ways that `groupby()` can simplify the process of visualizing groups of data. First of all, `groupby()` makes it easy to visualize one group at a time using the `plot` method. The following visualization improves on Figure 3.1 by grouping mammals by their diets.

```
# Plot histograms of 'sleep_total' for two separate groups.
>>> vores.get_group("carni").plot(kind="hist", y="sleep_total", legend=False,
                                 title="Carnivore Sleep Data")
>>> plt.xlabel("Hours")
>>> vores.get_group("herbi").plot(kind="hist", y="sleep_total", legend=False,
                                 title="Herbivore Sleep Data")
>>> plt.xlabel("Hours")
```

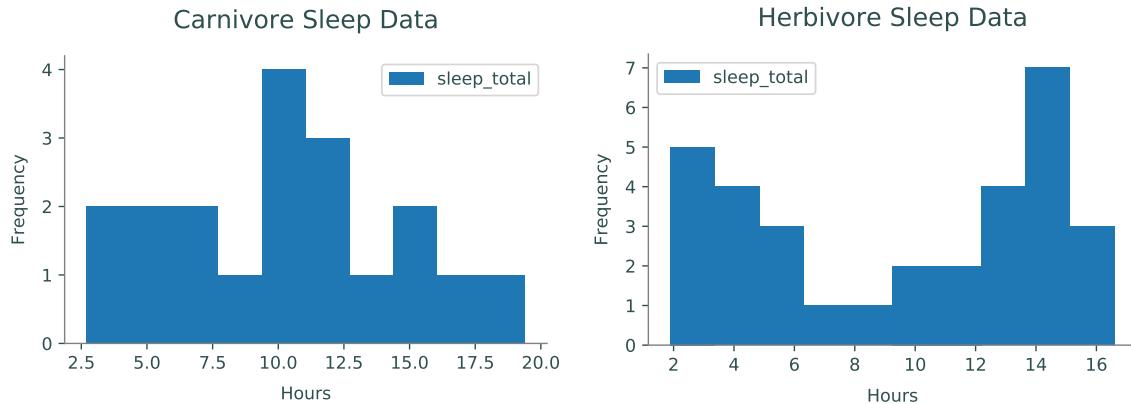
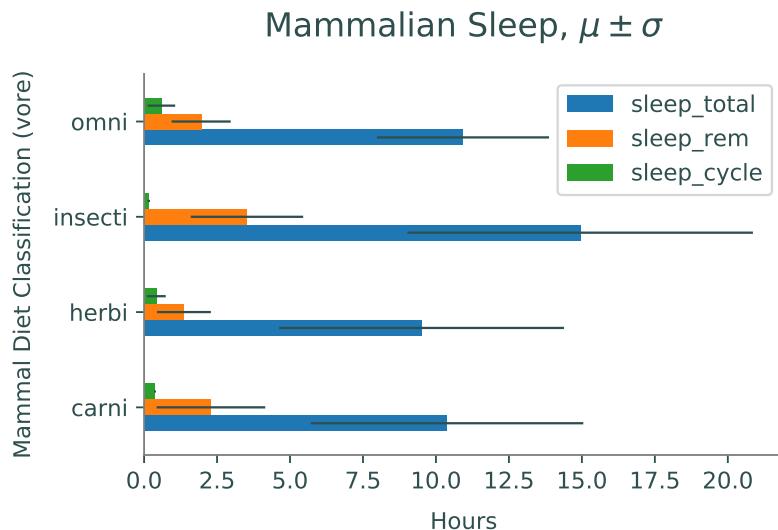


Figure 3.2: "`sleep_total`" histograms for two groups in the mammalian sleep data set.

The statistical summaries from the `GroupBy` object's `mean()`, `std()`, or `describe()` methods also lend themselves well to certain visualizations for comparing groups.

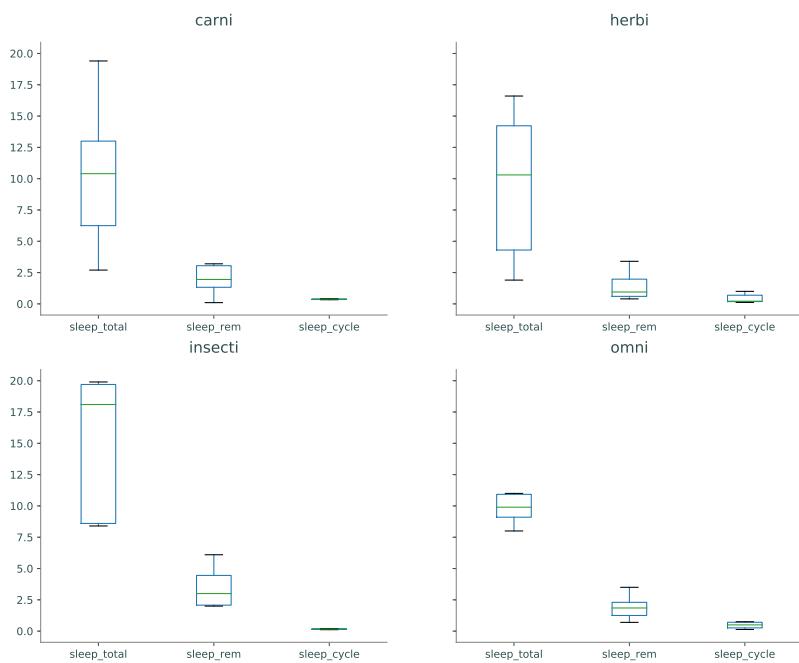
```
>>> vores[["sleep_total", "sleep_rem", "sleep_cycle"]].mean().plot(kind="barh",
                           xerr=vores.std(), title=r"Mammalian Sleep, $\mu\pm\sigma$")
>>> plt.xlabel("Hours")
```

```
>>> plt.ylabel("Mammal Diet Classification (vore)")
```



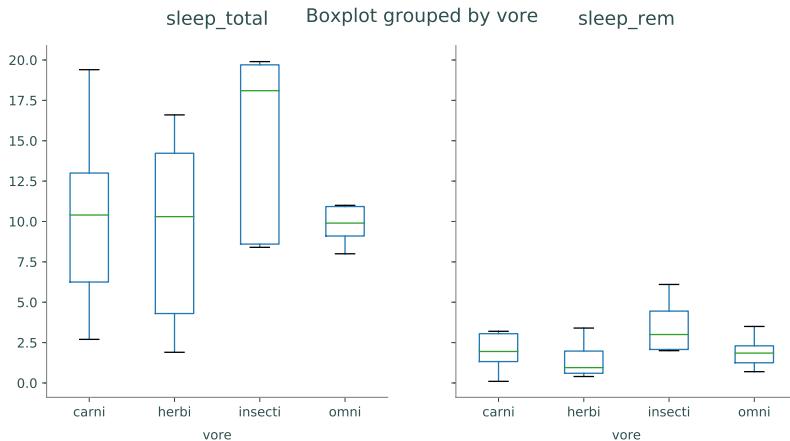
Box plots are well suited for comparing similar distributions. The `boxplot()` method of the `GroupBy` class creates one subplot **per group**, plotting each of the columns as a box plot.

```
# Use GroupBy.boxplot() to generate one box plot per group.
>>> vores.boxplot(grid=False)
>>> plt.tight_layout()
```



Alternatively, the `boxplot()` method of the `DataFrame` class creates one subplot **per column**, plotting each of the columns as a box plot. Specify the `by` keyword to group the data appropriately.

```
# Use DataFrame.boxplot() to generate one box plot per column.
>>> msleep.boxplot(["sleep_total", "sleep_rem"], by="vore", grid=False)
```



Like `groupby()`, the `by` argument can be a single column label or a list of column labels. Similar methods exist for creating histograms (`GroupBy.hist()` and `DataFrame.hist()` with `by` keyword), but generally box plots are better for comparing multiple distributions.

Problem 2. Create visualizations that give relevant information answering the following questions (using `college.csv`):

1. How do the number of applicants, number of accepted students, and number of enrolled students compare between private and public universities?
2. How does the range of money spent on room and board compare between private and public universities?

Pivot Tables

One of the downfalls of `groupby()` is that a typical `GroupBy` object has too much information to display coherently. A pivot table intelligently summarizes the results of a `groupby()` operation by aggregating the data in a specified way. The standard tool for making a pivot table is the `pivot_table()` method of the `DataFrame` class. As an example, consider the "`HairEyeColor`" data set from `pydataset`.

```
>>> from pydataset import data
>>> hec = data("HairEyeColor") # Load and preview the data.
>>> hec.sample(5)
   Hair    Eye      Sex  Freq
3    Red  Brown    Male    10
```

```

1  Black  Brown   Male    32
14 Brown  Green   Male    15
31  Red   Green   Female   7
21  Black  Blue   Female   9

>>> for col in ["Hair", "Eye", "Sex"]:
...     print("{}: {}".format(col, ", ".join(set(str(x) for x in hec[col]))))
...
Hair: Brown, Black, Blond, Red
Eye: Brown, Blue, Hazel, Green
Sex: Male, Female

```

There are several ways to group this data with `groupby()`. However, since there is only one entry per unique hair-eye-sex combination, the data can be completely presented in a pivot table.

```

>>> hec.pivot_table(values="Freq", index=["Hair", "Eye"], columns="Sex")
Sex      Female  Male
Hair  Eye
Black Blue      9    11
        Brown     36    32
        Green      2     3
        Hazel     5    10
Blond Blue     64    30
        Brown     4     3
        Green     8     8
        Hazel     5     5
Brown Blue     34    50
        Brown     66    53
        Green     14    15
        Hazel    29    25
Red   Blue      7    10
        Brown     16    10
        Green     7     7
        Hazel     7     7

```

Listing the data in this way makes it easy to locate data and compare the female and male groups. For example, it is easy to see that brown hair is more common than red hair and that about twice as many females have blond hair and blue eyes than males.

Unlike "`HairEyeColor`", many data sets have more than one entry in the data for each grouping. An example in the previous dataset would be if there were two or more rows in the original data for females with blond hair and blue eyes. To construct a pivot table, data of similar groups must be aggregated together in some way.

By default entries are aggregated by averaging the non-null values. You can use the keyword argument `aggfunc` to choose among different ways to aggregate the data. For example, if you use `aggfunc='min'`, the value displayed will be the minimum of all the values. Other arguments include `'max'`, `'std'` for standard deviation, `'sum'`, or `'count'` to count the number of occurrences. You also may pass in any function that reduces to a single float, like `np.argmax` or even `np.linalg.norm` if you wish. A list of functions can also be passed into the `aggfunc` keyword argument.

Consider the Titanic data set found in `titanic.csv`². For this analysis, take only the "Survived", "Pclass", "Sex", "Age", "Fare", and "Embarked" columns, replace null age values with the average age, then drop any rows that are missing data. To begin, we examine the average survival rate grouped by sex and passenger class.

```
>>> titanic = pd.read_csv("titanic.csv")
>>> titanic = titanic[["Survived", "Pclass", "Sex", "Age", "Fare", "Embarked"]]
>>> titanic["Age"].fillna(titanic["Age"].mean(),)

>>> titanic.pivot_table(values="Survived", index="Sex", columns="Pclass")
Pclass      1.0      2.0      3.0
Sex
female    0.965   0.887   0.491
male      0.341   0.146   0.152
```

Note

The `pivot_table()` method is a convenient way of performing a potentially complicated `groupby()` operation with aggregation and some reshaping. The following code is equivalent to the previous example.

```
>>> titanic.groupby(["Sex", "Pclass"])["Survived"].mean().unstack()
Pclass      1.0      2.0      3.0
Sex
female    0.965   0.887   0.491
male      0.341   0.146   0.152
```

The `stack()`, `unstack()`, and `pivot()` methods provide more advanced shaping options.

Among other things, this pivot table clearly shows how much more likely females were to survive than males. To see how many entries fall into each category, or how many survived in each category, aggregate by counting or summing instead of taking the mean.

```
# See how many entries are in each category.
>>> titanic.pivot_table(values="Survived", index="Sex", columns="Pclass",
...                      aggfunc="count")
Pclass  1.0  2.0  3.0
Sex
female  144  106  216
male    179  171  493

# See how many people from each category survived.
>>> titanic.pivot_table(values="Survived", index="Sex", columns="Pclass",
...                      aggfunc="sum")
Pclass      1.0      2.0      3.0
```

²There is a "Titanic" data set in `pydataset`, but it does not contain as much information as the data in `titanic.csv`.

```
Sex
female 137.0 94.0 106.0
male    61.0 25.0  75.0
```

Problem 3. The file `Ohio_1999.csv` contains data on workers in Ohio in the year 1999. Use pivot tables to answer the following questions:

1. Which race.sex combination has the highest `Usual Weekly Earnings` in total?
2. Which race.sex combination has the lowest cumulative `Usual Hours Worked`?
3. What race.sex combination has the highest average `Usual Hours Worked`?

Return a tuple for each question (in order of the questions) where the first entry is the numerical code corresponding to the race and the second entry is corresponding to the sex.

Some useful keys in understand the data are as follows:

1. In column `Sex`, {1: `male`, 2: `female`}.
2. In column `Race`, {1: `White`, 2: `African-American`, 3: `Native American/Eskimo`, 4: `Asian`}.

Discretizing Continuous Data

In the Titanic data, we examined survival rates based on sex and passenger class. Another factor that could have played into survival is age. Were male children as likely to die as females in general? We can investigate this question by multi-indexing, or pivoting, on more than just two variables, by adding in another index.

In the original dataset, the `"Age"` column has a floating point value for the age of each passenger. If we add `"Age"` as another pivot, then the table would create a new row for `each` age present. Instead, we partition the `"Age"` column into intervals with `pd.cut()`, thus creating a categorical that can be used for grouping. Notice that when creating the pivot table, the index uses the categorical age instead of the column name `"Age"`.

```
# pd.cut() maps continuous entries to discrete intervals.
>>> pd.cut([1, 2, 3, 4, 5, 6, 7], [0, 4, 8])
[(0, 4], (0, 4], (0, 4], (0,4], (0, 4], (4, 8], (4, 8], (4, 8]]
Categories (2, interval[int64]): [(0, 4] < (4, 8]]

# Partition the passengers into 3 categories based on age.
>>> age = pd.cut(titanic['Age'], [0, 12, 18, 80])

>>> titanic.pivot_table(values="Survived", index=["Sex", age],
                        columns="Pclass", aggfunc="mean")
Pclass          1.0      2.0      3.0
Sex   Age
female (0, 12]  0.000  1.000  0.467
```

	(12, 18]	1.000	0.875	0.607
	(18, 80]	0.969	0.871	0.475
male	(0, 12]	1.000	1.000	0.343
	(12, 18]	0.500	0.000	0.081
	(18, 80]	0.322	0.093	0.143

From this table, it appears that male children (ages 0 to 12) in the 1st and 2nd class were very likely to survive, whereas those in 3rd class were much less likely to. This clarifies the claim that males were less likely to survive than females. However, there are a few oddities in this table: zero percent of the female children in 1st class survived, and zero percent of teenage males in second class survived. To further investigate, count the number of entries in each group.

```
>>> titanic.pivot_table(values="Survived", index=["Sex", "age"],
                        columns="Pclass", aggfunc="count")
Pclass      1.0  2.0  3.0
Sex   Age
female (0, 12]    1   13   30
          (12, 18]   12    8   28
          (18, 80]  129   85  158
male   (0, 12]    4   11   35
          (12, 18]   4   10   37
          (18, 80]  171  150  420
```

This table shows that there was only 1 female child in first class and only 10 male teenagers in second class, which sheds light on the previous table.

Achtung!

The previous pivot table brings up an important point about partitioning datasets. The Titanic dataset includes data for about 1300 passengers, which is a somewhat reasonable sample size, but half of the groupings include less than 30 entries, which is **not** a healthy sample size for statistical analysis. Always carefully question the numbers from pivot tables before making any conclusions.

Pandas also supports multi-indexing on the columns. As an example, consider the price of a passenger tickets. This is another continuous feature that can be discretized with `pd.cut()`. Instead, we use `pd.qcut()` to split the prices into 2 equal quantiles. Some of the resulting groups are empty; to improve readability, specify `fill_value` as the empty string or a dash.

```
# pd.qcut() partitions entries into equally populated intervals.
>>> pd.qcut([1, 2, 5, 6, 8, 3], 2)
[(0.999, 4.0], (0.999, 4.0], (4.0, 8.0], (4.0, 8.0], (4.0, 8.0], (0.999, 4.0]]
Categories (2, interval[float64]): [(0.999, 4.0] < (4.0, 8.0]]

# Cut the ticket price into two intervals (cheap vs expensive).
>>> fare = pd.qcut(titanic["Fare"], 2)
>>> titanic.pivot_table(values="Survived",
                        index=["Sex", "age"], columns=[fare, "Pclass"],
```

		aggfunc="count", fill_value='-')		
Fare	(-0.001, 14.454]	(14.454, 512.329]		
Pclass		1.0	2.0	3.0
Sex	Age			
female	(0, 12]	-	-	7
	(12, 18]	-	4	23
	(18, 80]	-	31	101
male	(0, 12]	-	-	8
	(12, 18]	-	5	26
	(18, 80]	8	94	350
				1 13 23
				12 4 5
				129 54 57
				4 11 27
				4 5 11
				163 56 70

Not surprisingly, most of the cheap tickets went to passengers in 3rd class.

Problem 4. Use the employment data from Ohio in 1999 to answer the following questions:

1. The column `Educational Attainment` contains numbers 0-46. Any number less than 39 means the person did not get any form of degree. 39-42 refers to either a high-school or associate's degree. A number greater than or equal to 43 means the person got at least a bachelor's degree. Out of these categories, which degree type is the most common among the workers in this dataset?
2. Partition the `Age` column into 6 equally-sized groups using `pd.qcut()`. Which interval has the highest average `Usual Hours Worked`?
3. Using the partitions from the first two parts, what age/degree combination has the lowest yearly salary on average?

Return the answer to each question (in order) as an `Interval`. For part three, the answer should be a tuple where the first entry in the `Interval` of the age and the second is the `Interval` of the degree.

An `Interval` is the object returned by `pd.cut()` and `pd.qcut()`. These can also be obtained from a pivot table, as in the example below.

```
>>> # Create pivot table used in last example with titanic dataset
>>> table = titanic.pivot_table(values="Survived",
                                 index=[age], columns=[fare, "Pclass"],
                                 aggfunc="count")
>>> # Get index of maximum interval
>>> table.sum(axis=1).idxmax()
Interval(0, 12, closed='right')
```

Problem 5. Examine the college dataset using `pivot tables` and `groupby` objects. Determine the answer to the following questions. If the answer is yes, save the answer as `True`. If the answer is no, save the answer as `False`. For the last question, save the answer as a string giving your explanation. Return a tuple containing your answers to the questions in order.

1. Is there a correlation between the percent of alumni that donate and the amount the school spends per student in BOTH private and public universities?
2. Partition **Grad.Rate** into evenly spaced intervals of 20%. Is the partition with the greatest number of schools the same for private and public universities?
3. Does having a lower acceptance rate correlate with having more students from the top 10 percent of their high school class being admitted on average for BOTH private and public universities?
4. Why is the average percentage of students admitted from the top 10 percent of their high school class so high in private universities with very low acceptance rates? Use only the data to explain why; do not extrapolate.

4

Information Theory and Wordle

Lab Objective: Use the information theory concept of entropy to create an algorithm for playing the popular word game Wordle.

Wordle

Wordle is a popular word game¹ where you have 6 guesses to guess a five-letter word. Every time a guess is made, you receive some information about how close your guess is to the correct answer. Letters in the guess that are in the correct location are colored green; letters that are present in the secret word but not in the correct location are colored yellow; and letters that aren't present in the secret word are colored gray. An example game is given in Figure ??.



Figure 4.1: An example game of Wordle.

The secret word is chosen at random from a fixed list of 2309 words. While it is possible to only select guesses from these words, it is not necessarily the best strategy; in many situations, there is a word that can be guessed that gives more information about what the secret word is than guessing any possible secret word would. Additionally, there is a list of 12953 words that are allowed to be used as guesses; the guess we make cannot be any arbitrary string of 5 characters, but must always must be one of these words.

¹Specifically, it went viral on the internet in 2022.

There are a few technicalities with how the guess is evaluated when there are duplicates of a letter. If the secret word were “speak” and a guess of “bevel” was made, the first e would be colored yellow and the second gray. Letters that are marked green take priority over this; if “bevel” is guessed and the secret word were “ashes” instead, the first e in the guess will be gray and the second green. With the same guess, if the secret word were “steel”, then the first e would be yellow and the second green. So, if there are more of a given letter in the guess than in the secret word, only as many will be marked yellow or green in the guess as there are in the secret word.

Problem 1. Write a function that accepts the secret word and a guess, and returns the colors of the guess as an array. Label correct letters with the number 2, letters in the wrong location with 1, and incorrect letters with 0. For instance, with the secret word "`pages`" and the guess "`green`", your function should return `array([1,0,0,2,0])`.

Hint: Find some way to keep track of which letters in the secret word have been matched to. Since strings are immutable, it may also be helpful in this case to cast the guess and secret word into arrays if you need to modify them.

Problem 2. In order to efficiently implement our strategy for Wordle, we need to know what the result of each guess is for every possible secret word. We only need to make this computation once for each pair, so we will do them all at once and store them in an array.

Load the lists of possible secret words and allowed guesses from `possible_words.txt` and `allowed_words.txt`. Write a function `get_all_guess_results()` that finds the result of making a guess for each pair of secret word and allowed guess. Store the results in a 3-dimensional numpy array, where the first axis corresponds to the guess, the second to the secret word, and the third to the letter.

This computation on the whole set of words will take several minutes at the very least, so test your function on a smaller subset of the word lists. Using the first three secret words and the first two possible guesses, your function should output the following:

```
>>> get_all_guess_results(possible_words[:3], allowed_words[:2])
array([[ [2, 1, 0, 0, 0],
         [2, 1, 0, 1, 0],
         [2, 1, 0, 1, 0]],

        [[2, 1, 0, 0, 0],
         [2, 1, 0, 0, 0],
         [2, 1, 0, 0, 0]]])
```

Compute the array for the full word lists. Use `np.save` to save the array to a file, to avoid needing to recompute it.

Note

Three-dimensional numpy arrays behave similarly to two-dimensional ones, and can be accessed and sliced in the same way. The only difference is that indexing only a single axis will give a two-dimensional array. We illustrate this by showing how to access some useful subsets of the final array. Here, `all_guess_results` is the output of Problem 2, `i=1388` is the index of a given guess (“boxes”), and `j=1914` is the index of a given secret word (“steel”).

```
# This 2-D array is the result of the i-th guess for every secret word
>>> all_guess_results[i,:]
array([[1, 0, 0, 0, 0],
       [1, 0, 0, 1, 1],
       [1, 0, 0, 1, 0],
       ...,
       [1, 0, 0, 1, 0],
       [0, 0, 0, 1, 1],
       [0, 2, 0, 0, 0]])

# This is equivalent to all_guess_results[i] and all_guess_results[i,:,:]

# This 2-D array is the result of every guess for the j-th secret word
>>> all_guess_results[:,j]
array([[0, 0, 0, 2, 0],
       [0, 0, 1, 0, 0],
       [0, 0, 0, 0, 0],
       ...,
       [0, 0, 0, 0, 0],
       [0, 0, 0, 2, 1],
       [0, 0, 0, 0, 0]])

# This 1-D array is the result of the i-th guess on the j-th secret word
>>> all_guess_results[i,j]
array([0, 0, 0, 2, 1])
```

Achtung!

We will use this array frequently in the next few problems and modify it in several ways; however, be sure to keep the original array, as it will still be needed. Remember that arrays are mutable, so do not do modifications in-place on the original array. If you accidentally modify the original array, reload it from your file with `np.load`.

Our objective is to create some strategy to play Wordle as effectively as possible. Simply choosing the word that is most likely to be the secret word is completely ineffective, as there is no reason to prefer any word over another, as long as both are consistent with the information we have. A much better strategy is to maximize the amount of information each of our guesses gives us, which we will quantify by using entropy.

Information and Entropy

Entropy is the expected amount of information we would gain by knowing the result of a variable. A natural way to define the information of an event A is as $-\log_2 P(A)$.² The entropy of a random variable X , which we denote $H(X)$, is then defined as

$$H(X) = \mathbb{E}[-\log_2 P(X = x)] = -\sum_x P(X = x) \log_2 P(X = x),$$

where the sum version comes from the Law of the Unconscious Statistician. A loose interpretation is that if a random variable has lower entropy, then we know more about what its value will be even if we haven't observed it yet, and observing it usually will give little information. At one extreme, if a discrete random variable has zero entropy, then it is in fact necessarily constant. On the other hand, if a random variable has higher entropy, then we know less about its result and observing it typically will give more information.

For Wordle, since we don't know the secret word, it is reasonable to consider it as a random variable; this gives the secret word a value of entropy, which can be used to choose a guess that is likely to give more information. Denote the secret word as W and the result of making a guess g as $R(g)$; since we don't know the secret word, this is also a random variable. There are two approaches we can take to make a strategy out of this. First, we can think of trying to minimize the entropy of the variable $W|R(g)$. This essentially is trying to find the guess that will on average minimize how much we don't know about the secret word after we know the result of the guess. Second, we can think of trying to maximize the entropy of the variable $R(g)$. This amounts to finding which guess is expected to give the most information.

These two approaches are in fact equivalent, as

$$H(W|R(g)) = H((W, R(g))) - H(R(g)) = H(W) - H(R(g)),$$

where $H((W, R(g)))$ denotes the entropy of the joint random variable $(W, R(g))$ (not the cross entropy). To see this equality, note that for random variables X, Y we have

$$-\log_2 P(X|Y) = -\log_2 \frac{P(X, Y)}{P(Y)} = -\log_2 P(X, Y) + \log_2 P(Y);$$

taking the expectation of both sides implies that

$$H(X|Y) = H((X, Y)) - H(Y).$$

Additionally, the value of $R(g)$ is completely determined by W , so $H((W, R(g))) = H(W)$.

The result of all of this is that minimizing the entropy of $W|R(g)$ is equivalent to maximizing the entropy of $R(g)$, which we will use to select a good guess. Since the entropy of $R(g)$ is more straightforward to calculate, this is the approach we take for the remainder of the lab.

We now seek to calculate the entropy of $R(g)$, the result of the guess, for each guess g we can make. This is given by

$$\begin{aligned} H(R(g)) &= -\sum_r P(R(g) = r) \log_2 P(R(g) = r) \\ &= -\sum_r P(R(g, W) = r) \log_2 P(R(g, W) = r). \end{aligned}$$

²This choice of definition has a number of desirable properties for information: information is non-negative, the information of two independent events is the sum of their individual informations, and information is a continuous function of the probability of an event. In fact, it can be shown that such a function is the negative logarithm of the probability of an event, for some logarithm base; refer to the Volume 3 textbook for more details. The base-2 logarithm is commonly used because it can be thought of as representing the number of bits needed to encode the information.

Since we assumed a uniform distribution over the set of possible secret words, the probability $P(R(g, W) = r)$ can be calculated simply as the proportion of secret words that yield the same result r given the same guess g . We already computed the result of making each acceptable guess with each possible secret word in Problem 2. So, to find the entropy of a guess, we need only to compute the probability of each unique guess result, and then apply the equation above. This sum will need to be evaluated for each individual guess that we can make.

As an example, suppose that we know the secret word is one of the words “boney”, “disco”, “marsh”, “stock”, or “visor”, and we are evaluating the guess “boxes”. The result of this guess for each of these words is as follows:

Word	Guess result
boney	(2,2,0,2,0)
disco	(0,1,0,0,1)
marsh	(0,0,0,0,1)
stock	(0,1,0,0,1)
visor	(0,1,0,0,1)

There are three distinct possible results: $(2,2,0,2,0)$, with probability $\frac{1}{5}$; $(0,1,0,0,1)$, with probability $\frac{3}{5}$; and $(0,0,0,0,1)$, with probability $\frac{1}{5}$. Using the above formula gives the entropy of this guess as

$$-\frac{1}{5} \log_2 \frac{1}{5} - \frac{3}{5} \log_2 \frac{3}{5} - \frac{1}{5} \log_2 \frac{1}{5} \approx 1.3710$$

Problem 3. Write a function that accepts the multidimensional array created in Problem 2 and calculates the entropy of each guess. Return the guess with the highest entropy. Also return the index of this guess in the list of allowed guesses, to avoid needing to find it later.

In order to simplify determining the numbers of each unique guess result, we can first condense the result of each guess into a single number. A simple way to do this is interpreting the five numbers of the result as a ternary (base 3) number. For instance, we can convert the array $[1, 0, 2, 2, 1]$ to the number $1 \cdot 3^0 + 0 \cdot 3^1 + 2 \cdot 3^2 + 2 \cdot 3^3 + 1 \cdot 3^4 = 154$. This step is also simple to vectorize, which makes it a relatively quick operation. Applying this to the whole array from Problem 2 will result in a 2-dimensional array, whose axes correspond to the possible secret words and the allowed guesses. Be sure not to overwrite the original array.

Hint: `np.unique` with the argument `return_counts=True` will return an array with the number of occurrences of each of the different values in a one-dimensional array. By looping over the allowed guesses, you can use this function to compute the entropy quickly. Applying this function directly to multidimensional arrays results in different behavior, however.

After we make a guess, we would like to compute the effect of knowing the result on the probability distribution of the secret word. Bayes’ Rule gives

$$P(W = w|R(g) = r) = \frac{P(R(g) = r|W = w)P(W = w)}{P(R(g) = r)}.$$

First, we look at the term $P(R(g) = r|W = w)$. If we know the secret word W , then for any guess g , the result $R(g)$ is uniquely determined. Thus, this probability is either 0 or 1, depending on whether the guess result we observed is the result that would be seen if w is the secret word. For instance, with the secret word $w = \text{“steel”}$ and the guess $g = \text{“boxes”}$, the only value of r for which the probability is not zero is $r = [0, 0, 0, 2, 1]$. This is precisely the result of making that guess for that word.

Now, also note that $P(W = w)$ is a constant, and $P(R(g) = r)$ is constant for all secret words that have $P(R(g) = r|W = w) \neq 0$, since these all have the same value of $R(g)$. So, the posterior distribution is just a uniform distribution over the set of words that give the same result for our guess as we observed. Finding the optimal next guess to make is then equivalent to repeating the same process as before with a smaller initial list of possible secret words.

Problem 4. Create a function that filters the list of possible words after making a guess. Since we already computed the result of all guesses for all possible words in Problem 2, we will use this array instead of recomputing the results. Accept this array, the list of possible words, a guessed word's index in the list of allowed guesses, and the guess's result. Return a filtered list of possible words that are still possible after knowing the result of a guess. Also return a filtered version of the array of all guess results that only contains the results for the secret words still possible after making the guess. This smaller array will be used to compute the entropies for the following guess.

Hint: The most efficient way to do this problem is with boolean masking. If A is any numpy array and `mask` is a 1-D array of True/False values, then $A[mask]$ will return the portion of A where `mask` is true. This can be used even if A is multidimensional, and on dimensions other than the first; for instance, $A[:, mask]$ will use the mask for the second dimension of the array.

Note

Note that, while we filter down the list of possible secret words, we do not do anything similar for the list of allowed guesses. The reason for this is that, as the game goes on and we make more guesses, the list of words that could still be the secret word shrinks, while the list of words we are allowed to guess stays the same. In general, it is beneficial to allow ourselves to guess words that cannot be the secret word, because in some cases we will get more information that way.

Before we assemble our algorithm for playing Wordle, we would like a benchmark. A simple strategy to compare to is to select an allowed guess at random until we know the secret word.

Problem 5. The file `wordle.py` contains a class called `WordleGame` object that can be used to simulate games of Wordle.^a Instantiate one of these, use the `start_game()` function to start a game, and use the `make_guess()` function to make a guess.

Write a function that accepts a `WordleGame` and starts and plays a game using the strategy of randomly selecting words. At each step:

- If we know the word, guess it; otherwise, choose a guess at random from the list of allowed guesses.
- Filter the list of possible words to only those that are still possible; this allows us to determine if we know what the secret word is
- Repeat until the secret word has been guessed

Return the number of guesses needed to guess the secret word. To visualize this algorithm, pass the argument `display=True`, and the `WordleGame` will print out each word as it is guessed.

^aThis class uses the `colorama` package to format terminal output. This package is included with the Anaconda distribution of python, but can easily be installed with `pip` if needed.

Problem 6. Write a function that accepts a `WordleGame` object and starts and plays a game using the strategy of maximizing the entropy of each guess. At each step:

- If we know the secret word, guess that word
- Otherwise, compute the entropies, and make the guess that has the highest entropy
- Filter the possible words to only those that are possible after the guess
- Repeat until the secret word has been guessed

Return the number of guesses needed to guess the secret word.

Problem 7. Write a function that accepts an integer n and simulates that many games of Wordle using each of the above algorithms. Return the average number of guesses each required to find the secret word. Compare their performance; the approach using the entropy should require about half as many guesses on average.

The `WordleGame` object also has a version you can play in the terminal, which can be started using the `play_game_interactive()` method. You can use this to also compare your own performance to that of your algorithm.

5

GeoPandas

Lab Objective: GeoPandas is a package designed to organize and manipulate geographic data. It combines the data manipulation tools of pandas with the geometric capabilities of the Shapely package. In this lab, we explore the basic data structures of GeoSeries and GeoDataFrames and their functionalities.

Installation

GeoPandas is a new package designed to combine the functionality of pandas with Shapely, a package used for geometric manipulation. Using GeoPandas with geographic data is very useful as it allows the user to not only compare numerical data, but also geometric attributes. GeoPandas can be installed via pip:

```
>>> pip install geopandas
```

However, Geopandas can be notoriously difficult to install. This is especially the case if the python environment is not carefully maintained. Some of its dependencies can also be very difficult to install on certain systems. Because of this, using Colab for this lab is recommended; its environment is set up in a way that makes GeoPandas very easy to install. Otherwise, the GeoPandas documentation contains some additional options that can be used if installation difficulties occur: <https://geopandas.org/install.html>.

GeoSeries

A GeoSeries is a pandas Series where each entry is a set of geometric objects. There are three classes of geometric objects inherited from the Shapely package:

1. Points / Multi-Points
2. Lines / Multi-Lines
3. Polygons / Multi-Polygons

A point is used to identify objects like coordinates, where there is one small instance of the object. A line could be used to describe objects such as roads. A polygon could be used to identify regions, such as a country. Multipoints, multilines, and multipolygons contain lists of points, lines, and polygons, respectively.

Since each object in the GeoSeries is also a Shapely object, the GeoSeries inherits many methods and attributes of Shapely objects. Some of the key attributes and methods are listed in Table 5.1. These attributes and methods can be used to calculate distances, find the sizes of countries, and determine whether coordinates are within country's boundaries. The example below uses the attribute `bounds` to find the maximum and minimum coordinates of Egypt in a built-in GeoDataFrame.

Method/Attribute	Description
<code>distance(other)</code>	returns minimum distance from GeoSeries to <code>other</code>
<code>contains(other)</code>	returns <code>True</code> if shape contains <code>other</code>
<code>intersects(other)</code>	returns <code>True</code> if shape intersects <code>other</code>
<code>area</code>	returns shape area
<code>convex_hull</code>	returns convex shape around all points in the object
<code>bounds</code>	returns the bounding x- and y-coordinates of the object

Table 5.1: Attributes and Methods for GeoSeries

```
>>> import geopandas as gpd
>>> world = gpd.read_file(gpd.datasets.get_path('naturalearth_lowres'))
# Get GeoSeries for Egypt
>>> egypt = world[world['name']=='Egypt']

# Find bounds of Egypt
>>> egypt.bounds
      minx     miny         maxx     maxy
47  24.70007   22.0    36.86623  31.58568
```

Creating GeoDataFrames

The main structure used in GeoPandas is a GeoDataFrame, which is similar to a pandas DataFrame. A GeoDataFrame has one special column called `geometry`, which must be a GeoSeries. This GeoSeries column is used when a spatial method, like `distance()`, is used on the GeoDataFrame. Therefore all attributes and methods used for GeoSeries can also be used on GeoDataFrame objects.

A GeoDataFrame can be made from a pandas DataFrame. At least one of the columns in the DataFrame should contain geometric information. This column containing geometric information can be converted to a GeoSeries using the `apply()` method. At this point, the Pandas DataFrame can be cast as a GeoDataFrame. Assign which column will be the `geometry` using either the `geometry` keyword in the constructor or the `set_geometry()` method afterwards.

```
>>> import pandas as pd
>>> import geopandas as gpd
>>> from shapely.geometry import Point, Polygon
```

```

# Create a Pandas DataFrame
>>> df = pd.DataFrame({'City': ['Seoul', 'Lima', 'Johannesburg'],
...                      'Country': ['South Korea', 'Peru', 'South Africa'],
...                      'Latitude': [37.57, -12.05, -26.20],
...                      'Longitude': [126.98, -77.04, 28.04]})

# Create geometry column
>>> df['Coordinates'] = list(zip(df.Longitude, df.Latitude))

# Make geometry column Shapely objects
>>> df['Coordinates'] = df['Coordinates'].apply(Point)

# Cast as GeoDataFrame
>>> gdf = gpd.GeoDataFrame(df, geometry='Coordinates')

# Equivalently, specify the geometry after construction
# Note that set_geometry() returns a new GeoDataFrame
>>> gdf = gpd.GeoDataFrame(df)
>>> gdf = gdf.set_geometry('Coordinates')

# Display the GeoDataFrame
>>> gdf
      City      Country  Latitude  Longitude           Coordinates
0    Seoul  South Korea    37.57   126.98  POINT (126.98000 37.57000)
1      Lima        Peru   -12.05   -77.04  POINT (-77.04000 -12.05000)
2  Johannesburg  South Africa   -26.20    28.04  POINT (28.04000 -26.20000)

# Create a polygon with all three cities as points
>>> city_polygon = Polygon(list(zip(df.Longitude, df.Latitude)))

```

A GeoDataFrame can also be made directly from a dictionary. If the dictionary already contains geometric objects, the corresponding column can be directly set as the `geometry` in the constructor. Otherwise, a column containing geometry data can be created as in the above example and then set as the `geometry` with the `set_geometry()` method.

```

# Both of these methods create the same GeoDataFrame as above
# Directly create the GeoDataFrame from the dictionary
>>> gdf = gpd.GeoDataFrame({'City': ['Seoul', 'Lima', 'Johannesburg'],
...                           'Country': ['South Korea', 'Peru', 'South Africa'],
...                           'Latitude': [37.57, -12.05, -26.20],
...                           'Longitude': [126.98, -77.04, 28.04]})

# Create geometry column and set as the geometry
>>> gdf['Coordinates'] = list(zip(gdf.Longitude, gdf.Latitude))
>>> gdf['Coordinates'] = gdf['Coordinates'].apply(Point)
# inplace=True modifies gdf itself rather than returning a copy
>>> gdf.set_geometry('Coordinates', inplace=True)

# Equivalently, using a dictionary that already contains geometry objects
>>> gdf = gpd.GeoDataFrame({'City': ['Seoul', 'Lima', 'Johannesburg'],
...
```

```
...           'Country': ['South Korea', 'Peru', 'South Africa'],
...
...           'Coordinates': [Point(126.98,37.57),
...                             Point(-77.04,-12.05), Point(28.04,-12.05)]},
...
...           geometry='Coordinates')
```

Method/Attribute	Description
<code>abs()</code>	returns series/dataframe with absolute numeric value of each element
<code>add(other)</code>	returns addition of dataframe and <code>other</code> element-wise
<code>affine_transform(matrix)</code>	returns <code>GeoSeries</code> with translated geometries
<code>append(other)</code>	returns new object with appended rows of <code>other</code> to the end of caller
<code>dot(other)</code>	returns dataframe of matrix multiplication with <code>other</code>
<code>equals(other)</code>	tests if the two objects contain the same elements

Table 5.2: Attributes and Methods for GeoDataFrame

Note

Longitude is the angular measurement starting at the Prime Meridian, 0° , and going to 180° to the east and -180° to the west. Latitude is the angle between the equatorial plane and the normal line at a given point; a point along the Equator has latitude 0, the North Pole has latitude $+90^\circ$ or $90^\circ N$, and the South Pole has latitude -90° or $90^\circ S$.

Plotting GeoDataFrames

Information from a GeoDataFrame is plotted based on the geometry column. Data points are displayed as geometry objects. The following example plots the shapes in the `world` GeoDataFrame.

```
# Plot world GeoDataFrame
>>> world.plot()
```

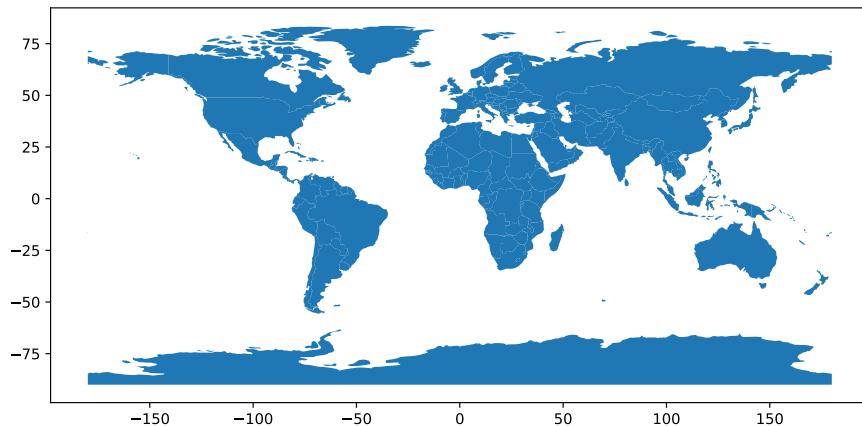


Figure 5.1: World map

Multiple GeoDataFrames can be plotted at once. This can be done by setting one GeoDataFrame as the base of the plot and ensuring that each layer uses the same axes. In the following example, the file `airports.csv`, containing the coordinates of world airports, is loaded into a GeoDataFrame and plotted on top of the boundary of the `world` GeoDataFrame.

```
# Set outline of world countries as base
>>> fig,ax = plt.subplots(figsize=(10,7), ncols=1, nrows=1)
>>> base = world.boundary.plot(edgecolor='black', ax=ax, linewidth=1)

# Load airport data and convert to a GeoDataFrame
>>> airports = pd.read_csv('airports.csv')
>>> airports['Coordinates'] = list(zip(airports.Longitude, airports.Latitude))
>>> airports['Coordinates'] = airports.Coordinates.apply(Point)
>>> airports = gpd.GeoDataFrame(airports, geometry='Coordinates')

# Plot airports on top of world map
>>> airports.plot(ax=base, marker='o', color='green', markersize=1)
>>> ax.set_xlabel('Longitude')
>>> ax.set_ylabel('Latitude')
>>> ax.set_title('World Airports')
```

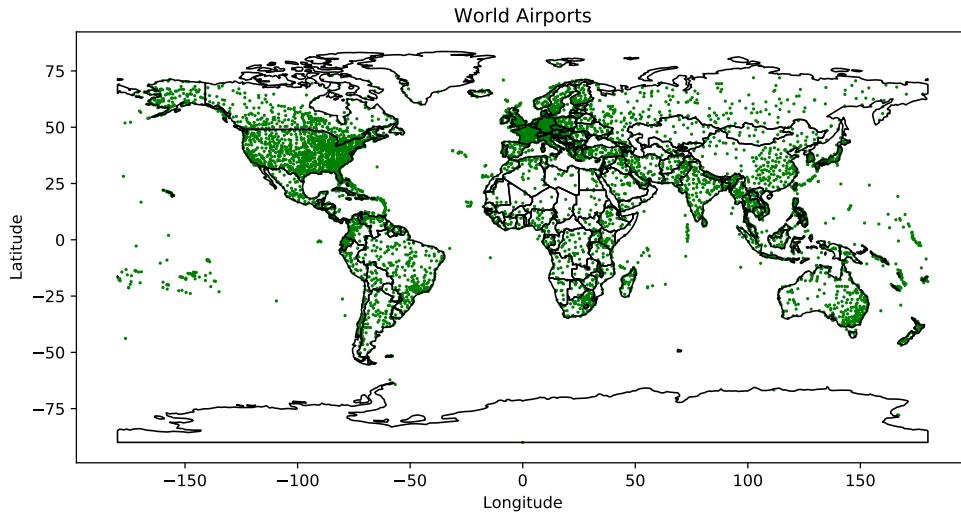


Figure 5.2: Airport map

Problem 1. Read in the file `airports.csv` as a pandas DataFrame. Create three convex hulls around the three sets of airports listed below. This can be done by passing in lists of the airports' coordinates to a `shapely.geometry.Polygon` object.

Create a new GeoDataFrame with these three Polygons as entries. Plot this GeoDataFrame on top of an outlined world map.

- Maio Airport, Scatsta Airport, Stokmarknes Skagen Airport, Bekily Airport, K. D. Matanzima Airport, RAF Ascension Island
- Oiapoque Airport, Maio Airport, Zhezkazgan Airport, Walton Airport, RAF Ascension Island, Usiminas Airport, Piloto Osvaldo Marques Dias Airport
- Zhezkazgan Airport, Khanty Mansiysk Airport, Novy Urengoy Airport, Kalay Airport, Biju Patnaik Airport, Walton Airport

Working with GeoDataFrames

As previously mentioned, GeoDataFrames contain many of the functionalities of pandas DataFrames. For example, to create a new column, define a new column name in the GeoDataFrame with the needed information for each GeoSeries.

```
# Create column in the world GeoDataFrame for gdp_per_capita
>>> world['gdp_per_cap'] = world.gdp_md_est / world.pop_est
```

GeoDataFrames can utilize many pandas functionalities, and they can also be parsed by geometric manipulations. For example, a useful way to index GeoDataFrames is with the `cx` indexer. This splits the GeoDataFrame by the coordinates of each geometric object. It is used by calling the method `cx` on a GeoDataFrame, followed by a slicing argument, where the first element refers to the longitude and the second refers to latitude.

```
# Create a GeoDataFrame containing the northern hemisphere
>>> north = world.cx[:, 0:]

# Create a GeoDataFrame containing the southeastern hemisphere
>>> south_east = world.cx[0:, :0]
```

GeoSeries objects in a GeoDataFrame can also be dissolved, or merged, together into one GeoSeries based on their geometry data. For example, all countries on one continent could be merged to create a GeoSeries containing the information of that continent. The method designed for this is called `dissolve`. It receives two parameters, `by` and `aggfunc`. `by` indicates which column to dissolve along, and `aggfunc` tells how to combine the information in all other columns. The default `aggfunc` is `first`, which returns the first application entry. In the following example, we use `sum` as the `aggfunc` so that each continent is the combination of its countries.

```
>>> world = world[['continent', 'geometry', 'gdp_per_cap']]

# Dissolve world GeoDataFrame by continent
>>> continent = world.dissolve(by = 'continent', aggfunc='sum')
```

Projections and Coloring

When plotting, GeoPandas uses the CRS (coordinate reference system) of a GeoDataFrame. This reference system indicates how coordinates should be spaced on a plot. Two of the most commonly used CRSs are EPSG:4326 and EPSG:3395. EPSG:4326 is the standard latitude-longitude projection used by GPS. EPSG:3395, also known as the Mercator projection, is the standard navigational projection.

When creating a new GeoDataFrame, it is important to set the `crs` attribute of the GeoDataFrame. This allows any plots to be shown correctly. Furthermore, GeoDataFrames being layered need to have the same CRS. To change the CRS, use the method `to_crs()`.

```
# Check CRS of world GeoDataFrame
>>> print(world.crs)
epsg:4326

# Change CRS of world to Mercator
# inplace=True ensures that we modify world instead of returning a copy
>>> world.to_crs(3395, inplace=True)
>>> print(world.crs)
epsg:3395
```

GeoPandas accepts many different CRSs; a reference can be found at www.spatialreference.org. Additionally, inspecting a given CRS object in the terminal without using `print()` or `str()` can be used to get additional information about a specific CRS:¹

```
>>> world.crs
```

¹This can also be accomplished using `print(repr(crs))`.

```
<Projected CRS: EPSG:3395>
Name: WGS 84 / World Mercator
Axis Info [cartesian]:
- E[east]: Easting (metre)
- N[north]: Northing (metre)
Area of Use:
- name: World between 80°S and 84°N.
- bounds: (-180.0, -80.0, 180.0, 84.0)
Coordinate Operation:
- name: World Mercator
- method: Mercator (variant A)
Datum: World Geodetic System 1984
- Ellipsoid: WGS 84
- Prime Meridian: Greenwich
```

GeoDataFrames can also be plotted using the values in the other attributes of the GeoSeries. The map plots the color of each geometry object according to the value of the column selected. This is done by passing in the parameter `column` into the `plot()` method.

```
>>> fig, ax = plt.subplots(1, figsize=(10,4))
# Plot world based on gdp
>>> world.plot(column='gdp_md_est', cmap='OrRd', legend=True, ax=ax)
>>> ax.set_title('World Map based on GDP')
>>> ax.set_xlabel('Longitude')
>>> ax.set_ylabel('Latitude')
>>> plt.show()
```

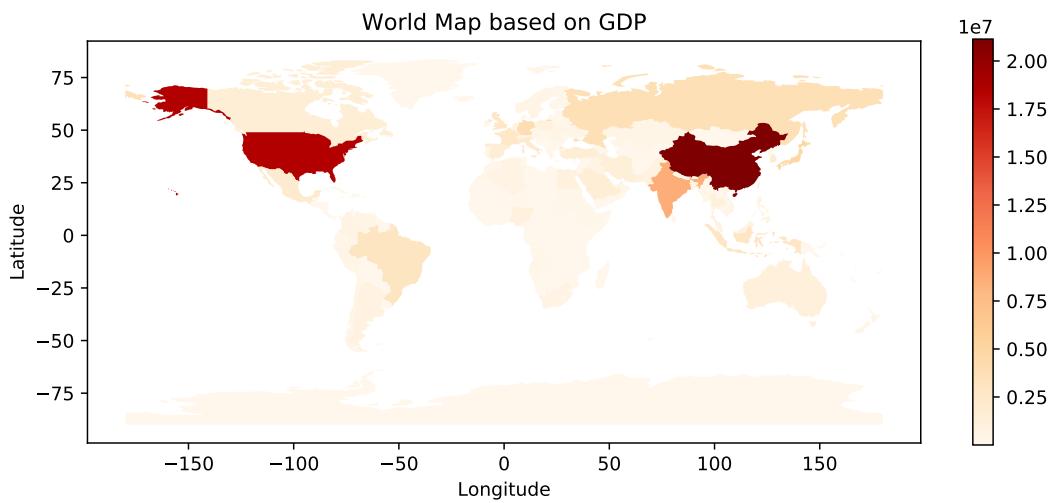


Figure 5.3: World Map Based on GDP

Problem 2. The file `county_data.gpkg.zip` contains information about US counties.^a After unzipping, use the command `gpd.read_file('county_data.gpkg')` to create a GeoDataFrame of this information. Each county's shape is stored in the `geometry` column. Use this to plot all US counties two times, first using the default CRS and then using EPSG:5071.

Next, create a new GeoDataFrame that merges all counties within a single state. Drop regions with the following STATEFP codes: 02, 15, 60, 66, 69, 72, 78. Plot this GeoDataFrame to see an outline of all 48 contiguous states. Ensure a CRS of EPSG:5071.

^aSource: http://www2.census.gov/geo/tiger/GENZ2016/shp/cb_2016_us_county_5m.zip

Note

.gpkg files are actually structured as a directory that contains several files that each contain parts of the data. For instance, `county_data.gpkg` consists of the files `county_data.cpg`, `county_data.dbf`, `county_data.prj`, `county_data.shp`, and `county_data.shx`. Be sure that these files are placed directly in the first level of folders, and not in further subdirectories.

To use this file in Google Colab, upload the zipped file and extract it with the following code:

```
county = files.upload()
!unzip county_data.gpkg.zip
```

It then can be loaded:

```
county_df = gpd.read_file('county_data.gpkg')
```

Merging GeoDataFrames

Just as multiple pandas DataFrames can be merged, multiple GeoDataFrames can be merged with attribute joins or spatial joins. An attribute join is similar to a merge in pandas. It combines two GeoDataFrames on a column (not the geometry column) and then combines the rest of the data into one GeoDataFrame.

```
>>> world = gpd.read_file(geopandas.datasets.get_path('naturalearth_lowres'))
>>> cities = gpd.read_file(geopandas.datasets.get_path('naturalearth_cities'))

# Create subsets of the world and cities GeoDataFrames
>>> world = world[['continent', 'name', 'iso_a3']]
>>> cities = cities[['name', 'iso_a3']]

# Merge the GeoDataFrames on their iso_a3 code
>>> countries = world.merge(cities, on='iso_a3')
```

A spatial join merges two GeoDataFrames based on their geometry data. The function used for this is `sjoin`. `sjoin` accepts two GeoDataFrames and then direction on how to merge. It is imperative that two GeoDataFrames have the same CRS. In the example below, we merge using an `inner` join with the option `intersects`. The `inner` join means that we will only use keys in the intersection of both geometry columns, and we will retain only the left geometry column. `intersects` tells the GeoDataFrames to merge on GeoSeries that intersect each other. Other options include `contains` and `within`.

```
# Combine countries and cities on their geographic location
>>> countries = gpd.sjoin(world, cities, how='inner', op='intersects')
```

Problem 3. Load in the file `nytimes.csv`^a as a DataFrame. This file includes county-level data for the cumulative cases and deaths of Covid-19 in the US, starting with the first case in Snohomish County, Washington, on January 21, 2020. Begin by converting the `date` column into a `DatetimeIndex`.

Next, use county FIPS codes to merge your GeoDataFrame from Problem 2 with the DataFrame you just created. A FIPS code is a 5-digit unique identifier for geographic locations. Ignore rows in the Covid-19 DataFrame with unknown FIPS codes as well as all data from Hawaii and Alaska.

Note that the `fips` column of the Covid-19 DataFrame stores entries as floats, but the county GeoDataFrame stores FIPS codes as strings, with the first two digits in the `STATEFP` column and the last three in the `COUNTYFP` column.

Once you have completed the merge, plot the cases from March 21, 2020 on top of your state outline map from Problem 2, using the CRS of EPSG:5071. Finally, print out the name of the county with the most cases on March 21, 2020 along with its case count.

^aSource: <https://raw.githubusercontent.com/nytimes/covid-19-data/master/us-counties.csv>

Logarithmic Plotting Techniques

The color scheme of a graph can also help to communicate information clearly. A good list of available colormaps can be found at https://matplotlib.org/3.2.1/gallery/color/color_map_reference.html. Note also that you can reverse any colormap by adding `_r` to the end. The following example demonstrates some plotting features, using country GDP as in Figure 5.3.

```
>>> fig, ax = plt.subplots(figsize=(15,7), ncols=1, nrows=1)
>>> world.plot(column='gdp_md_est', cmap='plasma_r',
...             ax=ax, legend=True, edgecolor='gray')

# Add title and remove axis tick marks
>>> ax.set_title('GDP on Linear Scale')
>>> ax.set_yticks([])
>>> ax.set_xticks([])
>>> plt.show()
```

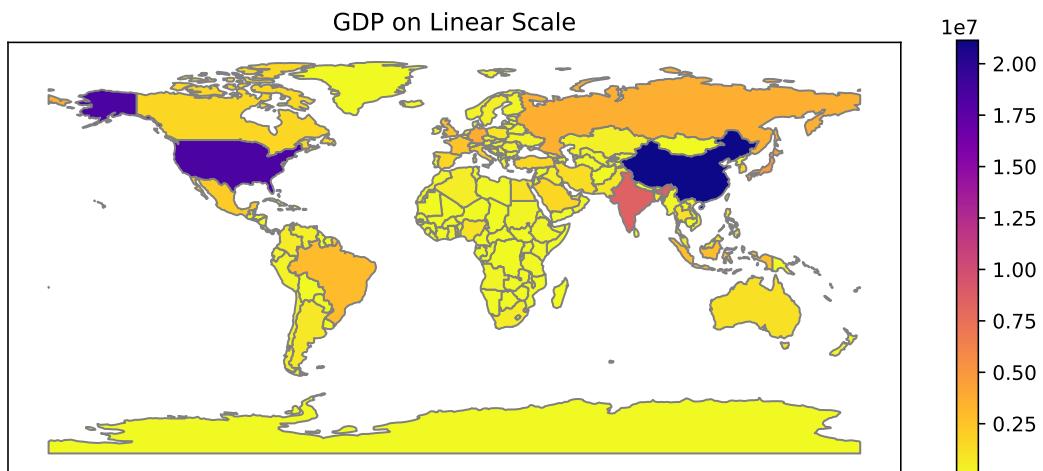


Figure 5.4: World map showing country GDP

Sometimes data can be much more informative when plotted on a logarithmic scale. See how the world map changes when we add a `norm` argument in the code below. Depending on the purpose of the graph, Figure 5.5 may be more informative than Figure 5.4.

```
>>> from matplotlib.colors import LogNorm
>>> from matplotlib.cm import ScalarMappable
>>> fig, ax = plt.subplots(figsize=(15,6), ncols=1, nrows=1)

# Set the norm using data bounds
>>> data = world.gdp_md_est
>>> norm = LogNorm(vmin=min(data), vmax=max(data))

# Plot the graph using the norm
>>> world.plot(column='gdp_md_est', cmap='plasma_r', ax=ax,
...             edgecolor='gray', norm=norm)

# Create a custom colorbar
>>> cbar = fig.colorbar(ScalarMappable(norm=norm, cmap='plasma_r'),
...                      ax=ax, orientation='horizontal', pad=0, label='GDP')

>>> ax.set_title('Country Area on a Log Scale')
>>> ax.set_yticks([])
>>> ax.set_xticks([])
>>> plt.show()
```

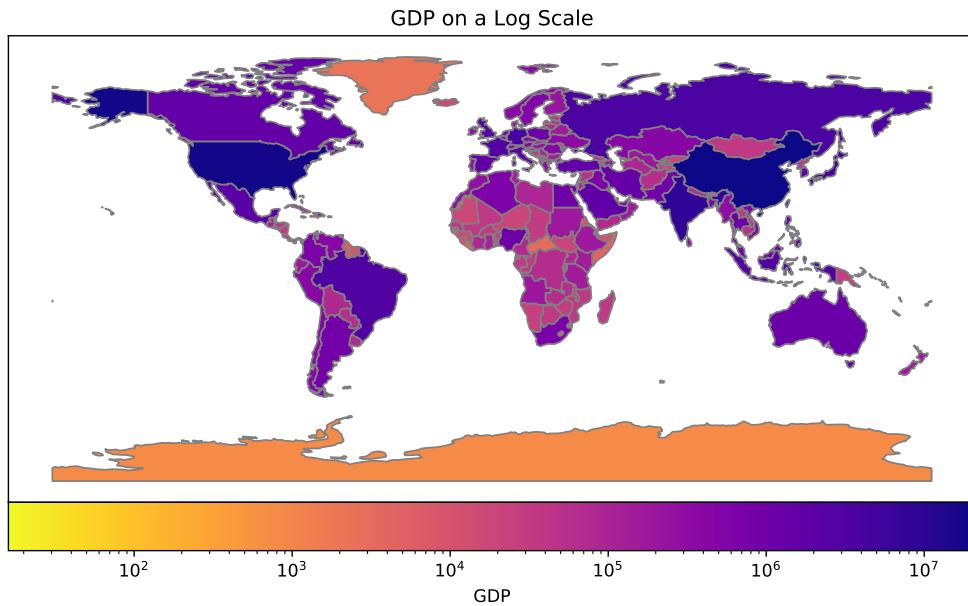


Figure 5.5: World map showing country GDP using a log scale

Problem 4. As in Problem 3, plot your state outline map from Problem 2 on top of a map of the Covid-19 cases from March 21, 2020. This time, however, use a log scale. Use EPSG:5071 for the CRS. Pick a good colormap (the counties with the most cases should generally be darkest) and be sure to display a colorbar.

Problem 5. In this problem, you will create an animation of the spread of Covid-19 through US counties from January 21, 2020 to June 21, 2020. Use a log scale and a good colormap, and be sure that you're using the same norm and colorbar for the whole animation. Use EPSG:5071 for the projection.

As a reminder, below is a summary of what you will need in order to animate this map. You may also find it helpful to refer to the animation section included with the Volume 4 lab manual.

1. Set up your figure and norm. Be sure to use the highest case count for your `vmax` so that the scale remains uniform.
2. Write your `update` function. This should plot the cases from a given day.
3. Set up your colorbar. Do this outside the `update` function to avoid adding a new colorbar each day.
4. Create the animation and embed it.

6

Data Cleaning

Lab Objective: The quality of a data analysis or model is limited by the quality of the data used. In this lab we learn techniques for cleaning data, creating features, and determining feature importance.

Almost every dataset has problems that make it unsuitable for regression or other modeling. Some problems will throw errors in your code and so will be easily observed. Other data problems are less noticeable. If code executes on poorly maintained data, the results could vary significantly from the true results which might have been obtained if the underlying dataset were better prepared.

Data cleaning is the process of identifying and correcting bad data. This could be data that is missing, duplicated, irrelevant, inconsistent, incorrect, in the wrong format, or otherwise does not make sense. Though it can be tedious, data cleaning is the most important step of data analysis. Without accurate and legitimate data, any results or conclusions are suspect and may be incorrect.

We will demonstrate common issues with data and how to correct them using the following dataset. It consists of family members and some basic details.

```
# Example dataset
>>> df = pd.read_csv('toy_dataset.csv')

>>>
      Name   Age      name        DOB Marital_Status
0    John Doe   30     john  01/01/2010      Divorcee
1  Jane Doe    29     jane  12/02/1990      Divorced
2  Jill smith   40      NaN  03/04/1980      married
3  Jill smith   40     jill  03/04/1980      married
4  jack smith  100     jack  4/4/1980  marrieed
5 Jenny Smith    5      NaN  05/05/2015       NaN
6 JAMES Smith    2      NaN  20/06/2018      single
7      Rover    2      NaN  05/05/2018       NaN

      Height  Weight  Marriage_Len      Spouse
0    72.0     175          5       NaN
1     5.5     125          5  John Doe
```

2	64.0	120		10	Jack Smith
3	64.0	120		NaN	jack smith
4	1.8	220		10	jill smith
5	105.0	40		NaN	NaN
6	27.0	25	Not Applicable		NaN
7	36.0	50		NaN	NaN

Inspection

The first step of data cleaning is to analyze the quality of the data. If the quality is poor, the data might not be worth using. Knowing the quality of the data will also give you an idea of how long it will take to clean it. A quality dataset is one in which the data is valid, accurate, complete, consistent, and uniform. Some of these issues, like uniformity, are fairly easy to fix during cleaning, while other aspects like accuracy are more difficult, if not impossible, to address.

Validity is the degree that the data conforms to given rules. If a column corresponds to the temperature in Salt Lake City, measured in degrees Fahrenheit, then a value over 110 or below 0 should make you suspicious, since those would be extreme values for Salt Lake City. In fact, checking the all-time temperature records for Salt Lake shows that the values in this column should never be more than 107 and never less than –30. Any values outside that range are almost certainly errors and should probably be reset to *NaN*, unless you have special information that allows you to impute more accurate values.

Some standard considerations when determining the validity of a dataset are:

- **data type:** The data types of each column should all be the same.
- **data range:** The data of a column, typically numbers or dates, should all be within some valid range.
- **mandatory constraints:** Certain columns cannot have missing entries.
- **unique constraint:** Entries in certain columns must be unique.
- **regular expression patterns:** A text column must be in the same format (for example, forcing phone numbers to be formatted as 999-999-9999).
- **cross-field validation:** Conditions must hold across multiple columns (for example, a hospital discharge date can't be earlier than the admittance date).
- **duplicated data:** Rows or columns that are repeated. In some cases, they may not be exact.

We can check the data type in Pandas using `dtype`. A `dtype` of `object` means that the data in that column contains either strings or mixed `dtypes`. These fields should be investigated to determine if they contain mixed datatypes. In our toy example, we would expect that `Marriage_Len` is numerical, so an object `dtype` is suspicious. Looking at the data, we see that `James` has `Not Applicable`, which is a string.

```
# Check validity of data
# Check Data Types
>>> df.dtypes
Name          object
```

```

Age           int64
name          object
DOB           object
Marital_Status object
Height         float64
Weight          int64
Marriage_Len   object
Spouse          object
dtype: object

```

Duplicates

Duplicates can be easily identified in Pandas using the `duplicated()` function. When no parameters are passed, it returns a DataFrame of the first duplicates. We can identify rows that are duplicated in only some columns by passing in the column names. The `keep` parameter has three possible values, first, last, and False. False keeps all duplicated values, while first and last keep only the first and last instances, respectively.

```

# Display duplicated rows
>>> df[df.duplicated()]
Empty DataFrame
Columns: [Name, Age, name, DOB, Marital_Status, Height, Weight, Marriage_Len, ←
          Spouse]
Index: []

# Display rows that have duplicates in some columns
>>> df[df.duplicated(['Name','DOB','Marital_Status'],keep=False)]
      Name  Age  name        DOB  Marital_Status  Height  Weight  ←
      Marriage_Len  Spouse
2  Jill  smith    40  03/04/1980      married    64.0     120  ←
   10  Jack  Smith
3  Jill  smith    40   jill  03/04/1980      married    64.0     120  ←
   NaN  jack  smith

```

Range

We can check the range of values in a numeric column using the `min` and `max` attributes. Other options for looking at the values include line plots, histograms, and boxplots. Some other useful Pandas commands for evaluating the breadth of a dataset include `df.nunique()` (which returns a series giving the name of each column and the number of unique values in each column), `pd.unique()` (which returns an array of the unique values in a series), and `value_counts()` (which counts the number of instances of each unique value in a column, like a histogram).

```

# Count the number of unique values in each row
>>> df.nunique()
Name      7
Age       6

```

```

name          2
DOB           7
Marital_Status 5
Height         7
Weight          7
Marriage_Len    4
Spouse          4
dtype: int64

# Print the unique Marital_Status values
>>> pd.unique(df['Marital_Status'])
array(['Divorcee', 'Divorced', 'married', 'married', nan, 'single'],
      dtype=object)

# Count the number of each Marital_Status values
>>> df['Marital_Status'].value_counts()
married      2
single       1
married     1
Divorcee     1
Divorced     1
Name: Marital_Status, dtype: int64

```

Accuracy

The accuracy of the data, how close the data is to reality, is harder to confirm. Just because a data point is valid, doesn't mean that it is true. For example, a valid street address doesn't have to exist, or a person might lie about their weight. The first case could be checked using mapping software, but the second could be unverifiable.

Missing Data

The percentage of missing data is the completeness of the data. All uncleaned data will have missing values, but datasets with large amounts of missing data, or lots of missing data in key columns, are not going to be as useful. Pandas has several functions to help identify and count missing values. In Pandas, all missing data is considered a `NaN` and does not affect the `dtype` of a column. `df.isna()` returns a boolean DataFrame indicating whether each value is missing. `df.notnull()` returns a boolean DataFrame with `True` where a value is not missing.

```

# Count number of missing data points in each column
>>> df.isna().sum()
Name          0
Age           0
name          6
DOB           0
Marital_Status 2

```

```

Height          0
Weight          0
Marriage_Len    2
Spouse          4
dtype: int64

```

Consistency

Consistency measures how cohesive the data is, both within the dataset and across multiple datasets. For example, in our toy dataset `Jack Smith` is 100 years old, but his birth year is 1980. Data is inconsistent across datasets when the data points should be the same and are different. This could be due to incorrect entries or syntax errors.

Uniformity

Lastly, uniformity is the measure of how similarly the data is formatted. Data that has the same units of measure and syntax are considered uniform. Looking at the `Height` column in our dataset, we see values ranging from 1.8 to 105. This is likely the result of different units of measure.

Uniformity also matters across multiple datasets. For example, if you use multiple finance datasets to build a predictive model then the dates in each dataset should have the same format so that they can all be used equally in the model.

No Set Rules

When looking at the quality of the data, there are no set rules on how to measure these concepts or at what point the data is considered bad data. Sometimes the only data available is of poor quality and so must still be used. Other times, higher quality data may be available elsewhere. An idea of the quality of the data will inform you of which cleaning steps are needed and will influence the strength of your final analysis.

You should always investigate the quality of the dataset they wish to use because a model is only as good as the data it relies on. Such an investigation should include statistics summarizing the principles discussed in this section, visualizations to identify outliers in the data, and written descriptions of the mitigating steps taken to improve the data set. Using various data visualizations can also give one a general sense of the quality of their data. Using histograms, box plots, and hexbins can identify outliers in the data. Outliers should be investigated to determine if they are accurate. Removing outliers will improve your model, but you should only remove an outlier if you have a legitimate reason. Columns that have a small distribution or variance, or consist of one value, could be worth removing since they might contribute little to the model.

Problem 1. The `g_t_results.csv` file is a set of parent-reported scores on their child's Gifted and Talented tests. The two tests, OLSAT and NNAT, are used by NYC to determine if children are qualified for gifted programs. The OLSAT Verbal has 16 questions for Kindergartners and 30 questions for first, second, and third graders. The NNAT has 48 questions. Using this dataset, answer the following questions.

1. What column has the highest number of null values and what percent of its values are null? Print the answer as a tuple with (column name, percentage). Make sure the second value is a percent.
2. List the columns that should be numeric that aren't. Print the answer as a tuple.
3. How many third graders have scores outside the valid range for the OLSAT Verbal Score? Print the answer
4. How many data values are missing (NaN)? Print the number.
Each part is one point.

Cleaning

After the data has been inspected, it's time to start cleaning. There are many aspects and methods of cleaning; not all of them will be used in every dataset. Which ones you choose should be based on your dataset and the goal of the project.

Unwanted Data

Removing unwanted data typically falls into two categories, duplicated data and irrelevant data. Duplicated observations usually occur when data is scraped, combined from multiple datasets, or submitted twice by a user. Irrelevant data consists of observations that don't fit the specific problem you are trying to solve or don't have enough variation to affect the model. We can drop duplicated data using the `duplicated()` function described above with `drop()` or using `drop_duplicates`, which has the same parameters as `duplicated`.

Validity Errors

After removing unwanted data, we correct any validity errors found during inspection. All features should have a consistent type, standard formatting (like capitalization), and the same units. Syntax errors should be fixed, and white space at the beginning and ends of strings should be removed. Some data might need to be padded so that it's all the same length.

Method	Description
<code>series.str.lower()</code>	Convert to all lower case
<code>series.str.upper()</code>	Convert to all upper case
<code>series.str.strip()</code>	Remove all leading and trailing white space
<code>series.str.lstrip()</code>	Remove leading white space
<code>series.str.replace(" ", "")</code>	Remove all spaces
<code>series.str.pad()</code>	Pad strings

Table 6.1: Pandas String Formatting Methods

Validity also includes correcting or removing contradicting values. This might be two values in a row or values across datasets. For example, a child shouldn't have a marital status of married. Or, if two columns should sum to a third but don't, then your data has invalid values which may need to be removed.

Missing Data

There will always be missing data in any uncleaned dataset. Some commonly suggested methods for handling data are removing the missing data and setting the missing values to some value based on other observations. However, missing data can be informative and removing or replacing missing data erases that information. Removing missing values from a dataset might result in losing significant amounts of data or even in a less accurate model. Retaining the missing values can help increase accuracy.

We have several options to deal with missing data:

- Dropping missing data is the easiest method. Dropping rows should only be done if there are a small number of missing data points in a column or if the row is missing a significant amount of data. If a column is very sparse, consider dropping the entire column. If dropping missing data is inappropriate, you may instead choose to estimate the missing values. There are many ways to do this including mean, mode, median, randomly choosing from a distribution, linear regression, and hot-decking, to name a few.
- Hot-decking is when you fill in the data based on similar observations. It can be applied to numerical and categorical data, unlike most of the other options listed above. Sequential hot-decking sorts the column with missing data based on an auxiliary column and then fills in the data with the value from the next available data point. K-Nearest Neighbors can also be used to identify similar data points.
- The last option is to flag the data as missing. This retains the information from missing data and removes the missing data (by replacing it). For categorical data, simply replace the data with a new category. For numerical data, we can fill the missing data with 0, or some value that makes sense, and add an indicator variable for missing data.

```
## Replace missing data
import numpy as np

# Add an indicator column based on missing Marriage_Len
>>> df['missing_ML'] = df['Marriage_Len'].isna()

# Fill in all missing data with 0
>>> df['Marriage_Len'] = df['Marriage_Len'].fillna(0)

# Change all other NaNs to missing
>>> df = df.fillna('missing')

# Change Not Applicable row to NaNs
>>> df = df.replace('Not Applicable',np.nan)

# Drop rows with NaNs
>>> df = df.dropna()

>>> df
      Name   Age          DOB Marital_Status
0    John   30  1985-07-15        Single
1    Jane   25  1988-01-01       Married
2    Tom   45  1970-05-20        Single
3    Alice  22  1995-03-12       Married
4    Bob   38  1975-09-10        Single
5    Carol  50  1965-12-25       Married
6    David  28  1989-07-14        Single
7    Emma  20  1998-04-01       Married
8    Frank  42  1972-02-12        Single
9    Grace  35  1984-09-17       Married
10   Helen  55  1960-06-10        Single
11   Ivan  32  1986-11-05       Married
12   Julia  25  1993-08-10        Single
13   Kevin  40  1978-03-15       Married
14   Linda  38  1981-07-20        Single
15   Michael  28  1990-01-01       Married
16   Nancy  45  1973-05-12        Single
17   Oliver  35  1987-09-18       Married
18   Paul  52  1968-02-14        Single
19   Quinn  22  1996-06-01       Married
20   Robert  48  1976-12-22        Single
21   Sarah  30  1985-07-15       Married
22   Thomas  38  1978-01-01        Single
23   Victoria  25  1994-04-10       Married
24   William  50  1965-12-25        Single
```

0	JOHN DOE	30	01/01/2010	divorcee
1	JANE DOE	29	12/02/1990	divorced
2	JILL SMITH	40	03/04/1980	married
3	JACK SMITH	40	4/4/1980	married
4	JENNY SMITH	5	05/05/2015	missing
Height	Weight	Marriage_Len	Spouse	missing_ML
72.0	175	5	missing	False
68.0	125	5	John Doe	False
64.0	120	10	Jack Smith	False
71.0	220	10	jill smith	False
41.0	40	0	missing	True

Nonnumerical Values Misencoded as Numbers

Recording data as a numerical data type (`float` or `int`) when no numerical meaning applies to the situation causes errors which can be extremely difficult to debug. Some data should be recorded in data types that cannot be multiplied or summed.

Missing data should always be stored in a form that cannot accidentally be incorporated into the model. Typically this is done by storing missing values as *Nan*. However, some algorithms will not run on data with *Nan* values, in which case you may choose to fill missing data with a string '`missing`'. Unfortunately, many datasets have recorded missing values with a 0 or some other number. You should verify that this does not occur in your dataset. Similarly, a survey with a scale from 1 to 5 will sometimes have the additional choice of "N/A" (meaning "not applicable"), which could be coded as 6, not because the value 6 is meaningful, but just because that is the next thing after 5. Again, this should be fixed so that the "N/A" choice cannot accidentally be used for any computations.

Categorical data are also often encoded as numerical values. These values should not be left as numbers that can be computed with. For example, postal codes are shorthand for locations, and there is no numerical meaning to the code. It makes no sense to add, subtract, or multiply postal codes, so it is important not to let those accidentally be added, subtracted, or multiplied, for example by inadvertently including them in the design matrix (unless they are one-hot encoded or given some other meaningful numerical value). It is good practice to convert postal codes, area codes, ID numbers, and other non-numeric data into strings or other data types that cannot be computed with.

Ordinal Data

Ordinal data is data that has a meaningful order but the differences between the values aren't consistent, or maybe aren't even meaningful at all. For example, a survey question might ask about your level of education, with 1 being high-school graduate, 2 bachelor's degree, 3 master's degree, and 4 doctoral degree. These values are called ordinal data because it is meaningful to talk about an answer of 1 being less than an answer of 2. However, the difference between 1 and 2 is not necessarily the same as the difference between 3 and 4, and it would not make sense to compute an average answer—the average of a high school diploma and a masters degree is not a bachelor's degree, despite the fact that the average of 1 and 3 is 2. Treating these like categorical data loses the information of the ordering, but treating it like regular numerical data implies that a difference of 2 has the same meaning whether it comes as $3 - 1$ or $4 - 2$. If that difference of 2 has approximately the same meaning, then it may be ok to treat these data as numerical in your model, but if that assumption is not correct then it may be better to treat the variable as categorical.

Problem 2. `imdb.csv` contains a small set of information about 99 movies. Clean the data set by doing the following in order:

1. Remove duplicate rows by dropping the first or last. Print the shape of the dataframe after removing the rows.
2. Drop all rows that contain missing data. Print the shape of the dataframe after removing the rows.
3. Remove rows that have data outside valid data ranges and explain briefly how you determined your ranges for each column.
4. Identify and drop columns with three or fewer different values. Print a tuple with the names of the columns dropped.
5. Convert the titles to all lower case.

Print the first five rows of your dataframe.

Feature Engineering

One often needs to construct new columns, commonly referred to as `features` in the context of machines learning, for a dataset, because the dependent variable is not necessarily a linear function of the features in the original dataset. Constructing new features is called feature engineering. Once new features are created, we can analyze how much a model depends on each feature. Features with low importance probably do not contribute much and could potentially be removed.

Fognets are fine mesh nets that collect water that condenses on the netting. These are used in some desert cities in Morocco to produce drinking water. Consider a dataset measuring the amount of water Y collected from fognets, where one of the features `WindDir` is the wind direction, measured in degrees. This feature is not likely to contribute meaningfully in a linear model because the direction 359 is almost the same as the direction 0, but no nonzero linear multiple of `WindDir` will reflect this relation. One way to improve the situation is to replace the `WindDir` with two new (engineered) features: $\sin(\frac{\pi}{180}\text{WindDir})$ and $\cos(\frac{\pi}{180}\text{WindDir})$.

Discrete Fourier transforms and wavelet decomposition often reveal important properties of data collected over time (called time-series), like sound, video, economic indicators, etc. In many such settings it is useful to engineer new features from a wavelet decomposition, the DFT, or some other function of the data.

Problem 3. `basketball.csv` contains data for all NBA players between 2001 and 2018. Each row represents a player's stats for a year. The features in this data set are

- player (str): the player's name
- age (int): the player's age
- team_id (cat): the player's team
- per (float): player efficiency rating, how much a player produced in one minute of play
- ws (float): win shares, an estimate of how much the player contributed to
- bpm (float): box plus/minus is the estimated number of points a player contributed to over 100 possessions
- year (int): the year

(float):

Create two new features:

- career_length (int): number of years player has been playing (start at 0).
- target (str): The target team if the player is leaving. If the player is retiring, the target should be 'retires'. A player is retiring if their name doesn't exist the next year. (Set the players in 2019 to NaN).

Remove all duplicates of a player in each year. Remove all rows except those where a player changes team, that is, target is not null nor 'retires'. Drop the player, year, and team_id columns.

Return the first ten lines of the dataframe.

Engineering for Categorical Variables

Categorical features are those that take only a finite number of values, and usually no categorical value has a numerical meaning, even if it happens to be number. For example in an election dataset, the names of the candidates in the race are categorical, and there is no numerical meaning (neither ordering nor size) to numbers assigned to candidates based solely on their names.

Consider the following election data.

Ballot number	For Governor	For President
001	Herbert	Romney
002	Cooke	Romney
003	Cooke	Obama
004	Herbert	Romney
005	Herbert	Romney
006	Cooke	Stein

A common mistake occurs when someone assigns a number to each categorical entry (say 1 for Cooke, 2 for Herbert, 3 for Romney, etc.). While this assignment is not, in itself, inherently incorrect, it is incorrect to use the value of this number in a statistical model. Any such model would be fundamentally wrong because a vote for Cooke cannot, in any reasonable way, be considered half of a vote for Herbert or a third of a vote for Romney. Many researchers have accidentally used categorical data in this way (and some have been very publicly embarrassed) because their categorical data was encoded numerically, which made it hard to recognize as categorical data.

Whenever you encounter categorical data that is encoded numerically like this, immediately change it either to non-numerical form (“Cooke,” “Herbert,” “Romney,”...) or apply a one-hot encoding as described below.

In order to construct a meaningful model with categorical data, one normally applies a one-hot encoding or dummy variable encoding.¹ To do this construct a new feature for every possible value of the categorical variable, and assign the value 1 to that feature if the variable takes that value and zero otherwise. Pandas makes one-hot encoding simple:

```
# one-hot encoding
df = pd.get_dummies(df, columns=['For President'])
```

The previous dataset, when the presidential race is one-hot encoded, becomes

Ballot number	Governor	Romney	Obama	Stein
001	Herbert	1	0	0
002	Cooke	1	0	0
003	Cooke	0	1	0
004	Herbert	1	0	0
005	Herbert	1	0	0
006	Cooke	0	0	1

Note that the sum of the terms of the one-hot encoding in each row is 1, corresponding to the fact that every ballot had exactly one presidential candidate.

When the gubernatorial race is also one-hot encoded, this becomes

Ballot number	Cooke	Herbert	Romney	Obama	Stein
001	0	1	1	0	0
002	1	0	1	0	0
003	1	0	0	1	0
004	0	1	1	0	0
005	0	1	1	0	0
006	1	0	0	0	1

¹Yes, these are silly names, but they are the most common names for it. Unfortunately, it is probably too late to change these now.

Now the sum of the terms of the one-hot encodings in each row is 2, corresponding to the fact that every ballot had two names—one gubernatorial candidate and one presidential candidate.

Summing the columns of the one-hot-encoded data gives the total number of votes for the candidate of that column. So the numerical values in the one-hot encodings are actually numerically meaningful, and summing the entries gives meaningful information. One-hot encoding also avoids the pitfalls of incorrectly using numerical proxies for categorical data.

The main disadvantage of one-hot encoding is that it is an inefficient representation of the data. If there are C categories and n datapoints, a one-hot encoding takes an $n \times 1$ -dimensional feature and turns it into an $n \times C$ sparse matrix. But there are ways to store these data efficiently and still maintain the benefits of the one-hot encoding.

Achtung!

When performing linear regression, it is good practice to add a constant column to your dataset and to remove one column of the one-hot encoding of each categorical variable. (Adding a constant column should only be done in linear regression).

To see why, notice that summing terms in one row corresponding to the one-hot encoding of a specific categorical variable (for example the presidential candidate) always gives 1. If the dataset already has a constant column (which you really always should add if it isn't there already), then the constant column is a linear combination of the one-hot encoded columns. This causes the matrix to fail to be invertible and can cause identifiability problems.

The standard way to deal with this is to remove one column of the one-hot embedding for each categorical variable. For example, with the elections dataset above, we could remove the Cooke and Romney columns. Doing that means that in the new dataset a row sum of 0 corresponds to a ballot with a vote for Cooke and a vote for Romney, while a 1 in any column indicates how the ballot differed from the base choice of Cooke and Romney.

When using pandas, you can drop the first column of a one-hot encoding by passing in `drop_first=True`.

Problem 4. Load `housing.csv` into a dataframe with `index=0`. Descriptions of the features are in `housing_data_description.txt`. The goal is to construct a regression model that predicts `SalePrice` using the other features of the dataset. Do this as follows:

1. Identify and handle the missing data. Hint: Dropping every row with some missing data is not a good choice because it gives you an empty dataframe. What can you do instead?
2. Identify the variable with nonnumerical values that are misencoded as numbers. One-hot encode it. Hint: don't forget to remove one of the encoded columns to prevent collinearity with the constant column).
3. Add a constant column to the dataframe.
4. Save a copy of the dataframe.
5. Choose four categorical features that seem very important in predicting `SalePrice`. One-hot encode these features, and remove all other categorical features.

6. Run an OLS regression on your model.

Print the ten features that have the highest coefficient in your model and the summary.

To run an OLS model in python, use the following code.

```
import statsmodels.api as sm

>>> results = sm.OLS(y, X).fit()

# Print the summary
>>> results.summary()

# Convert the summary table to a dataframe
>>> results_as_html = a.tables[1].as_html()
>>> result_df = pd.read_html(results_as_html, header=0, index_col=0)[0]
```

Problem 5. Using the copy of the dataframe you created in Problem 4, one-hot encode all the categorical variables. Print the shape of your database, and Run OLS.

Print the ten features that have the highest coefficient in your model and the summary.
Write a couple of sentences discussing which model is better and why.

7

Finding Patterns in Data: LSI and more about Scikit-Learn

Lab Objective: Understand the basics of principal component analysis and latent semantic indexing. Learn more about scikit-learn and implement a machine learning pipeline.

Principal Component Analysis

Principal Component Analysis (PCA) is a multivariate statistical tool used to change the basis of a set of samples from the basis of original features (which may be correlated) into a basis of uncorrelated variables called the principal components. It is a direct application of the singular value decomposition (SVD). The first principal component will account for the greatest variance in the samples, the second principal component will be orthogonal to the first and account for the second greatest variance, etc. By projecting the samples onto the space spanned by the first few principal components, we can reduce the dimensionality of the data while preserving most of the variance.

Take a matrix X with samples as rows and features as columns. The first step in PCA is to preprocess the data, which usually includes translating the columns of X to have mean 0. Some datasets require additional scaling based on variance and units of measurement. Call the new pre-processed matrix Y .

We next compute the truncated SVD of our centered data, $Y = U\Sigma V^T$, where the columns of V are the principal components and form an orthonormal basis for the space spanned by the samples. The variance captured by each principal component can be calculated by the equation below, where σ_i is the i -th nonzero singular value and there are k total singular values.

$$\frac{\sigma_i^2}{\sum_{j=1}^k \sigma_j^2} \quad (7.1)$$

In general, we are only interested in the first several principal components. But just how many principal components should we keep? One method is to keep the first two principal components so that we can project the data into 2-dimensional space. Another is to only keep the set of principal components accounting for a certain percentage of the variance, using the equation above.

Once we have decided how many principal components to keep (say the first l), we can project the samples from the original feature space onto the principal component space by computing

$$\hat{Y} = U_{:,l}\Sigma_{:l,:l} = YV_{:,l}$$

Problem 1. The breast cancer dataset from scikit-learn has 569 samples with 30 features each. Each sample is labeled as 0 (malignant) or 1 (benign). With 30 features, this data can't be directly visualized, so we will use PCA to graph the first two principal components, which account for nearly all of the variance in the data.

You can load this data using the following code.

```
>>> cancer = sklearn.datasets.load_breast_cancer()
>>> X = cancer.data
>>> y = cancer.target # Class labels (0 or 1)
```

Write a function that performs PCA on the breast cancer dataset. Graph the first two principal components, with the first along the x-axis. Your graph should resemble Figure 7.1 below. Include in the graph title the amount of variance captured by the first two principal components, calculated with Equation 7.1.

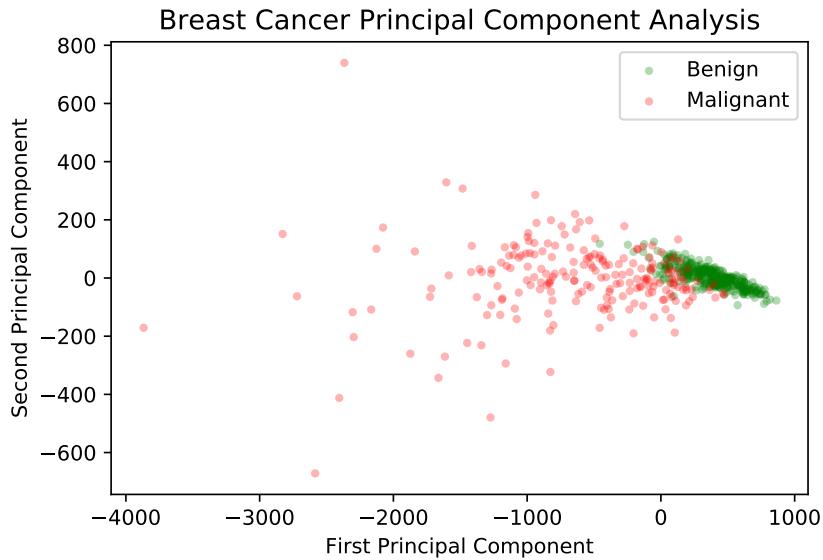


Figure 7.1: First two principal components of the transformed breast cancer data

Latent Semantic Indexing

Latent Semantic Indexing (LSI) is an application of PCA to the realm of natural language processing. In particular, LSI employs the SVD to reduce the dimensionality of a large corpus of text documents in order to enable us to evaluate the similarity between two documents. Many information-retrieval systems used in government and in industry are based on LSI.

To motivate the problem, suppose we have a large collection of documents about various topics. How can we find an article about BYU? We might consider simply choosing the article that contains the acronym the greatest number of times, but this is a crude method. A better way is to use a form of PCA on the collection of documents.

In order to do so, we need to represent the documents as numerical vectors. A standard way of doing this is to define an ordered set of words occurring in the collection of documents (called the vocabulary) and then to represent each document as a vector of word counts from the vocabulary. More formally, let our vocabulary be $V = \{w_1, w_2, \dots, w_m\}$. Then a document is a vector $x = (x_1, x_2, \dots, x_m) \in \mathbb{R}^m$ such that x_i is the number of occurrences of word w_i in the document. In this setup, we represent the entire collection of m documents as an $n \times m$ matrix X , where m is the number of vocabulary words and n is the number of documents in our collection, so each row is a document vector. As expected, we let $X_{i,j}$ be the number of times term j occurs in document i . Note that X is often a sparse matrix, as any single document likely does not contain most of the vocabulary words. This mode of representation is called the bag of words model for documents.

We calculate the SVD of X without centering or scaling the data so that we may retain the sparsity. This is unique to this particular problem. We now have $X = U\Sigma V^T$. If we are keeping l principal components, we can represent the corpus of documents by the matrix

$$\hat{X} = U_{:,l}\Sigma_{:,l}V_{:,l} = X V_{:,l}$$

Note that \hat{X} will no longer be a sparse matrix, but will have dimension $n \times l$.

Now that we have our documents represented in terms of the first l principal components, we can find the similarity between two documents. Our measure for similarity is simply the cosine of the angle between the vectors; a small angle (large cosine) indicates greater similarity, while a large angle (small cosine) indicates greater dissimilarity. Recall that we can use the inner product to find the cosine of the angle between two vectors. Under this metric, the similarity between document i and document j (represented by the i -th and j -th row of \hat{X} , notated \hat{X}_i and \hat{X}_j , respectively) is just

$$\frac{\langle \hat{X}_i, \hat{X}_j \rangle}{\|\hat{X}_i\| \|\hat{X}_j\|}.$$

To find the document most similar to document i , we simply compute

$$\operatorname{argmax}_{j \neq i} \frac{\langle \hat{X}_i, \hat{X}_j \rangle}{\|\hat{X}_i\| \|\hat{X}_j\|}.$$

Problem 2. Create a function `similar` that takes in a sparse matrix `Xhat` and an index `i` and returns the indices of the most similar and the least similar documents.

Application: State of the Union

We now discuss some practical issues involved in creating the bag of words representation X from the raw text. Our dataset will consist of the US State of the Union addresses from 1945 through 2013, each contained in a separate text file in the folder `Addresses`. We would like to avoid loading in all of the text into memory at once, and so we will stream the documents one at a time.

The first thing we need to establish is the vocabulary set, i.e. the set of unique words that occur throughout the collection of documents. A Python set object automatically preserves the uniqueness of the elements, so we will create a set and then iteratively read through the documents, adding the unique words of each document to the set. As we read in each document, we will remove punctuation and numerical characters and convert everything to lower case. The following code, found in the function `document_converter()`, will accomplish this task.

```
# Get list of file paths to each text file in the folder
>>> folder = "./Addresses/"
>>> paths = [folder+p for p in os.listdir(folder) if p.endswith(".txt")]

# Helper function to get list of words in a string
>>> def extractWords(text):
...     ignore = string.punctuation + string.digits
...     cleaned = "".join([t for t in text.strip() if t not in ignore])
...     return cleaned.lower().split()

# Initialize vocab set, then read each file and add to the vocab set.
>>> vocab = set()
>>> for p in paths:
...     with open(p, 'r') as infile:
...         for line in infile:
...             vocab.update(extractWords(line))
```

We now have a set containing all of the unique words in the corpus. However, many of the most common words do not provide important information. We call these stop words. Examples in English include the, a, an, and, I, we, you, it, there, etc; a list of common English stop words is given in `stopwords.txt`. We remove the stop words from our vocabulary set as follows and then fix an ordering to the vocabulary by creating a dictionary with key-value pairs of the form (word, index).

```
# Load stopwords.
>>> with open("stopwords.txt", 'r') as f:
...     stops = set([w.strip().lower() for w in f.readlines()])

# Remove stopwords from vocabulary, create ordering.
>>> vocab = {w:i for i, w in enumerate(vocab.difference(stops))}
```

We are now ready to create the word count vectors for each document, and we store these in a sparse matrix X . It is convenient to use the `Counter` object from the `collections` module, as this object automatically counts the occurrences of each distinct element in a list.

```
>>> from collections import Counter
>>> counts = [] # holds the entries of X
>>> doc_index = [] # holds the row index of X
>>> word_index = [] # holds the column index of X

# Iterate through the documents.
>>> for doc, p in enumerate(paths):
...     with open(p, 'r') as f:
...         # create the word counter
...         ctr = Counter()
...         for line in f:
...             ctr.update(extractWords(line))
...         # Iterate through the word counter, storing counts
...         for word, count in ctr.items():
...             counts.append((doc, word, count))
...             if len(counts) % 1000 == 0:
...                 print(len(counts))
```

```

...
    if word in vocab:
        word_index.append(vocab[word])
        counts.append(count)
    doc_index.append(doc)

# Create sparse matrix holding these word counts.
>>> X = sparse.csr_matrix((counts, [doc_index, word_index]),
...                         shape=(len(paths), len(vocab)), dtype=np.float)

```

Problem 3. Applying the techniques of LSI discussed above to the word count matrix X , and keeping the first 7 principal components, write a function that takes in the path to a single State of the Union address `speech` and returns a tuple of the addresses that are most and least similar to `speech`. For Ronald Reagan's 1984 speech, the input would be '`./Addresses/1984-Reagan.txt`', and your output should be ('`1988-Reagan`', '`1946-Truman`'). Be sure to format the strings properly.

Since X is a sparse matrix, you will need to use the SVD method found in `scipy.sparse.linalg`. This method operates slightly differently than the SVD method found in `scipy.linalg`, so be sure to read the documentation.

The simple bag of words representation is a bit crude, as it fails to consider how some words may be more important than others in determining the similarity of documents. Words appearing in few documents tend to provide more information than words occurring in every document. For example, while the word war might not be considered a stop word, it is likely to appear in many more addresses than the word Afghanistan. Two speeches sharing the word Afghanistan are probably more closely related than two speeches sharing the word war. So while $X_{i,j}$ is a good measure of the importance of term j in document i , we also need to consider some kind of global weight for each term j , indicating how important the term is over the entire collection. There are a number of different weights we could choose; we will employ the following approach. Define

$$p_{i,j} = \frac{X_{i,j}}{\sum_j X_{i,j}}.$$

We then let

$$g_j = 1 + \sum_{i=1}^m \frac{p_{i,j} \log(p_{i,j} + 1)}{\log m},$$

where m is the number of documents in the collection. We call g_j the global weight of term j . We replace each term frequency in the matrix X by weighting it globally. Specifically, we define a matrix A with entries

$$A_{i,j} = g_j \log(X_{i,j} + 1).$$

We can now perform LSI on the matrix A , whose entries are both locally and globally weighted.

Problem 4. Use the equation above to edit the function `weighted_document_converter()` to calculate the sparse matrix A . Similar to the function `document_converter()`, this function should return A and a list of file paths.

Scikit-Learn

Scikit-learn is one of the fundamental tools Python offers for machine learning. It includes classifiers, such as `RandomForestClassifier` and `KNeighborsClassifier`, as well as transformers, which preprocess data before classification. In the remainder of this lab, we will discuss transformers, validation tools, how to find optimal hyperparameters, and how to build a machine learning pipeline.

Transformers

A scikit-learn transformer processes data to make it better suited for classification. This may involve shifting or scaling data, dropping columns, replacing missing values, and so on. The function from Problem 4 is an example of a transformer, as is PCA.

Note

A hyperparameter is not dependent on data. Hyperparameters are declared in the constructor `__init__()`, before data is even passed in. Parameters set during the `fit()` method are often called model parameters and do depend on specific data. For example, a `StandardScaler` transformer shifts and scales data to have a mean of 0 and a standard deviation of 1.

Scikit-learn's transformers have three main methods: `fit_transform()`, which fits model parameters and also transforms given data; `fit()`, which sets model parameters but does not perform a transformation; and `transform()`, which transforms data according to pre-fitted model parameters. Model parameters are fitted according to training data, and they are not refitted to testing data, so a `StandardScaler` will shift and scale testing data according to the mean and variance of the training data; the transformed test data likely will not have mean 0 and variance 1.

Scikit-learn has a built-in PCA package. Its hyperparameters include the desired number of principal components and the type of SVD solver to use. Its `fit_transform()` method takes in an array of data and returns the decomposition with `n_components`.

```
>>> from sklearn.decomposition import PCA
>>> pca = PCA(n_components=5) # Create the PCA transformer with hyperparameters
>>> Xhat = pca.fit_transform(X) # Fit the transformer and transform the data
```

Problem 5. Repeat Problem 3 using your weighted document converter function and scikit-learn's built-in PCA decomposition. Do your answers seem more reasonable than before? For Bill Clinton's 1993 speech, your code should return ('./Addresses/1994-Clinton.txt', './Addresses/1951-Truman.txt').

Hint: Scikit-learn's PCA does not accept sparse matrices.

Validation Tools

We now turn our attention from transformers to classifiers. A classifier is trained to predict how a new sample should be classified or labeled. Knowing how to determine whether or not a classifier performs well is an essential part of machine learning. This often turns out to be a surprisingly sophisticated issue that largely depends on the type of problem being solved and the kind of data that is available for training. Scikit-learn has validation tools for many situations; for brevity, we restrict our attention to the simple (but important) case of binary classification, where the possible labels are only 0 or 1.

The `score()` method of a scikit-learn classifier returns the accuracy of the model, or the percent of labels predicted correctly. However, accuracy isn't always the best measure of success. Consider the confusion matrix for a classifier, the matrix where the (i, j) th entry is the number of samples with actual label i but that are classified with label j . Call the class with label 0 the negatives and the class with label 1 the positives. Then the confusion matrix is as follows.

$$\begin{array}{cc} & \text{Predicted: 0} & \text{Predicted: 1} \\ \text{Actual: 0} & \left[\begin{array}{cc} \text{True Negatives (TN)} & \text{False Positives (FP)} \\ \text{False Negatives (FN)} & \text{True Positives (TP)} \end{array} \right] \\ \text{Actual: 1} & & \end{array}$$

With this terminology, we define the following metrics.

- Accuracy: $\frac{TN + TP}{TN + FN + FP + TP}$, the percent of labels predicted correctly.
- Precision: $\frac{TP}{TP + FP}$, the percent of predicted positives that are actually correct.
- Recall: $\frac{TP}{TP + FN}$, the percent of actual positives that are predicted correctly.

Precision is useful in situations where false positives are dangerous or costly, while recall is important when avoiding false negatives takes priority. For example, an email spam filter should avoid filtering out an email that isn't actually spam; here a false positive is more dangerous, so precision is a valuable metric for the filter. On the other hand, recall is more important in disease detection: it is better to test positive and not have the disease than to test negative when the disease is actually present. Focusing on a single metric often leads to skewed results, so the following metric is also common.

$$F_\beta \text{ Score} : (1 + \beta^2) \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}} = \frac{(1 + \beta^2)TP}{(1 + \beta^2)TP + FP + \beta^2FN}.$$

Choosing $\beta < 1$ weighs precision more than recall, while $\beta > 1$ prioritizes recall over precision. The choice of $\beta = 1$ yields the common F_1 score, which weighs precision and recall equally. This is an important alternative to accuracy when, for example, the training set is heavily unbalanced with respect to the class labels.

Scikit-learn implements all of these metrics in `sklearn.metrics`. The general syntax for such functions is `some_score(actual_labels, predicted_labels)`. We will be using the function `classification_report()`, which returns precision, recall, and F_1 scores for each label. Each row in the report corresponds to a specific label and gives the scores with its label as the "positive" classification. For example, in binary classification, the row corresponding to 1 gives the scores as they would normally be calculated, with 1 as "positive."

```

>>> from sklearn.neighbors import KNeighborsClassifier
>>> from sklearn.metrics import confusion_matrix, classification_report
>>> from sklearn.model_selection import train_test_split

# Split the data into training and testing sets
>>> X_train, X_test, y_train, y_test = train_test_split(X, y)

# Fit the estimator to training data and predict the test labels.
>>> knn = KNeighborsClassifier(n_neighbors=2)
>>> knn.fit(X_train, y_train)
>>> knn_predicted = knn.predict(X_test)

# Compute the confusion matrix by comparing actual labels to predicted labels.
>>> CM = confusion_matrix(y_test, knn_predicted)
>>> CM
array([[44,  5],
       [10, 84]])

# Get precision, recall, and F1 scores all at once.
# The row labeled 1 gives these scores as we normally calculate them.
>>> print(classification_report(y_test, knn_predicted))
      precision    recall  f1-score   support

          0       0.81      0.90      0.85       49
          1       0.94      0.89      0.92       94

   accuracy                           0.90      143
  macro avg       0.88      0.90      0.89      143
weighted avg       0.90      0.90      0.90      143

```

Problem 6. For this problem, use the cancer dataset from Problem 1 to compare a `RandomForestClassifier` and a `KNeighborsClassifier`, using the default parameters for each.

Use `train_test_split()` with `random_state=2` to split up the data. Fit the classifiers with the training set and predict the labels for the testing set. Print out a classification report for each classifier, making sure to clearly label which report corresponds to which classifier.

Write a few sentences explaining which of these classifiers would be better to use in this situation and why, using the information from the report as evidence. Remember that in this dataset, the label 1 means benign and 0 means malignant.

Grid Search

Finding the optimal hyperparameters for a given model is a challenging and active area of research.¹ However, brute-force searching over a small hyperparameter space is simple in scikit-learn: a `sklearn.model_selection.GridSearchCV` object is initialized with a classifier, a dictionary of hyperparameters, and some validation parameters. When its `fit()` method is called, it tests the given classifier with every possible hyperparameter combination.

For example, a `KNeighborsClassifier` has a few important hyperparameters that can have a significant impact on the speed and accuracy of the model. These include `n_neighbors`, the number of nearest neighbors allowed to vote, and `weights`, which specifies a strategy for weighting the distances between points. The code box below tests various combinations of these hyperparameters.

The cost of a grid search rapidly increases as the hyperparameter space grows. However, the outcomes of each trial are completely independent of each other, so the problem of training each classifier is embarrassingly parallel, meaning the trials can easily be computed simultaneously. To parallelize the grid search over n CPU cores, set the `n_jobs` parameter to n , or set it to -1 to divide the labor between as many cores as are available.

```
>>> from sklearn.model_selection import GridSearchCV

>>> knn = KNeighborsClassifier()
# Specify values for certain hyperparameters
>>> param_grid = {"n_neighbors": [2, 3, 4, 5, 6],
...                 "weights": ["uniform", "distance"]}
>>> knn_gs = GridSearchCV(knn, param_grid, scoring="f1", n_jobs=-1)

# Run the actual search. This may take some time.
>>> knn_gs.fit(X_train, y_train)

# After fitting, you can access data about the results.
>>> print(knn_gs.best_params_, knn_gs.best_score_, sep='\n')
{'n_neighbors': 5, 'weights': 'uniform'}
0.9532526583188765

# Access the model
>>> knn_gs.best_estimator_
KNeighborsClassifier(weights='distance')
```

In some circumstances, the parameter grid can be organized in a way that eliminates redundancy. For example, with a `RandomForestClassifier`, you could test each `max_depth` argument with entirely different sets of values for `min_samples_leaf`. To specify certain combinations of parameters, enter the parameter grid as a list of dictionaries.

Problem 7. Do a grid search on the breast cancer dataset using a `RandomForestClassifier`. Modify at least three parameters in your grid. Use `scoring="f1"` for the `GridSearchCV` object. Fit your model with the same train-test split as in Problem 6. Print out the best parameters and the best score.

¹Intelligent hyperparameter selection is sometimes called metalearning.

Next, use the `GridSearchCV` object to predict labels for your test set. Print out a confusion matrix using these values.

Pipelines

Most machine learning problems require at least a little data preprocessing before estimation in order to get good results. A scikit-learn pipeline chains together one or more transformers and one estimator (such as a classifier) into a single object, complete with `fit()` and `predict()` methods. This simplifies and automates the machine learning process so that when you get new data or make changes to various functions and features, you can easily rerun the new version from beginning to end.

The following example demonstrates how to use a pipeline with a `StandardScaler` transformer and a `KNeighborsClassifier`. Like classifiers, pipelines have `fit()`, `predict()`, and `score()` methods. Each member of the pipeline is declared as a tuple where the first element is a string naming the step and the second is the actual transformer or classifier.

```
>>> from sklearn.preprocessing import StandardScaler
>>> from sklearn.pipeline import Pipeline

# Chain together a StandardScaler transformer and a KNN classifier.
>>> pipe = Pipeline([("scaler", StandardScaler()), # "scaler" is the step name
...                   ("knn", KNeighborsClassifier())]) # "knn" is the step name
>>> pipe.fit(X_train, y_train)
>>> pipe.score(X_test, y_test)
0.972027972027972
```

Since `Pipeline` objects follow `fit()` and `predict()` conventions, they can be used with tools like `GridSearchCV`. To specify which hyperparameters belong to which steps of the pipeline, precede each hyperparameter name with `<stepname>_`. For example, `knn__n_neighbors` corresponds to the `n_neighbors` hyperparameter of the pipeline step labeled `knn`.

```
# Create the Pipeline, labeling each step.
>>> pipe = Pipeline([("scaler", StandardScaler()),
...                   ("knn", KNeighborsClassifier())])

# Specify the hyperparameters to test for each step.
>>> pipe_param_grid = {"scaler__with_mean": [True, False],
...                      "scaler__with_std": [True, False],
...                      "knn__n_neighbors": [2,3,4,5,6],
...                      "knn__weights": ["uniform", "distance"]}

# Pass the Pipeline object to the GridSearchCV and fit it to the data.
>>> pipe_gs = GridSearchCV(pipe, pipe_param_grid,
...                         n_jobs=-1).fit(X_train, y_train)

>>> print(pipe_gs.best_params_, pipe_gs.best_score_, sep='\n')
{'knn__n_neighbors': 6, 'knn__weights': 'distance',
 'scaler__with_mean': True, 'scaler__with_std': True}
```

0.971830985915493

Pipelines can also be used to compare different transformations or estimators. For example, a pipeline can end in either a `KNeighborsClassifier()` or a classifier called `SVC()`, even though they have different hyperparameters. Like before, you can use a list of dictionaries to specify the specific combinations of the hyperparameter space.

```
# Create the pipeline, using any classifier as a placeholder
>>> pipe = Pipeline([("scaler", StandardScaler()),
                    ("classifier", KNeighborsClassifier())])

# Create the grid
>>> pipe_param_grid = [
...     {"classifier": [KNeighborsClassifier()],      # Try a KNN classifier...
...      "classifier__n_neighbors": [2,3,4,5],
...      "classifier__weights": ["uniform", "distance"]},
...     {"classifier": [SVC(kernel="rbf")],           # ...and an SVM classifier.
...      "classifier__C": [.001, .01, .1, 1, 10, 100],
...      "classifier__gamma": [.001, .01, .1, 1, 10, 100]}]

# Fit using training data
>>> pipe_gs = GridSearchCV(pipe, pipe_param_grid,
...                           scoring="f1", n_jobs=-1).fit(X_train, y_train)

# Get the best hyperparameters
>>> params = pipe_gs.best_params_
>>> print("Best classifier:", params["classifier"])
Best classifier: SVC(C=10, cache_size=200, class_weight=None, coef0=0.0,
decision_function_shape='ovr', degree=3, gamma=0.01, kernel='rbf',
max_iter=-1, probability=False, random_state=None, shrinking=True,
tol=0.001, verbose=False)

# Check the best classifier against the test data
>>> confusion_matrix(y_test, pipe_gs.predict(X_test))
array([[48,  1],                                # Near perfect!
       [ 1, 93]])
```

Problem 8. The breast cancer dataset has 30 features. By using PCA, we can drastically reduce the dimensionality while still retaining predictive power.

Create a pipeline with a `StandardScaler`, `PCA`, and a `KNeighborsClassifier`. Use the same train-test split as before. Do a grid search on this pipeline, modifying at least six hyperparameters and using `scoring="f1"`. Use no more than 5 principal components. Print out your best parameters and best score. Attain a score of at least .96.

Hint: The documentation for `StandardScaler`, `PCA`, and `KNeighborsClassifier` can be found at these links.

8

Introduction to Parallel Computing

Lab Objective: Many modern problems involve so many computations that running them on a single processor is impractical or even impossible. There has been a consistent push in the past few decades to solve such problems with parallel computing, meaning computations are distributed to multiple processors. In this lab, we explore the basic principles of parallel computing by introducing the cluster setup, standard parallel commands, and code designs that fully utilize available resources.

Parallel Architectures

Imagine that you are in charge of constructing a very large building. You could, in theory, do all of the work yourself, but that would take so long that it simply would be impractical. Instead, you hire workers, who collectively can work on many parts of the building at once. Managing who does what task takes some effort, but the overall effect is that the building will be constructed many times faster than if only one person was working on it. This is the essential idea behind parallel computing.

A serial program is executed one line at a time in a single process. This is analogous to a single person creating a building. Since modern computers have multiple processor cores, serial programs only use a fraction of the computer's available resources. This is beneficial for smooth multitasking on a personal computer because multiple programs can run at once without interrupting each other.

For smaller computations, running serially is fine. However, some tasks are large enough that running serially could take days, months, or in some cases years. In these cases it is beneficial to devote all of a computer's resources (or the resources of many computers) to a single program by running it in parallel. Each processor can run part of the program on some of the inputs, and the results can be combined together afterwards. In theory, using N processors at once can allow the computation to run N times faster. Even though communication and coordination overhead prevents the improvement from being quite that good, the difference is still substantial.

A computer cluster or supercomputer is essentially a group of regular computers that share their processors and memory. There are several common architectures that are used for parallel computing, and each architecture has a different protocol for sharing memory, processors, and tasks between computing nodes, the different simultaneous processing areas. Each architecture offers unique advantages and disadvantages, but the general commands used with each are very similar.

In this lab, we will explore the usage and capabilities of parallel computing using Python's iPyParallel package. iPyParallel can be installed with either pip or conda:

```
$ pip install ipyparallel
```

```
$ conda install ipyparallel
```

The iPyParallel Architecture

There are three main parts of the iPyParallel architecture:

- Client: The main program that is being run.
- Controller: Receives directions from the client and distributes instructions and data to the computing nodes. Consists of a hub to manage communications and schedulers to assign processes to the engines.
- Engines: The individual processors. Each engine is like a separate Python terminal, each with its own namespace and computing resources.

Essentially, a Python program using iPyParallel creates a `Client` object connected to the cluster that allows it to send tasks to the cluster and retrieve their results. The engines run the tasks, and the controller manages which engines run which tasks.

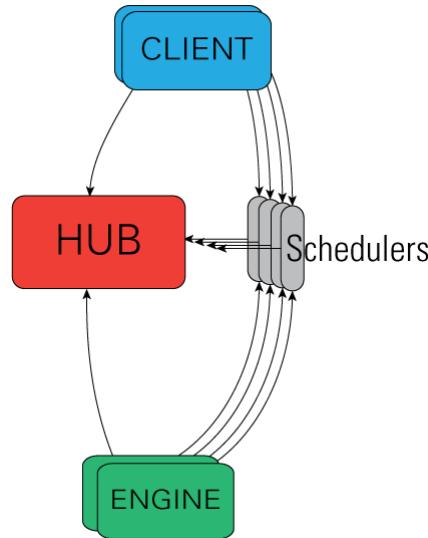


Figure 8.1: An outline of the iPyParallel architecture.

Setting up an iPyParallel Cluster

Before being able to use iPyParallel in a script or interpreter, it is necessary to start an iPyParallel cluster. We demonstrate here how to use a single machine with multiple processor cores as a cluster. Establishing a cluster on multiple machines requires additional setup, which is detailed in the Additional Material section. The following commands initialize parts or all of a cluster when run in a terminal window:

Command	Description
<code>ipcontroller start</code>	Initialize a controller process.
<code>ipengine start</code>	Initialize an engine process.
<code>ipcluster start</code>	Initialize a controller process and several engines simultaneously.

Each of these processes can be stopped with a keyboard interrupt (`Ctrl+C`). By default, the controller uses JSON files in `UserDirectory/.ipython/profile-default/security/` to determine its settings. Once a controller is running, it acts like a server, listening connections from clients and engines. Engines will connect automatically to the controller when they start running. There is no limit to the number of engines that can be started in their own terminal windows and connected to the controller, but it is recommended to only use as many engines as there are cores to maximize efficiency.

Achtung!

The directory that the controller and engines are started from matters. To facilitate connections, navigate to the same folder as your source code before using `ipcontroller`, `ipengine`, or `ipcluster`. Otherwise, the engines may not connect to the controller or may not be able to find auxiliary code as directed by the client.

Starting a controller and engines in individual terminal windows with `ipcontroller` and `ipengine` is a little inconvenient, but having separate terminal windows for the engines allows the user to see individual errors in detail. It is also actually more convenient when starting a cluster of multiple computers. For now, we use `ipcluster` to get the entire cluster started quickly.

```
$ ipcluster start          # Assign an engine to each processor core.
$ ipcluster start --n 4    # Or, start a cluster with 4 engines.
```

Note

Jupyter notebooks also have a **Clusters** tab in which clusters can be initialized using an interactive GUI. To enable the tab, run the following command. This operation may require root permissions.

```
$ ipcluster nbextension enable
```

The iPyParallel Interface

Once a controller and its engines have been started and are connected, a cluster has successfully been established. The controller will then be able to distribute messages to each of the engines, which will compute with their own processor and memory space and return their results to the controller. The client uses the `ipyparallel` module to send instructions to the controller via a `Client` object.

```
>>> from ipyparallel import Client
```

```
>>> client = Client()      # Only works if a cluster is running.
>>> client.ids
[0, 1, 2, 3]              # Indicates that there are four engines running.
```

Once the client object has been created, it can be used to create one of two classes: a `DirectView` or a `LoadBalancedView`. These views allow for messages to be sent to collections of engines simultaneously. A `DirectView` allows for total control of task distribution while a `LoadBalancedView` automatically tries to spread out the tasks equally on all engines. The remainder of the lab will be focused on the `DirectView` class.

```
>>> dview = client[:]    # Group all engines into a DirectView.
>>> dview2 = client[:2]  # Group engines 0,1, and 2 into a DirectView.
>>> dview2.targets       # See which engines are connected.
[0, 1, 2]
```

Since each engine has its own namespace, modules must be imported in every engine. There is more than one way to do this, but the easiest way is to use the `DirectView` object's `execute()` method, which accepts a string of code and executes it in each engine.

```
# Import NumPy in each engine.
>>> dview.execute("import numpy as np")
```

```
# Make sure to include client.close() after each function or else the test ←
     driver will time out
client.close()
```

Problem 1. Write a function that initializes a `Client` object, creates a `DirectView` with all available engines, and imports `scipy.sparse` as `sparse` on all engines. Return the `DirectView`. Note: Make sure to include `client.close()` after EVERY function or else the test driver will time out.

Managing Engine Namespaces

Before continuing, set the `DirectView` you are using to use blocking:

```
>>> dview.block = True
```

This affects the way that functions called using the `DirectView` return their values. Using blocking makes the process simpler, so we will use it initially. What blocking is will be explained later.

Push and Pull

The `push()` and `pull()` methods of a `DirectView` object manage variable values in the engines. Use `push()` to set variable values and `pull()` to get variables. Each method also has a shortcut via indexing.

```
# Initialize the variables 'a' and 'b' on each engine.
>>> dview.push({'a':10, 'b':5})           # OR dview['a'] = 10; dview['b'] = 5
[None, None, None]                      # Output from each engine

# Check the value of 'a' on each engine.
>>> dview.pull('a')                   # OR dview['a']
[10, 10, 10, 10]

# Put a new variable 'c' only on engines 0 and 2.
>>> dview.push({'c':12}, targets=[0, 2])
[None, None]
```

Problem 2. Write a function `variables(dx)` that accepts a dictionary of variables. Create a `Client` object and a `DirectView` and distribute the variables. Pull the variables back and make sure they haven't changed.

Scatter and Gather

Parallelization almost always involves splitting up collections and sending different pieces to each engine for processing. The process is called scattering and is usually used for dividing up arrays or lists. The inverse process of pasting a collection back together is called gathering and is usually used on the results of processing. This method of distributing a dataset and collecting the results is common for processing large data sets using parallelization.

```
>>> import numpy as np

# Send parts of an array of 8 elements to each of the 4 engines.
>>> x = np.arange(1, 9)
>>> dview.scatter("nums", x)
>>> dview["nums"]
[array([1, 2]), array([3, 4]), array([5, 6]), array([7, 8])]

# Scatter the array to only the first two engines.
>>> dview.scatter("nums_big", x, targets=[0,1])
>>> dview.pull("nums_big", targets=[0,1])
[array([1, 2, 3, 4]), array([5, 6, 7, 8])]

# Gather the array again.
>>> dview.gather("nums")
array([1, 2, 3, 4, 5, 6, 7, 8])

>>> dview.gather("nums_big", targets=[0,1])
array([1, 2, 3, 4, 5, 6, 7, 8])
```

Executing Code on Engines

Execute

The `execute()` method is the simplest way to run commands on parallel engines. It accepts a string of code (with exact syntax) to be executed. Though simple, this method works well for small tasks.

```
# 'nums' is the scattered version of np.arange(1, 9).
>>> dview.execute("c = np.sum(nums)")    # Sum each scattered component.
<AsyncResult: execute:finished>
>>> dview['c']
[3, 7, 11, 15]
```

Apply

The `apply()` method accepts a function and arguments to plug into it, and distributes them to the engines. Unlike `execute()`, `apply()` returns the output from the engines directly.

```
>>> dview.apply(lambda x: x**2, 3)
[9, 9, 9, 9]
>>> dview.apply(lambda x,y: 2*x + 3*y, 5, 2)
[16, 16, 16, 16]
```

Note that the engines can access their local variables in either of the execution methods.

Map

The built-in `map()` function applies a function to each element of an iterable. The iPyParallel equivalent, the `map()` method of the `DirectView` class, combines `apply()` with `scatter()` and `gather()`. Simply put, it accepts a dataset, splits it between the engines, executes a function on the given elements, returns the results, and combines them into one object.

```
>>> num_list = [1, 2, 3, 4, 5, 6, 7, 8]
>>> def triple(x):                      # Map a function with a single input.
...     return 3*x
...
>>> dview.map(triple, num_list)
[3, 6, 9, 12, 15, 18, 21, 24]

>>> def add_three(x, y, z):            # Map a function with multiple inputs.
...     return x+y+z
...
>>> x_list = [1, 2, 3, 4]
>>> y_list = [2, 3, 4, 5]
>>> z_list = [3, 4, 5, 6]
>>> dview.map(add_three, x_list, y_list, z_list)
[6, 9, 12, 15]
```

Blocking vs. Non-Blocking

Parallel commands can be implemented two ways. The difference is subtle but extremely important.

- Blocking: The main program sends tasks to the controller, and then waits for all of the engines to finish their tasks before continuing (the controller "blocks" the program's execution). This mode is usually best for problems in which each node is performing the same task.
- Non-Blocking: The main program sends tasks to the controller, and then continues without waiting for responses. Instead of the results, functions return an `AsyncResult` object that can be used to check the execution status and eventually retrieve the actual result.

Whether a function uses blocking is determined by default by the `block` attribute of the `DirectView`. The execution methods `execute()`, `apply()`, and `map()`, as well as `push()`, `pull()`, `scatter()`, and `gather()`, each have a keyword argument `block` that can instead be used to specify whether or not to use blocking. Alternatively, the methods `apply_sync()` and `map_sync()` always use blocking, and `apply_async()` and `map_async()` always use non-blocking.

```
>>> f = lambda n: np.sum(np.random.random(n))

# Evaluate f(n) for n=0,1,...,999 with blocking.
>>> %time block_results = [dview.apply_sync(f, n) for n in range(1000)]
CPU times: user 9.64 s, sys: 879 ms, total: 10.5 s
Wall time: 13.9 s

# Evaluate f(n) for n=0,1,...,999 with non-blocking.
>>> %time responses = [dview.apply_async(f, n) for n in range(1000)]
CPU times: user 4.19 s, sys: 294 ms, total: 4.48 s
Wall time: 7.08 s

# The non-blocking method is faster, but we still need to get its results.
# Both methods produced a list, although the contents are different
>>> block_results[10] # This list holds actual result values from each engine.
[3.833061790352166,
 4.8943956129713335,
 4.268791758626886,
 4.73533677711277]

>>> responses[10]          # This list holds AsyncResult objects.
<AsyncResult: <lambda>:finished>
# We can get the actual results by using the get() method of each AsyncResult
>>> %time nonblock_results = [r.get() for r in responses]
CPU times: user 3.52 ms, sys: 11 mms, total: 3.53 ms
Wall time: 3.54 ms          # Getting the responses takes little time.

>>> nonblock_results[10]    # This list also holds actual result values
[5.652608204341693,
 4.984164642641558,
 4.686288406810953,
 5.275735658763963]
```

When non-blocking is used, commands can be continuously sent to engines before they have finished their previous task. This allows them to begin their next task without waiting to send their calculated answer and receive a new command. However, this requires a design that incorporates checkpoints to retrieve answers and enough memory to store response objects.

Class Method	Description
<code>wait(timeout)</code>	Wait until the result is available or until <code>timeout</code> seconds pass.
<code>ready()</code>	Return whether the call has completed.
<code>successful()</code>	Return whether the call completed without raising an exception.
<code>get(timeout)</code>	Will raise <code>AssertionError</code> if the result is not ready. Return the result when it arrives. If <code>timeout</code> is not <code>None</code> and the result does not arrive within <code>timeout</code> seconds then <code>TimeoutError</code> is raised.

Table 8.1: All information from <https://ipyparallel.readthedocs.io/en/latest/details.html#AsyncResult>.

Table 8.1 details the methods of the `AsyncResult` object.

There are additional magic methods supplied by `iPyParallel` that make some of these operations easier. These methods are explained in the Additional Material section. More information on `iPyParallel` architecture, interface, and methods can also be found at <https://ipyparallel.readthedocs.io/en/latest/index.html>.

Problem 3. Write a function that accepts an integer n . Instruct each engine to make n draws from the standard normal distribution, then hand back the mean, minimum, and maximum draws to the client. Return the results in three lists.

If you have four engines running, your results should resemble the following:

```
>>> means, mins, maxs = problem3(1000000)
>>> means
[0.0031776784, -0.0058112042, 0.0012574772, -0.0059655951]
>>> mins
[-4.1508589, -4.3848019, -4.1313324, -4.2826519]
>>> maxs
[4.0388107, 4.3664958, 4.2060184, 4.3391623]
```

Problem 4. Use your function from Problem 3 to compare serial and parallel execution times. For $n = 1000000, 5000000, 10000000, 15000000$,

1. Time how long it takes to run your function.
2. Time how long it takes to do the same process serially. Make n draws and then calculate and record the statistics, but use a `for` loop with N iterations, where N is the number of engines running.

Plot the execution times against n . You should notice an increase in efficiency in the parallel version as the problem size increases.

Applications

Parallel computing, when used correctly, is one of the best ways to speed up the run time of an algorithm. As a result, it is very commonly used today and has many applications, such as the following:

- Graphic rendering
- Facial recognition with large databases
- Numerical integration
- Calculating discrete Fourier transforms
- Simulation of various natural processes (weather, genetics, etc.)
- Natural language processing

In fact, there are many problems that are only feasible to solve through parallel computing because solving them serially would take too long. With some of these problems, even the parallel solution could take years. Some brute-force algorithms, like those used to crack simple encryptions, are examples of this type of problem.

The problems mentioned above are well suited to parallel computing because they can be manipulated in such a way that running them on multiple processors results in a significant run time improvement. Manipulating an algorithm to be run with parallel computing is called parallelizing the algorithm. When a problem only requires very minor manipulations to parallelize, it is often called embarrassingly parallel. Typically, an algorithm is embarrassingly parallel when there is little to no dependency between results. Algorithms that do not meet this criteria can still be parallelized, but there is not always a significant enough improvement in run time to make it worthwhile. For example, calculating the Fibonacci sequence using the usual formula, $F(n) = F(n - 1) + F(n - 2)$, is poorly suited to parallel computing because each element of the sequence is dependent on the previous two elements.

Problem 5. The trapezoid rule is a simple technique for numerical integration:

$$\int_a^b f(x)dx \approx \frac{h}{2} \sum_{k=1}^N (f(x_k) + f(x_{k+1})),$$

where $a = x_1 < x_2 < \dots < x_N = b$ and $h = x_{n+1} - x_n$ for each n . See Figure 8.2.

Note that estimation of the area of each interval is independent of all other intervals. As a result, this problem is considered embarrassingly parallel.

Write a function that accepts a function handle to integrate, bounds of integration, and the number of points to use for the approximation. Parallelize the trapezoid rule in order to estimate the integral of f . That is, evenly divide the points among all available processors and run the trapezoid rule on each portion simultaneously. The sum of the results of all the processors will be the estimation of the integral over the entire interval of integration. Return this sum.

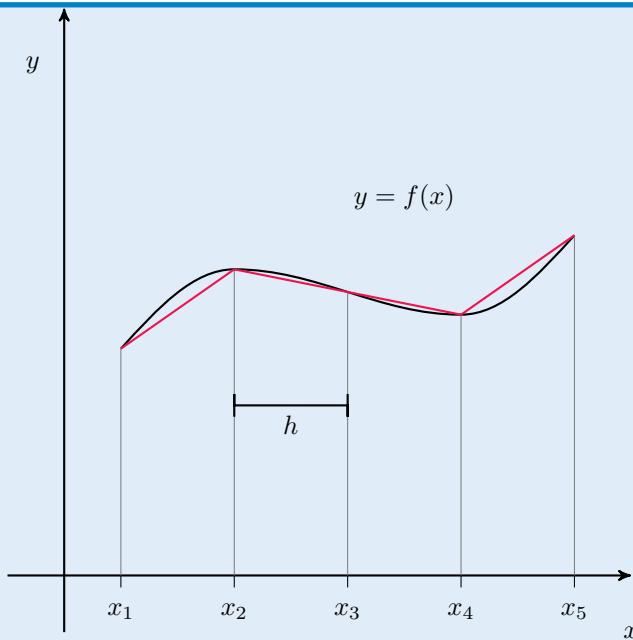


Figure 8.2: A depiction of the trapezoid rule with uniform partitioning.

Intercommunication

The phrase parallel computing refers to designing an architecture and code that makes the best use of computing resources for a problem. Occasionally, this will require nodes to be interdependent on each other for previous results. This contributes to a slower result because it requires a great deal of communication latency, but is sometimes the only method to parallelize a function. Although important, the ability to effectively communicate between engines has not been added to iPyParallel. It is, however, possible in an MPI framework and will be covered in the MPI lab.

Additional Material

Clusters of Multiple Machines

Though setting up a computing cluster with `iPyParallel` on multiple machines is similar to a cluster on a single computer, there are a couple of extra considerations to make. The majority of these considerations have to do with the network setup of your machines, which is unique to each situation. However, some basic steps have been taken from <https://ipyparallel.readthedocs.io/en/latest/process.html> and are outlined below.

SSH Connection

When using engines and controllers that are on separate machines, their communication will most likely be using an SSH tunnel. This Secure Shell allows messages to be passed over the network.

In order to enable this, an SSH user and IP address must be established when starting the controller. An example of this follows.

```
$ ipcontroller --ip=<controller IP> --user=<user of controller> --enginessh=<←
    user of controller>@<controller IP>
```

Engines started on remote machines then follow a similar format.

```
$ ipengine --location=<controller IP> --ssh=<user of controller>@<controller IP←
    >
```

Another way of affecting this is to alter the configuration file in `UserDirectory/.ipython/profile-default/security/ipcontroller-engine.json`. This can be modified to contain the controller IP address and SSH information.

All of this is dependent on the network feasibility of SSH connections. If there are a great deal of remote engines, this method will also require the SSH password to be entered many times. In order to avoid this, the use of SSH Keys from computer to computer is recommended.

Magic Methods & Decorators

To be more easily usable, the `iPyParallel` module has incorporated a few magic methods and decorators for use in an interactive iPython or Python terminal.

Magic Methods

The `iPyParallel` module has a few magic methods that are very useful for quick commands in iPython or in a Jupyter Notebook. The most important are as follows. Additional methods are found at <https://ipyparallel.readthedocs.io/en/latest/magics.html>.

%px - This magic method runs the corresponding Python command on the engines specified in `dview.targets`.

%autopx - This magic method enables a boolean that runs any code run on every engine until `%autopx` is run again.

Examples of these magic methods with a client and four engines are as follows.

```
# %px
In [4]: with dview.sync_imports():
....:     import numpy
....:
importing numpy on engine(s)
In [5]: \%px a = numpy.random.random(2)

In [6]: dview['a']
Out[6]:
[array([ 0.30390162,  0.14667075]),
 array([ 0.95797678,  0.59487915]),
 array([ 0.20123566,  0.57919846]),
 array([ 0.87991814,  0.31579495])]

# %autopx
In [7]: %autopx
%autopx enabled
In [8]: max_draw = numpy.max(a)

In [9]: print('Max_Draw: {}'.format(max_draw))
[stdout:0] Max_Draw: 0.30390161663280246
[stdout:1] Max_Draw: 0.957976784975849
[stdout:2] Max_Draw: 0.5791984571339429
[stdout:3] Max_Draw: 0.8799181411958089

In [10]: %autopx
%autopx disabled
```

Decorators

The `iPyParallel` module also has a few decorators that are very useful for quick commands. The two most important are as follows:

`@remote` - This decorator creates methods on the remote engines.

`@parallel` - This decorator creates methods on remote engines that break up element wise operations and recombine results.

Examples of these decorators are as follows.

```
# Remote decorator
>>> @dview.remote(block=True)
>>> def plusone():
...     return a+1
>>> dview['a'] = 5
>>> plusone()
[6, 6, 6, 6,]
```

```
# Parallel decorator
>>> import numpy as np

>>> @dview.parallel(block=True)
>>> def combine(A,B):
...     return A+B
>>> ex1 = np.random.random((3,3))
>>> ex2 = np.random.random((3,3))
>>> print(ex1+ex2)
[[ 0.87361929  1.41110357  0.77616724]
 [ 1.32206426  1.48864976  1.07324298]
 [ 0.6510846   0.45323311  0.71139272]]
>>> print(combine(ex1,ex2))
[[ 0.87361929  1.41110357  0.77616724]
 [ 1.32206426  1.48864976  1.07324298]
 [ 0.6510846   0.45323311  0.71139272]]
```


9

Linear Regression

Lab Objective: This section will introduce the basics of Linear Regression, feature selection methods, and regularization.

Introduction to Linear Regression

One of the first skills taught in basic algebra is to effectively plot the line $y = mx + b$ which can be done with two points. But what if we want to find the line that best fits a set of points?

In this case, we can use the simplest form of linear regression: Ordinary Least Squares (OLS). Given data as a set of points $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$ we wish to find the line that best fits the data. The line is given by $y = mx + b$ where m and b are unknown constants and x and y are the independent and dependent variables respectively. Using OLS, let

$$y_i = mx_i + b + \varepsilon_i$$

describe the i th point in D for each $i \in \{1, \dots, n\}$. Note that ε_i is the vertical distance from the i th point to the line given by $y = mx + b$ and is often called the residual or the error.

The n equations for each point in D can be written in vector notation. Let the x and y coordinates of D be represented by column vectors \mathbf{x} and \mathbf{y} respectively. In statistical science, the intercept (b) and slope (m) are denoted as β_0 and β_1 respectively and

$$\begin{bmatrix} b \\ m \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \boldsymbol{\beta}.$$

Additionally, the residuals are represented by a column vector $\boldsymbol{\varepsilon}$ and $\mathbf{1}$ is a column vector of ones. So we have

$$\mathbf{y} = m\mathbf{x} + \mathbf{1}b + \boldsymbol{\varepsilon} = [\mathbf{1}, \mathbf{x}] \cdot \begin{bmatrix} b \\ m \end{bmatrix} + \boldsymbol{\varepsilon}.$$

Denoting $X = [\mathbf{1}, \mathbf{x}]$, we have our final equation given as

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

This notation may seem excessive, but suppose we wanted to fit a model of the form $y = ax^3 + bx^2 + cx + d$. A little work can show that $X = [\mathbf{1}, \mathbf{x}, \mathbf{x}^2, \mathbf{x}^3]$ and $\boldsymbol{\beta} = [\beta_0, \beta_1, \beta_2, \beta_3]^T$, which is very easy to work with. Thus, this notation is actually the ideal way to generalize linear regression, especially when working with higher degree polynomials.

The solution to OLS is straight forward with some important assumptions. Sparing you the algebraic details and assuming that $\mathbf{y} \sim \mathcal{N}(X\beta, \sigma^2 \mathbf{I})$ and $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ and \mathbf{I} is the identity matrix, the least squares estimator for β is given as

$$\hat{\beta} = (X^T X)^{-1} X^T \mathbf{y}. \quad (9.1)$$

Problem 1. Write a function that takes as input \mathbf{X} and \mathbf{y} . In your function, add a column of ones to \mathbf{X} to account for β_0 . Call this function `ols`. This function should return the least squares estimator for β as a numpy array.

Hint: Use functions from `numpy` or `scipy` to calculate a matrix inverse

Problem 2. Use the following code to generate random data.

```
n = 100 # Number of points to generate
X = np.arange(100) # The input X for the function ols
eps = np.random.uniform(-10,10, size=(100,)) # Noise to generate random y ←
    coordinates
y = .3*X + 3 + eps # The input y for the function ols
```

Find the least squares estimator for β using this random data. Produce a plot showing the random data and the line of best fit determined by the least squares estimator for β . Your plot should include a title, axis labels, and a legend.

Hint: Since `ols` takes \mathbf{X} without a column of ones, slice \mathbf{X} when you call `ols`.

Rank-Deficient Models

Notice that in order to find the least squares estimator $\hat{\beta}$, we need $X^T X$ to be invertible. However, when X does not have full rank, the product $X^T X$ is singular and not invertible. We can no longer use the previous solution for the least squares estimator, but we can use the SVD and still compute a solution.

Recall that if $X \in M_{n \times d}$ has rank r , then the compact form of the SVD of X is

$$X = U \Sigma V^H$$

where $U \in M_{n \times r}$ and $V \in M_{r \times d}$ have orthonormal columns and $\Sigma \in M_{r \times r}$ is diagonal. In addition, if X is real, then the factors U , Σ , and V^H are also real. In this lab we assume X is real. As described in Volume 1, there is a unique solution for the least squares estimator given by

$$\hat{\beta} = V \Sigma^{-1} U^T \mathbf{y}. \quad (9.2)$$

Problem 3. Write a function that finds the least squares estimator for rank-deficient models using the SVD. The function should still take \mathbf{X} and \mathbf{y} as inputs. In your function, add a column of ones to \mathbf{X} to account for β_0 . Call the function `svd_ols` and return the least squares estimator for β as a numpy array.

Hint: Use `np.linalg.svd` to factor `X` and use the argument `full_matrices=False`.

Problem 4. Use the following code to generate random data:

```
x = np.linspace(-4, 2, 500)
y = x**3 + 3*x**2 - x - 3.5
eps = np.random.normal(0, 3, len(y)) # Create noise
y += eps # Add noise to randomize data
```

Now use your function `svd_ols` to find the least squares estimator for a cubic polynomial. Create a plot that shows a scatter plot of the data and a curve using the least squares estimator. Your plot should include a title, axis labels, and a legend.

Model Accuracy

Residual Sum of Squares

The Residual Sum of Squares (RSS) is a common choice of measure for the quality of a model. The formula for RSS is given by

$$RSS = \|\mathbf{y} - X\hat{\boldsymbol{\beta}}\|_2^2.$$

Notice that the RSS measures the variance in the error of the model. So relative to other models, a smaller RSS value indicates a more accurate model.

Coefficient of Determination

Another method of model accuracy is the Coefficient of Determination, denoted R^2 . In the case of linear regression,

$$R^2 = 1 - \frac{RSS}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

and $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ is the sample mean of \mathbf{y} . The intuition of R^2 is that the ratio of the average residual and biased sample variance of \mathbf{y} is approximately the total variance explained by the model. A larger R^2 corresponds to a model that fits better. However, R^2 comes with flaws such as being able to take negative values, rewarding overfitting, and punishing under-fit models. Because of this, we typically want to use other methods for model accuracy.

Python Example

There are various python packages that can be used to calculate R^2 , but we will use `statsmodels` in this lab. Below is an example of how to build a model and extract R^2 using `statsmodels`.

```
import statsmodels.api as sm
data = pd.read_csv("/filepath") # Read in data as pandas dataframe
y = data["dependent_variable"] # Extract dependent variable
temp_X = data[["var_1", ..., "var_n"]] # Extract independent variables
```

```
X = sm.add_constant(temp_X) # Add column of 1's
model = sm.OLS(y, X).fit() # Fit the linear regression model
print(model.rsquared) # Print the R squared value
```

Problem 5. The file `realestate.csv` contains transaction data from 2012-2013. It has columns for transaction date, house age, distance to nearest MRT station, number of convenience stores, latitude, longitude, and house price of unit area. ^a Each row in the array is a separate measurement.

Find the combination of variables that builds the model with the best R^2 value when predicting `house price of unit area`. Use `statsmodels` to build each model and calculate R^2 . Using the same combination of variables, time the methods `ols`, `svd_ols`, and `statsmodels`. Return a list with the first element being a tuple of times for each method and the second element being the best R^2 value from the first part of the problem.

Hint: The `combinations` method from the `itertools` package will be very helpful for finding all feature combinations.

^aSee <https://www.kaggle.com/datasets/quantbruce/real-estate-price-prediction?resource=download>.

Feature Selection

Every regression model consists of features or variables used to predict a dependent variable or result. An important question to ask when building regression models is, which features are the most important in predicting the dependent variable? In addition to being used for model accuracy, R^2 can also be used in feature selection, as it was in Problem 5. It still has the same pitfalls of rewarding overfitting and punishing under-fit models, but it can be a useful tool used in conjunction with the following tools for feature selection. While there are other methods for implementing feature selection, most incorporate the p-value and are not included in this lab.

Akaike's Information Criterion (AIC)

A simple motivation for AIC is based on balancing goodness of fit and prescribing a penalty for model complexity. A more rigorous motivation for AIC is given in Volume 3 using the Kullback-Leibler (KL) divergence. Given two models, f and g , the KL divergence is given by

$$KL(f, g) = \int f(z) \log \left(\frac{f(z)}{g(z)} \right) dz$$

and it measures the amount of information lost when g is used to model f . Thus, a lower AIC value indicates a better model. Additionally, AIC penalizes the size of the parameter space with a coefficient of 2 which allows for slightly more complex models.

Bayesian Information Criterion (BIC)

Instead of estimating the KL-divergence between the model in question and the true model, BIC has the property of being minimized precisely when the posterior probability of a model, given the data, is maximized. The equations for AIC and BIC only differ with one term: the coefficient weighting the size of the parameter space. The coefficient for BIC is $\log(n)$ which is generally much larger than 2. As a result, BIC penalizes complex models more than AIC. The difference in AIC and BIC values will grow from having more data points.

When using AIC or BIC for feature selection, you need to consider how you want to penalize features in your model. If you want to exclude irrelevant features, then use BIC. If you want to keep all features that are relevant, then use AIC. In other words, BIC is more likely to choose too small a model, and AIC is more likely to choose too large a model.

Python Example

There are multiple ways to calculate AIC and BIC with various python packages. We will use the package `statsmodels` for the following problem. When constructing X for `statsmodels`, do not add the column of 1's manually because `statsmodels` has a method that will do this for us.

```
import statsmodels.api as sm
data = pd.read_csv("/filepath") # Read in data as pandas dataframe
y = data["dependent_variable"] # Extract dependent variable
temp_X = data[["var_1", ..., "var_n"]] # Extract independent variables
X = sm.add_constant(temp_X) # Add column of 1's
model = sm.OLS(y, X).fit() # Fit the linear regression model
print(model.aic) # or print(model.bic)
```

Problem 6. Use the file `realestate.csv` and the Python Example above as a template for constructing y and X and calculating model AIC and BIC. For the dependent variable, use `house price of unit area`. For the independent variables, use `distance to the nearest MRT station, number of convenience stores, latitude, and longitude`.

Find the model that has the lowest AIC and the model that has the lowest BIC. Are they the same model? Print the features of the model with the lowest AIC as a list.

Hint: The `combinations` method from the `itertools` package will be very helpful for finding all feature combinations.

Regularization

Up to this point, we have been solving the problem

$$\min_{\beta} \|X\beta - \mathbf{y}\|_2^2.$$

However, we have also assumed independence among the features used to predict the dependent variable. The pitfall of multicollinearity arises when the features of X have dependence and X becomes nearly singular. As a result, the least squares estimator is susceptible to random noise or error. Multicollinearity typically occurs when data is collected with poor experimental design. It is important to have good experimental design, but regularization can be used to mitigate poor design. Another issue OLS faces is feature selection. While there are feature selection methods available, regularization can be used to minimize non-zero coefficients.

Ridge Regularization Regression

The problem posed by Ridge Regularization is

$$\min_{\beta} \|X\beta - \mathbf{y}\|_2^2 + \alpha\|\beta\|_2^2$$

where $\alpha \geq 0$. This essentially penalizes the size of the coefficients. The larger α is, the more the model resists multicollinearity.

Lasso Regularization Regression

The problem posed by Lasso Regularization is

$$\min_{\beta} \frac{1}{n} \|X\beta - \mathbf{y}\|_2^2 + \alpha\|\beta\|_1.$$

Note that α provides the same functionality here as it does in Ridge Regularization. However, the use of the 1-norm often results in sparse solutions. As a result, Lasso Regularization can be used for feature selection since it only includes the most important features.

Python Example

Since α is not a fixed value in Ridge and Lasso Regularization, it is best practice to perform a Grid-Search to find the best parameter value. The example below goes over the syntax for implementing Ridge Regularization. Note that the syntax for Lasso Regularization is similar.

```
>>> from sklearn import linear_model
>>> y = # dependent variable data
>>> X = # independent variable data with no column of ones
>>> reg = linear_model.RidgeCV(alphas=np.logspace(-6, 6, 13)) # Range for grid ←
    search
>>> reg.fit(X, y) # Fit the model
>>> reg.alpha_ # Best parameter value
```

Problem 7. Use Ridge and Lasso Regression to model `house price of unit area` from the file `realestate.csv`. First, do a grid search for the model parameter. Then use the grid search result to fit the model. Once you have fit the model, you can use the `score` method to get R^2 . Print R^2 for each model as a tuple. How do these models compare to the models in problem 6?

10

Logistic Regression

Lab Objective: Understand the basic principles of Logistic Regression and binary classifiers. Apply this to a dataset.

Linear regression is unsuitable for predicting probabilities, because the resulting model may take values in any fixed interval in \mathbb{R} , but a probability-predicting model can only take values in the interval $[0, 1]$. Logistic regression is a form of regression that always takes its values in the interval $[0, 1]$ and as such, is a popular method for predicting probabilities and for constructing classifiers. As in linear regression, in a classification problem we have a random variable Y , conditioned on an input $X \in \mathbb{R}^d$. However, in binary classification problems the random variable Y is binary, that is, $Y \in \{0, 1\}$. A binary classifier is any function f taking values in $\{0, 1\}$. For example, $\mathbf{x} \in \mathbb{R}^d$ could be the pixel intensities of an image, and the classifier f gives 1 if the image is a picture of a duck and 0 otherwise. The goal of a classification problem is to choose a classifier \hat{f} so that $(X, \hat{f}(X))$ is a good approximation for (X, Y) .

Logistic Regression

Logistic regression relies heavily on the logistic function, also known as the sigmoid function, $\text{sigm} : \mathbb{R} \rightarrow (0, 1)$ given by

$$\text{sigm}(x) = \frac{1}{1 + e^{-x}}. \quad (10.1)$$

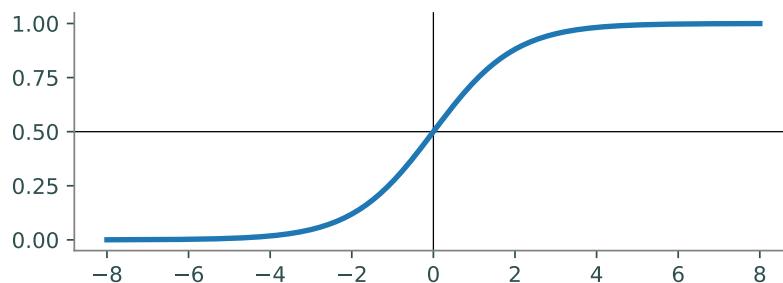


Figure 10.1: Sigmoid Function

This function works well for classifying objects based on probabilities, because it has some key properties that translate well into probability theory. Of particular note, the graph can be translated by adding a constant, giving the form $\text{sigm}(\beta_1 t + \beta_0)$. A larger value of β_1 makes the ramp up from 0 to 1 steeper, while a smaller value of β_1 makes it less steep. The trick behind logistic regression is to find the values of β_i such that the resulting sigmoid function best classifies the data.

In logistic regression models we have a random variable Y with support $\{0, 1\}$, where Y is conditioned on another random variable X , with support in \mathbb{R}^d . The distribution of Y , given X , is assumed to be Bernoulli

$$Y | X \sim \text{Bernoulli}(\text{sigm}(X^\top \boldsymbol{\beta})),$$

so that

$$P(Y | X) = \text{sigm}(X^\top \boldsymbol{\beta}) = \frac{1}{1 + \exp(-X^\top \boldsymbol{\beta})}.$$

As in the case of linear regression, we usually add a constant feature $X_0 = \mathbf{1}$ to X and a corresponding coefficient β_0 to $\boldsymbol{\beta}$, so that $X^\top \boldsymbol{\beta} = \beta_0 + \beta_1 X_1 + \dots + \beta_d X_d$. Given a draw of length n of the form $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ we wish to estimate $\boldsymbol{\beta}$. The maximum likelihood estimator is a good choice. To find this estimator, first observe that the likelihood of $\boldsymbol{\beta}$, given the data, is

$$\begin{aligned} L(\boldsymbol{\beta} | D) &= \prod_{i=1}^n P(Y = y_i, X = \mathbf{x}_i | \boldsymbol{\beta}) \\ &= \prod_{i=1}^n P(Y = y_i | X = \mathbf{x}_i, \boldsymbol{\beta}) P(X_i). \end{aligned}$$

which is equivalent to maximizing

$$\prod_{i=1}^n P(Y = y_i | X = \mathbf{x}_i, \boldsymbol{\beta}) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1 - y_i}$$

where

$$p_i = P(Y = 1 | \mathbf{x}_i, \boldsymbol{\beta}) = \text{sigm}(\mathbf{x}_i^\top \boldsymbol{\beta}) = \frac{1}{1 + \exp(-\mathbf{x}_i^\top \boldsymbol{\beta})}.$$

Taking the negative logarithm turns this into a convex minimization problem, and a little math shows that

$$\ell(\boldsymbol{\beta} | D) = \sum_{i=1}^n (y_i \log(1 + \exp(-\mathbf{x}_i^\top \boldsymbol{\beta})) + (1 - y_i) \log(1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta}))) . \quad (10.2)$$

The convexity of this problem implies there is a unique minimizer $\hat{\boldsymbol{\beta}}$ of $\ell(\boldsymbol{\beta} | D)$.

Problem 1. Create a Python classifier called `LogiReg` that accepts an $(n \times 1)$ array y of binary labels (0's and 1's) as well as an $(n \times d)$ array X of data points. Write a `fit()` method that uses equation 10.2 to find and save the optimal $\hat{\boldsymbol{\beta}}$.

Once the maximum likelihood estimate $\hat{\beta}$ is found, we have an estimate for the probability

$$P(Y = 1 \mid \mathbf{x}) \approx \text{sigm}(\mathbf{x}^\top \hat{\beta}).$$

From this, we can construct a classifier \hat{f} by setting $\hat{f}(x) = 1$ if $P(Y = 1 \mid \mathbf{x}) \geq \frac{1}{2}$ and $\hat{f}(x) = 0$ otherwise.

Problem 2. Write a method called `predict_prob()` for your classifier that accepts an $(n \times d)$ array x_test and returns $P(Y = 1 \mid x_test)$. Also write a method called `predict()` that calls `predict_prob()` and returns an array of predicted labels (0's or 1's) for the given array x_test .

Problem 3. To test your classifier, create training arrays X and y as well as testing array X_test . The code to generate X , y , and X_test is provided below. Both X and X_test have 100 random draws from a 2-dimensional multivariate normal distribution centered at $(1, 2)$, and another 100 draws from one centered at $(4, 3)$.

Train your classifier on X and y . Then generate a list of predicted labels using your trained classifier and X_test , and use it to plot X_test with a different color for each predicted label. Your plot should look similar to Figure 10.2.

```
>>> import numpy as np

>>> data = np.column_stack((
        # draw from 2 2-dim. multivariate normal dists.
        np.concatenate((
            np.random.multivariate_normal(np.array([1,2]), np.eye(2), 100),
            np.random.multivariate_normal(np.array([4,3]), np.eye(2), 100)
        )),
        # labels corresponding to each distribution
        np.concatenate((np.zeros(100), np.ones(100))) ))
>>> np.random.shuffle(data)
>>> # extract X and y from the shuffled data
>>> X = data[:, :2]
>>> y = data[:, 2].astype(int)

>>> X_test = np.concatenate((
        # draw from 2 identical 2-dim. multivariate normal dists.
        np.random.multivariate_normal(np.array([1,2]), np.eye(2), 100),
        np.random.multivariate_normal(np.array([4,3]), np.eye(2), 100)
    ))
>>> np.random.shuffle(X_test)
```

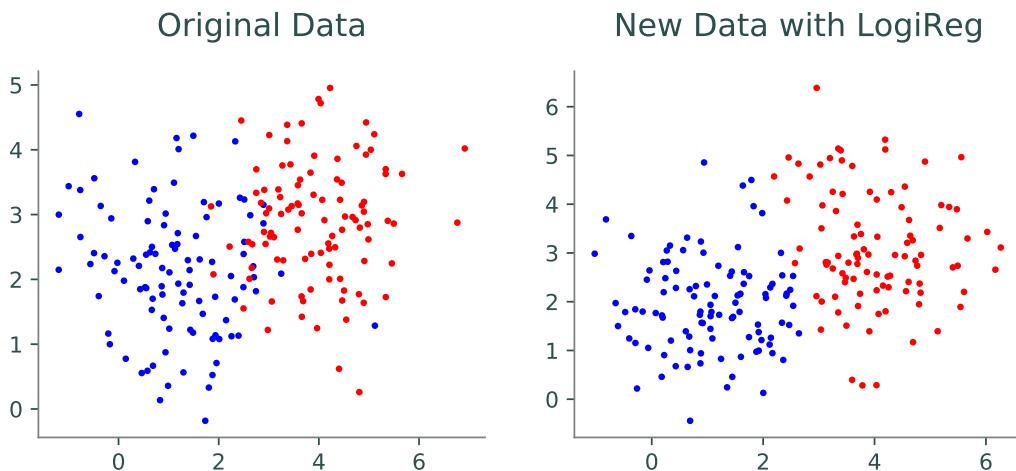


Figure 10.2: In reality, these two distributions overlap a little, but the logistic regression model makes a clean divide between the two.

Statsmodels and Sklearn

The module `statsmodels` contains a package that includes a logistic regression class called `Logit`. A simple example of this class being implemented is as follows.

```
>>> import statsmodels.api as sm

>>> model = sm.Logit(y, X).fit(disp=0) # setting disp=0 turns off printed info
>>> probs = model.predict(X_test) # list of probabilities, not labels
```

`Logit` does not add a constant feature (column of 1's) to X by default, so in order to do so, you must apply the function `sm.add_constant()` to both X and X_test . In addition, the `.fit()` method does not regularize the problem by default, which may lead to some errors involving singular matrices. To fix this, you can use the `.fit_regularized()` method instead of `.fit()`.

The module `sklearn` also has a package for logistic regression called `LogisticRegression`, which can be implemented as follows.

```
>>> from sklearn.linear_model import LogisticRegression

>>> model = LogisticRegression(fit_intercept=True).fit(X, y) # X before y
>>> labels = model.predict(X_test) # predicted labels of X_test
```

`LogisticRegression` already regularizes the problem by default. The parameter `fit_intercept` (which defaults to `False`) indicates whether you want to add a constant feature (column of 1's) to X and X_test .

You can also use `sklearn` to score a logistic regression model. After fitting an `sklearn` model, you can call `<model>.score(X_test, y_test)` to find the percentage of accuracy of the model's prediction for X_test , given the true labels in y_test . Alternatively, you can use `sklearn.metrics.accuracy_score` to find the percentage of accuracy between a list of predicted labels and the list of true labels.

```
>>> from sklearn.metrics import accuracy_score

>>> true_labels = [0, 1, 2, 3, 4]
>>> pred_labels = [0, 2, 2, 2, 4] # predicted labels from logistic regression
>>> accuracy_score(true_labels, predicted_labels)
0.6
```

Problem 4. The code to generate arrays X , y , X_test , and y_test is provided below. X and X_test are each composed of 200 draws from two 20-dimensional multivariate normal distributions, one centered at **0**, and the other centered at **2**.

Using each of `LogiReg`, `statsmodels`, and `sklearn`, train a logistic regression classifier on X and y to generate a list of predicted labels for X_test . Then, using y_test , print the accuracy scores for each trained model. Compare the accuracies and training/testing time for all three classifiers. Be sure to add a constant feature with each model.

```
>>> # redefine the true beta
>>> beta = np.random.normal(0, 7, 20)

>>> # X is generated from 2 20-dim. multivariate normal dists.
>>> X = np.concatenate((
    np.random.multivariate_normal(np.zeros(20), np.eye(20), 100),
    np.random.multivariate_normal(np.ones(20)*2, np.eye(20), 100)
))
>>> np.random.shuffle(X)
>>> # create y based on the true beta
>>> pred = 1. / (1. + np.exp(-X @ beta))
>>> y = np.array([1 if pred[i] >= 1/2 else 0
                 for i in range(pred.shape[0])])

>>> # X_test and y_test are generated similar to X and y
>>> X_test = np.concatenate((
    np.random.multivariate_normal(np.zeros(20), np.eye(20), 100),
    np.random.multivariate_normal(np.ones(20), np.eye(20), 100)
))
>>> np.random.shuffle(X_test)
>>> pred = 1. / (1. + np.exp(-X_test @ beta))
>>> y_test = np.array([1 if pred[i] >= 1/2 else 0
                     for i in range(pred.shape[0])])
```

Multiclass Logistic Regression

Sometimes we may want to classify data into more than two categories, but so far we've only used logistic regression as a binary classifier. The good news is that we can extend logistic regression to classify more than just two categories.

The more popular method for doing this is to generalize the logistic regression model to a multi-class setting. This method is called multinomial logistic regression or sometimes softmax regression. While standard logistic regression was based on the sigmoid function, multinomial logistic regression is based on the softmax function $\mathcal{S} : \mathbb{R}^k \rightarrow (0, 1)^k$, which is a multivariate version of the sigmoid function, given by

$$\mathcal{S}(t_1, \dots, t_k) = \left(\frac{e^{t_1}}{\sum_{j=1}^k e^{t_j}}, \dots, \frac{e^{t_k}}{\sum_{j=1}^k e^{t_j}} \right). \quad (10.3)$$

We will assume that $Y | X$ is categorically distributed as

$$\text{Cat}(p_1(X), \dots, p_k(X)) = \text{Cat}(\mathcal{S}(X^\top \boldsymbol{\beta}_1, \dots, X^\top \boldsymbol{\beta}_k))$$

for some choice of vectors $\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k$, which we will estimate from the data. Here

$$p_i(X) = P(Y = i | X) = \frac{e^{X^\top \boldsymbol{\beta}_i}}{\sum_{j=1}^k e^{X^\top \boldsymbol{\beta}_j}} = \frac{\text{sigm}(X^\top \boldsymbol{\beta}_i)}{\sum_{j=1}^k \text{sigm}(X^\top \boldsymbol{\beta}_j)}.$$

Given a draw of length n of the form $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, we wish to compute $\boldsymbol{\theta} = (\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_k)$ where, without loss of generality, we may assume $\boldsymbol{\beta}_k = \mathbf{0}$. The maximum likelihood estimate of $\boldsymbol{\theta}$ is computed in a manner similar to the way it was for standard logistic regression. A bit of math shows that

$$\begin{aligned} \ell(\boldsymbol{\theta} | D) &= - \sum_{i=1}^n \sum_{j=1}^k \delta_{c_j}(y_i) \log(p_j(\mathbf{x}_i)) \\ &= - \sum_{i=1}^n \sum_{j=1}^k \delta_{c_j}(y_i) \log \left(\frac{e^{\mathbf{x}_i^\top \boldsymbol{\beta}_j}}{\sum_{m=1}^k e^{\mathbf{x}_i^\top \boldsymbol{\beta}_m}} \right) \end{aligned}$$

where

$$\delta_{c_j}(y_i) = \begin{cases} 1 & \text{if } y_i = c_j, \text{ the jth class} \\ 0 & \text{otherwise.} \end{cases}$$

This is a convex minimization problem with unique minimizer $\hat{\boldsymbol{\theta}}$. Once $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\beta}}_1, \dots, \hat{\boldsymbol{\beta}}_k)$ is found, we have an estimate for the probability

$$P(Y = y | \mathbf{x}) \approx \frac{e^{\mathbf{x}^\top \hat{\boldsymbol{\beta}}_y}}{\sum_{j=1}^k e^{\mathbf{x}^\top \hat{\boldsymbol{\beta}}_j}}.$$

From this, we can construct a classifier \hat{f} by setting $\hat{f}(\mathbf{x}) = \text{argmax}_j P(Y = c_j | \mathbf{x})$.

Conveniently, `sklearn` has a very simple implementation of multinomial logistic regression that simply requires the argument `multi_class='multinomial'` when initiating a `LogisticRegression` model.

```
>>> from sklearn.linear_model import LogisticRegression

>>> model = LogisticRegression(
        multi_class='multinomial',
        fit_intercept=True).fit(X, y) # add constant feature
```

Problem 5. The Iris Dataset contains information taken from 150 samples of 3 different types of iris flowers (Setosa, Versicolor, and Virginica). The columns contain measurements for sepal length, sepal width, pedal length, and pedal width. Import the Iris Dataset and perform a train-test split on only the first two columns of the data with `test_size=0.4`. Train a multinomial logistic regression model using the training data with an added constant feature, and generate prediction labels for the test data. Plot the test data by color using your prediction labels.

Your plot should reflect Figure 10.3

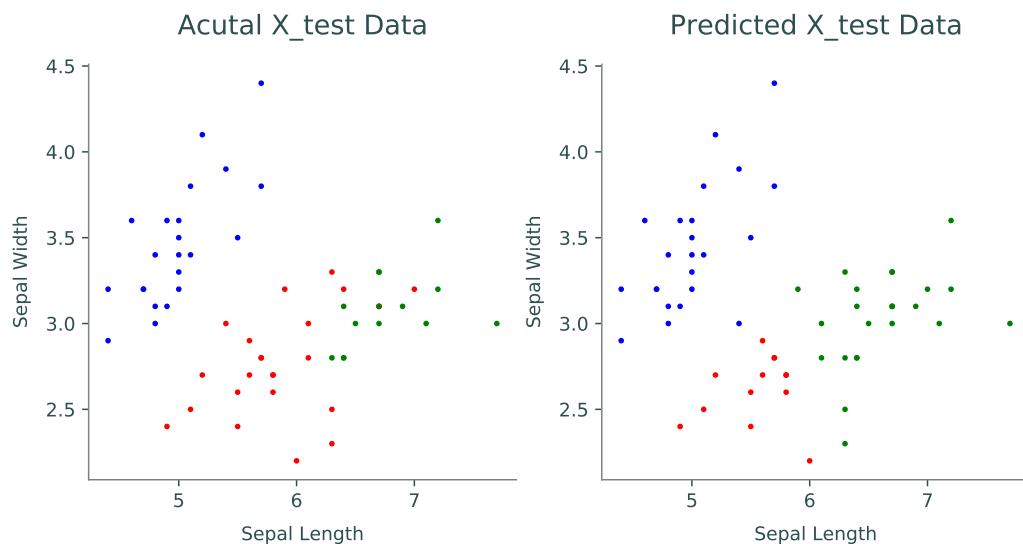


Figure 10.3: Multinomial logistic regression attempt to categorize the Iris Dataset.

11

Naive Bayes

Lab Objective: Create a Naïve Bayes Classifier. Use this classifier, and Sklearn's premade classifier to make an SMS spam filter.

About Naïve Bayes

Naïve Bayes classifiers are a family of machine learning classification methods that use Bayes' theorem to probabilistically categorize data. They are called naïve because they assume independence between the features. The main idea is to use Bayes' theorem to determine the probability that a certain data point belongs in a certain class, given the features of that data. Despite what the name may suggest, the naïve Bayes classifier is not a Bayesian method. This is because naïve Bayes is based on likelihood rather than Bayesian inference.

While naïve Bayes classifiers are most easily seen as applicable in cases where the features have, ostensibly, well defined probability distributions (such as classifying sex given physical characteristics), they are applicable in many other cases. While it is generally a bad idea to assume independence naïve Bayes classifiers are still very effective, even when we can be confident there is nonzero covariance between features.

The Classifier

You are likely already familiar with Bayes' Theorem, but we will review how we can use Bayes' Theorem to construct a robust machine learning model.

Given the feature vector of a piece of data we want to classify, we want to know which of all the classes is most likely. Essentially, we want to answer the following question

$$\operatorname{argmax}_{k \in K} P(C = k | \mathbf{x}), \quad (11.1)$$

where C is the random variable representing the class of the data. Using Bayes' Theorem, we can reformulate this problem into something that is actually computable. We find that for any $k \in K$ we have

$$P(C = k | \mathbf{x}) = \frac{P(C = k)P(\mathbf{x} | C = k)}{P(\mathbf{x})}.$$

Now we will examine each feature individually and use the chain rule to expand the following expression

$$\begin{aligned} P(C = k)P(\mathbf{x} | C = k) &= P(x_1, \dots, x_n, C = k) \\ &= P(x_1 | x_2, \dots, x_n, C = k)P(x_2, \dots, x_n, C = k) \\ &= \dots \\ &= P(x_1 | x_2, \dots, x_n, C = k)P(x_2 | x_3, \dots, x_n, C = k) \cdots P(x_n | C = k)P(C = k), \end{aligned}$$

and applying the assumption that each feature is independent we can drastically simplify this expression to the following

$$P(x_1 | x_2, \dots, x_n, C = k) \cdots P(x_n | C = k) = \prod_{i=1}^n P(x_i | C = k).$$

Therefore we have that

$$P(C = k | \mathbf{x}) = \frac{P(C = k)}{P(\mathbf{x})} \prod_{i=1}^n P(x_i | C = k),$$

which reforms Equation 11.1 as

$$\operatorname{argmax}_{k \in K} P(C = k | \mathbf{x}) = \operatorname{argmax}_{k \in K} P(C = k) \prod_{i=1}^n P(x_i | C = k). \quad (11.2)$$

We drop the $P(\mathbf{x})$ in the denominator since it is not dependent on k .

This problem is approximately computable, since we can use training data to attempt to find the parameters which describe $P(x_i | C = k)$ for each i, k combination, and $P(C = k)$ for each k . In reality, a naïve Bayes classifier won't often find the actual correct parameters for each distribution, but in practice the model does well enough to be robust. Something to note here is that we are actually computing $P(C = k | \mathbf{x})$ by finding $P(C = k, \mathbf{x})$. This means that naïve Bayes is a generative classifier, and not a discriminative classifier.

Spam Filters

A spam filter is just a special case of a classifier with two classes: spam and not spam (or ham). We can now describe in more detail how we are to calculate Equation 11.2 given that we know what the features are. We can use a labeled training set to determine $P(C = \text{spam})$ the probability of spam and $P(C = \text{ham})$ the probability of ham. To do this we will assume that the training set is a representative sample and define

$$P(C = \text{spam}) = \frac{N_{\text{spam}}}{N_{\text{samples}}}, \quad (11.3)$$

and

$$P(C = \text{ham}) = \frac{N_{\text{ham}}}{N_{\text{samples}}}. \quad (11.4)$$

Using a bag of words model, we can create a simple representation of $P(x_i | C = k)$ where x_i is the i^{th} word in a message, and therefore \mathbf{x} is the entire message. This results in the simple definition of

$$P(x_i | C = k) = \frac{\text{Noccurrences of } x_i \text{ in class } k}{N_{\text{words in class } k}}. \quad (11.5)$$

Note that the denominator in Equation 11.5 is not the number of unique words in class k , but the total number of occurrences of any word in class k . In the case we have some word x_u that is not found in the training set, we can choose $P(x_u | C = k)$ so that the computation is not effected, i.e. letting $P(x_u | C = k) = 1$ for unique words.

A First Model

When building a naïve Bayes classifier we need to choose what probability distribution we believe our features to have. For this first model, we will assume that the words are a categorically distributed random variable. This means the random variable may take on say N different values, each value has a certain probability of occurring. This distribution can be thought of as a Bernoulli trial with N outcomes instead of 2.

In our situation we may have N different words which we expect may occur in a spam or ham message, so we need to use the training data to find each word and its associated probability. In order to do this we will make a DataFrame that will allow us to calculate the probability of the occurrence of a certain word x_i based on what percentage of words in the training set were that word x_i . This DataFrame that will allow us to more easily compute Equation 11.5, assuming the words are categorically distributed. While we are creating this DataFrame, it will also be a good opportunity to compute Equations 11.3 and 11.4.

Throughout the lab we will use an SMS spam dataset contained in `sms_spam_collection.csv`. The following code makes full test and train sets, but we will also provide you with code to check against specific subsets.

```
>>> import pandas as pd
>>> from sklearn.model_selection import train_test_split

>>> # load in the sms dataset
>>> df = pd.read_csv('sms_spam_collection.csv')

>>> # separate the data into the messages and labels
>>> X = df.Message
>>> y = df.Label

>>> # split the data into test and train sets
>>> X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.7)
```

Training The Model

Problem 1. Create a class `NaiveBayesFilter`, with an `__init__()` method that is empty. Add a `fit()` method which takes as arguments `X`, the training data, and `y` the training labels. In this case `X` is a `pandas.Series` containing strings that are SMS messages. Create a new `DataFrame` with two rows and a column for each vocabulary word with `'spam'` and `'ham'` being the index. Each entry will be the number of times a word appears in spam or ham messages.

For example, `self.data.loc['ham', 'in']` is the number of times the word "in" appears in ham messages. Save this DataFrame as `self.data`.

Hint: be sure you count the number of occurrences of a word and not a string. For example, when searching the string `'find it in there'` for the word `'in'`, make sure you get 1 and not 2 (because of the `'in'` in `'find'`). The methods `pd.Series.str.split()` and `count()` may be helpful.

```
>>> # checkout what the DataFrame looks like
```

```
>>> NB = NaiveBayesFilter()
>>> NB.fit(X[:300], y[:300])
>>> NB.data.loc['ham','in']
47
>>> NB.data.loc['spam','in']
4
```

Predictions

Now that we have implemented the `fit()` method, we can begin to classify new data. We will do this with two methods, the first will be a method that calculates $P(S | \mathbf{x})$ and $P(H | \mathbf{x})$, and the other will determine the more likely of the two and assign a label. While it may seem like we should have $P(C = S | \mathbf{x}) = 1 - P(C = H | \mathbf{x})$, we do not. This would only be true if we assume the S and H are independent of \mathbf{x} . It is clear that we shouldn't make this assumption, because we are trying to determine the likelihood of S and H based on what \mathbf{x} tells us. Therefore we must compute both $P(C = S | \mathbf{x})$ and $P(C = H | \mathbf{x})$.

Problem 2. Implement the `predict_proba()` method in your naïve Bayes classifier. It should take as an argument \mathbf{X} , the data that needs to be classified, and it will compute the product portion of equation 11.2.

Notice that $P(x_i | C)$ is the same for every repeated instance of word x_i in message \mathbf{x} . To save time, we only want to calculate this probability once. To do this, find

$$\prod_{i=1}^l P(x_i | C)^{n_i}$$

for each message \mathbf{x} in \mathbf{X} where l is the number of unique words in the message and n_i is the number of times the i^{th} unique word (x_i) occurs.

The method should return an $(N, 2)$ array, where N is the length of \mathbf{X} , whose entries are the probabilities of each message \mathbf{x} in \mathbf{X} belonging to each category. The first column corresponds to $P(C = H | \mathbf{x})$, and the second to $P(C = S | \mathbf{x})$.

Problem 3. Implement the `predict()` method in your naïve Bayes classifier. This should take as an argument \mathbf{X} , the data that needs to be classified. Using `predict_proba()`, finish implementing equation 11.2 and return an array that classifies each message \mathbf{x} in \mathbf{X} .

```
>>> # create the filter
>>> NB = NaiveBayesFilter()

>>> # fit the filter to the first 300 data points
>>> NB.fit(X[:300], y[:300])

>>> # test the predict function
```

```
>>> NB.predict(X[530:535])
array(['ham', 'spam', 'ham', 'ham', 'ham'], dtype=object)
```

Underflow

There are some issues that we encounter given this implementation. Notice that in the following example, the likelihoods for both spam and ham are 0 for each message.

```
>>> # find the likelihoods for messages 1085 and 2010
>>> NB.predict_proba(X[[1085, 2010]])
array([[0., 0.],
       [0., 0.]])
```

This is because the messages are long, and thus involve the product of many numbers that are between 0 and 1. Because of this, we have encountered what is called underflow, where a number becomes so small it is not machine representable. Therefore, we should work in logspace, as to avoid inevitable underflow caused by long messages. If we take the log of equation 11.2 have

$$\operatorname{argmax}_{k \in K} \ln(P(C = k)) + \sum_{i=1}^n \ln(P(x_i | C = k)), \quad (11.6)$$

and this problem is still valid since logarithms are monotonically increasing. However, if any of the $P(x_i | C = k)$ are close to 0, we risk getting an overall value of $-\infty$. To prevent this from happening, we can perform Laplace add-one smoothing by adding 1 to the numerator of $P(x_i | C = k)$ and 2 to its denominator. This method is equivalent to using a Bayesian method for training. Thus, equation 11.5 becomes

$$P(x_i | C = k) = \frac{N_{\text{occurrences of } x_i \text{ in class } k} + 1}{N_{\text{words in class } k} + 2}. \quad (11.7)$$

Problem 4. Implement `predict_log_proba()` and `predict_log()` using equations 11.6 and 11.7. These methods will take the same arguments and return the same object types as the methods `predict_proba()` and `predict()`, respectively.

Notice how `X[[1085, 2010]]` is now classifiable.

The Poisson Model

Now that we've examine one way to constructing a naïve Bayes classifier, let us look at one more method. In the Poisson model we assume that each word is Poisson random variable, occurring with potentially different frequencies among spam and ham messages. Because each of the messages is a different length, we can reparameterize the Poisson PMF to the following

$$P(n_i = x) = \frac{(rn)^x e^{-rn}}{x!} \quad (11.8)$$

where n_i is the number of times word i occurs in a message, n is the length of the message, and $\lambda = rn$ is the classical Poisson rate. In this case r represents the number of events per unit time/space/etc.

We could easily refactor this model to use Bayesian inference to determine r , which would allow greater control over the model. This would also create a condition where the training data doesn't completely determine the model's viability. However, in this lab we will use maximum likelihood estimation to determine r .

Training the Model

Similar to the other classifier that we made, training the model amounts to using the training data to determine how $P(x_i | C = k)$ is computed, as well as computing $P(C = k)$. As stated earlier, we will attempt to find the most likely value of r for each word that appears in the training set. To do this we will use maximum likelihood estimation. The parameter we choose is the one that maximizes the likelihood function

$$\hat{r} = \operatorname{argmax}_r L(r | \mathbf{x}) = \operatorname{argmax}_r P(\mathbf{x} | r).$$

In this case, since we are using a Poisson distribution (11.8) for each word, we will solve the following problem for both the spam class and the ham classes

$$r_{i,k} = \operatorname{argmax}_{r \in [0,1]} \frac{(rN_k)^{n_i} e^{-rN_k}}{n_i!}, \quad (11.9)$$

where $r_{i,k}$ is the Poisson rate for the i^{th} word in class k (either spam or ham), N_k is the total number of words in class k , and n_i is the number of times the i^{th} word occurs in class k . We have $r \in [0, 1]$ because a word cannot occur more than once per word in the message. If we take the derivative of the right side of equation 11.9 with respect to r , set it equal to 0, and solve for the maximizing r , we find that $r_{i,k} = n_i/N_k$.

Predictions

Making predictions with this model is exactly the same as we did earlier. To clarify the calculation, let's reformulate 11.6 to fit the Poisson case better. This gives

$$\operatorname{argmax}_{k \in K} \ln(P(C = k)) + \sum_{i=1}^l \ln \left(\frac{(r_{i,k}n)^{n_i} e^{-r_{i,k}n}}{n_i!} \right), \quad (11.10)$$

with l being the number of unique words in the message, n_i the number of times the i^{th} word occurs in the message, n the total number of words in the message, and $r_{i,k}$ the Poisson rate of the i^{th} word in class k . Notice, if $r_{i,k}$ is close to 0, we'll risk getting a total value of $-\infty$. We can fix this by using the Laplace add-one smoothing method as we did before, but this time on $r_{i,k}$. Thus, our new Poisson rate for the i^{th} word in class k becomes

$$r_{i,k} = \frac{n_i + 1}{N_k + 2}, \quad (11.11)$$

which has a Bayesian interpretation, as it did before.

Problem 5. Create a new class called `PoissonBayesFilter` with an `__init__()` method that may be empty. Add a `fit()` method which takes as arguments training data `X` and training labels `y`.

Implement `fit()` by finding the MLE found in equation 11.11 to predict r for each word in both the spam and ham classes, thereby training the model. Store these computed rates in dictionaries called `self.spam_rates` and `self.ham_rates`, where the key is the word and the value is the associated r .

For example, `self.ham_rates['in']` will give the computed r value for the word "in" found in ham messages.

```
>>> #create a poisson bayes object to examine it
>>> PB = PoissonBayesFilter()
>>> PB.fit(X[:300], y[:300])

>>> # check spam and ham rate of 'in'
>>> PB.ham_rates['in']
0.012588512981904013
>>> PB.spam_rates['in']
0.004166666666666667
```

Problem 6. Implement the `predict_log_proba()` and `predict()` methods using equation 11.10. These methods will take the same arguments and return the same object types as the methods `predict_proba()` and `predict()` in the `NaiveBayesFilter` class, respectively. You may use `scipy.stats.poisson.pmf` if you wish.

Naive Bayes with Sklearn

Now that we've explored a few ways to implement our own naïve Bayes classifier, we can examine some robust tools from the `sklearn` library that will accomplish all the things we've coded up in a very simple manner.

The first thing we need to do is create a dictionary and transform the training data, which is what our first `fit()` method did. We instantiate a `CountVectorizer` model from `sklearn.feature_extraction.text`, and then use the `fit_transform()` method to create the dictionary and transform the training data.

```
>>> from sklearn.feature_extraction.text import CountVectorizer

>>> vectorizer = CountVectorizer()
>>> train_counts = vectorizer.fit_transform(X_train)
```

Now we can use the transformed training data to fit a `MultinomialNB` model from `sklearn.naive_bayes`

```
>>> from sklearn.naive_bayes import MultinomialNB

>>> clf = MultinomialNB()
>>> clf = clf.fit(train_counts, y_train)
```

Testing data we want to classify must first be transformed by our vectorizer with the `transform()` method (not the `fit_transform()` method). We can then classify the data using the `predict()` method of the `MultinomialNB` model.

```
>>> test_counts = vectorizer.transform(X_test)
>>> labels = clf.predict(test_counts)
```

This naïve Bayes model uses the multinomial distribution where we have used the categorical and poisson distributions. Multinomial is very similar to the categorical implementation, as the multinomial distribution models the outcome of n categorical trials (in the same way that the binomial distribution models n Bernoulli trials).

Problem 7. Write a function that will classify messages. It will take as arguments training data `X_train` and `y_train`, and test data `X_test`. In this function use the `CountVectorizer` and `MultinomialNB` from `sklearn` and return the predicted classification of the model.

The results of Problem 7 can help you test the two Bayes Filters you created in this lab. Using the `accuracy_score` method of `sklearn.metrics`, you can compare your predicted labels with the ones from Problem 7. You should have very high accuracy, as demonstrated below.

```
>>> from sklearn.metrics import accuracy_score

>>> # labels returned by Problem 7
>>> actual_labels = sklearn_method(X_train, y_train, X_test)

>>> # test against NaiveBayesFilter
>>> NB = NaiveBayesFilter()
>>> NB.fit(X_train, y_train)
>>> NB_labels = NB.predict_log(X_test)
>>> accuracy_score(actual_labels, NB_labels)
0.9769289925660087

>>> # test against PoissonBayesFilter
>>> PB = PoissonBayesFilter()
>>> PB.fit(X_train, y_train)
>>> PB_labels = PB.predict(X_test)
>>> accuracy_score(acutal_labels, PB_labels)
0.9782107152012305
```

References

Rish, Irina. (2001). An Empirical Study of the Naïve Bayes Classifier. IJCAI 2001 Work Empir Methods Artif Intell. 3.

Data from: <http://www.dt.fee.unicamp.br/~tiago/smsspamcollection/>

12

Random Forests

Lab Objective: Understand how to build and use a classification tree and a random forest.

Classification Trees

Decision Classification trees are a class of decision trees used in a wide variety of settings where labeled training data is available. The desired outcome is a model that can accurately assign labels to unlabeled data. Decision trees are widely used because they have a fast run time, low computation cost, and can handle irrelevant, missing, and noisy data easily.

We begin with a data set of samples, such as information about customers from a certain store. Each sample contains a variety of features, such as if the individual is married or has children. The sample also has a classification label, such as whether or not the person made a specific purchase.

A classification tree is composed of many nodes, which ask a question (i.e. “Is income ≥ 85 ?”) and then split the data based on the answers. If the response is `True`, then the sample is “pushed” down the tree to the left child node. If the response is `False`, then the sample is “pushed” down the tree to the right child node. A leaf node is a node that has no child node. Upon arrival at a leaf, an unlabeled sample is labeled with the classification that matches the majority of labeled samples at that leaf. Table 12.1 includes information about 10 individuals and then an indicator of whether or not they made a certain purchase. To simplify construction of the tree, all data is numeric, so 1=Yes and 0=No for yes/no questions.

Suppose we wanted to guess whether a single college student making under \$30,000 would purchase this item. Starting at the top of the tree, we compare our sample to the question and first choose the right branch, and then we compare with the second question and choose the right branch again. Now we reach a leaf with the dictionary `{0:1}`. The key 0 corresponds to the label, and the value 1 means one of our original samples is at this leaf with that label. Since 100% of samples at this leaf are labeled with 0, our new sample college student will be predicted to share the label 0.

If we arrived instead at a leaf with the dictionary `{0:1, 1:4}`, then one of our original samples at this leaf would be labeled 0 and four would be labeled 1, so the majority vote would assign the label 1 to our new sample.

Married (Y/N)	Children	Income (\$1000)	Purchased (Y/N)
0	5	125	0
1	0	100	0
0	0	70	0
1	3	120	0
0	0	95	1
1	0	60	0
0	2	220	1
0	0	85	1
1	0	75	0
0	0	90	1

Table 12.1: Customer data with 3 features (Married, Children, Income) and a label (Purchase) indicating whether or not the customer bought the item.

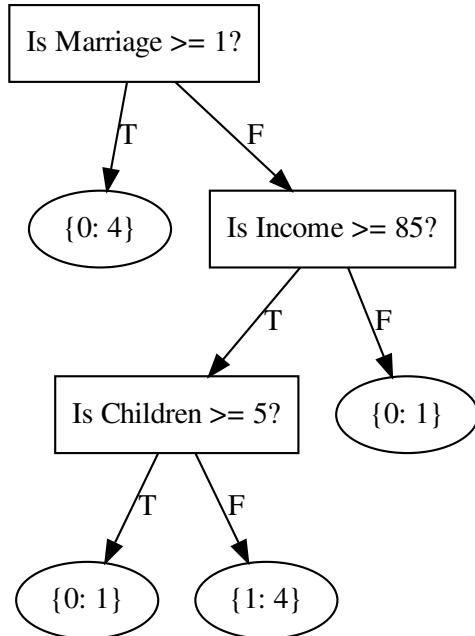


Figure 12.1: A classification tree built using Table 12.1. Each leaf includes a dictionary of the label (0 or 1) and how many individuals from the data match the classification. In this example, each leaf contains individuals with only one label.

Problem 1. At each node in a classification tree, a question indicates which branch a sample belongs to. Write a `match` method for the class `Question` that accepts a sample and returns `True` or `False` depending on how the sample's features compare to the question. This method will handle one feature at a time. For example, in the example above, a single college student making \$20,000 would be a sample represented by the array `[0, 0, 20]`.

Next, write a `partition` function that partitions a data set for a given question into two `numpy` arrays: `left` and `right`. Note that `left` will contain samples where the `match` method returns `True` and `right` will contain samples where the `match` method returns `False`. Return the left and right regions of the partition in that order. If one region is empty, return it as `None`.

Measures

To use the `partition` function from Problem 1, we need to know which question to ask at each node. Usually, the question is determined by the split that maximizes either the Gini impurity or the information gain. Gini impurity measures how often a sample would be mislabeled based on the distribution of labels. It is a measure of homogeneity of labels, so it is 0 when all samples at a node have the same label.

Definition 12.1. Let D be a data set with K different class labels and N different samples. Let N_k be the number of samples labeled class k for each $1 \leq k \leq K$, and let $f_k = \frac{N_k}{N}$. We define the Gini impurity to be

$$G(D) = 1 - \sum_{k=1}^K f_k^2.$$

Information gain is based on the concept of Information Theory entropy. It measures the difference between two probability distributions. If the distributions are equal, then the information gain is 0. We will use a modified version of information gain for simplicity.

Definition 12.2. Let $s_D(p, x) = D_1, D_2$ be a partition of data D . We define the information gain of this partition to be

$$I(s_D(p, x)) = G(D) - \sum_{i=1}^2 \frac{|D_i|}{|D|} \cdot G(D_i)$$

where $|D|$ represents the number of samples (or rows) in D .

Problem 2. Write a function `gini()` that computes the Gini impurity of an array of data with the class labels in the last column. Write another function `info_gain()` that computes the information gain for a given split of data. Make sure these functions account for the case of the data array containing only a single sample.

The file `animals.csv` contains information about 7 features for 100 animals. The last column, the class labels, indicates whether or not an animal lives in the ocean. You may use this file to test your functions. To test your functions, your values should match those below.

```
>>> import numpy as np
# Load in the data
>>> animals = np.loadtxt('animals.csv', delimiter=',')
# Load in feature names
>>> features = np.loadtxt('animal_features.csv', delimiter=',', dtype=str,
...                         comments=None)
```

```
# Load in sample names
>>> names = np.loadtxt('animal_names.csv', delimiter=',', dtype=str)

# Test your functions
>>> gini(animals)
0.4758
# split animals into two sets with fifty animals in each
>>> info_gain(animals[:50], animals[50:], gini(animals))
0.1457999999999999
```

Optimal Split

The optimal split of a data set can be chosen by maximizes either the Gini impurity or the information gain. We will optimize the information gain, so the optimal split is

$$s_D^* = s_D(p^*, x^*),$$

where

$$p^*, x^* = \operatorname{argmax}_{p,x} I(s_D(p, x)).$$

Sometimes the partition to split on may separate the data into very small subsets with only a few samples each. This can make the classification tree vulnerable to overfitting and noisy data. For this reason, classification trees include an argument to specify the smallest allowable leaf size, or the minimum number of samples at the node. This number depends on the size of the whole data set; for example, data with 10,000 samples would have a larger minimum leaf than our first example using data with only 10 samples.

Problem 3. Write a function `find_best_split()` that computes the optimal split of a data set by checking through all possible `Questions` associated with the data (each unique value in each feature (column)). Recall that the final column has the class label and will have no possible questions associated with it. Include a minimum leaf argument defaulting to 5. Do not allow the best split to include a leaf smaller than this size. Return the information gain and question associated with the best split. If two splits have the same information gain, choose the first split.

The output for the animals data set should be

`(0.12259833679833688, Is # legs/tentacles >= 2.0?).`

Building the Tree

Once the optimal split is determined, the node is defined to be a Leaf node or a Decision node. As described earlier, leaf nodes have no children nodes and is where the classification for a sample is made. If the optimal split returns a left and right tree, then the node is a decision node and has a question associated with it to determine which path a sample should follow. The next two problems will walk through building a classification tree using the functions and classes from the previous problems.

Problem 4. Write the class `Leaf`. It should have an attribute `prediction` that is the dictionary of how many samples at the leaf belong to each label, as shown in the leaves of Figure 12.1.

Next, write the class `Decision_Node`. This should have three attributes: an associated `Question`, a left branch, and a right branch. The branches will be `Leaf` or `Decision_Node` objects. Name these three attributes `question`, `left`, and `right`.

In addition to having a minimum leaf size, it's also important to have a maximum depth for trees. Without restricting the depth, the tree can become very large; if there is no minimum leaf size, it can be one less than the number of training samples. Limiting the depth can stop the tree from having too many splits, preventing it from becoming too complex and overfitting the training data. It's also important to not have too shallow of a tree because then the tree will underfit the data.

Problem 5. Write a function `build_tree()` that uses your previous functions to build a classification tree. Include a minimum leaf argument defaulting to 5 and a maximum depth argument defaulting to 4. Start counting depth at 0. For comparison, the tree in Figure 12.1 has depth 3. You will probably want to build this tree recursively.

Make a `Leaf` if the remaining data has too few samples, if the depth is too much, or if the information gain is 0. Otherwise, make a partition and build a new tree for each branch, returning those as `Decision_Nodes`.

The last column in the `animals.csv` file indicates whether or not the animal lives in the ocean; this is the class label for this data set. Test your classifier with this file and the function `draw_tree`. This will display and save a pdf of the graph. Examine the figure and test various parameters to check if your functions are working properly.

```
# How to draw a tree
>>> my_tree = build_tree(animals, features)
>>> draw_tree(my_tree)
```

Achtung!

The function `draw_tree` relies on the `graphviz` package. There are two options to aid in installing the `graphviz` package.

- You can try downloading by typing `conda install -c conda-forge python-graphviz` if you have the Anaconda distribution. If `draw_tree` returns an error about pdf being an unrecognized file type, try typing `dot -c` in your terminal.
- If you get an error related to a PATH problem you may need to download `graphviz` to your computer by following the instructions found at this link: Download `graphviz`.

Predicting

It's important to test your tree to ensure that it predicts class labels fairly accurately and so that you can adjust the minimum leaf and maximum depth parameters as needed. It is customary to randomly assign some of your labeled data to a training set that you use to fit your tree and then use the rest of your data as a testing set to check accuracy.

Problem 6. Write a function `predict_tree` that returns the predicted class label for a new sample given a trained tree. You will probably have to make this recursive in order to traverse the branches and reach a `Leaf` node with prediction information.

Next, write a function `analyze_tree` that accepts a labeled data set (with the labels in the last column, as in `animals.csv`) and a trained classification tree and returns the proportion of samples that the tree labels correctly.

Test your function with the `animals.csv` file. Shuffle the data set with `np.random.shuffle()` and use 80 samples to train your classification tree. Use the other 20 samples as the test set to see how accurately your tree classifies them. Your tree should be able to classify this set with roughly 80% accuracy on average, given the default parameters.

Random Forest

As noted, one of the main issues with Decision Trees is their tendency to overfit. Random forests are a way of mitigating overfitting that cannot be fixed by restricting the depth and leaf size. A random forest is just what it sounds like—a collection of trees. Each tree is trained randomly, meaning that at each node, only a small, random subset of the features is available by which to determine the next split. The size of this subset should be small relative to the total number of features present. Let n be the total number of features in the data set. One common method, and the one we will use here, is to split on \sqrt{n} features, rounding down where applicable.

When predicting the label of a new sample, each trained tree in the forest casts a vote, determined as above, and the sample is labeled according to the majority vote of the trees.

Problem 7. Add an argument `random_subset` to `build_tree()` and `find_best_split()`, defaulting to `False`, that indicates whether or not the tree should be trained randomly. When `True`, each node should be restricted to a random combination of \sqrt{n} features to use in its split, where n is the total number of features (note that class labels are not features).

Next, write a function `predict_forest()` that accepts a new sample and a trained forest (as a list of trees). It should return the assigned label, found by majority vote of the trees.

Finally, write a function `analyze_forest()` that accepts a labeled data set and a trained forest and analyzes the accuracy of the forest's predictions.

Test your functions out on the `animals.csv` file. Examine the graphs of the individual trees to see how they compare to the non-randomized versions.

Scikit-Learn

Next, we'll compare our implementation to scikit-learn's `RandomForestClassifier`. Rather than accepting all the data as a single array, as in our implementation, this package accepts the feature data as the first argument and all of the labels as the second argument.

```
>>> from sklearn.ensemble import RandomForestClassifier

# Create the forest with the appropriate arguments and 200 trees
>>> forest = RandomForestClassifier(n_estimators=200, max_depth=4,
...                                 min_samples_leaf=5)

# Shuffle the data
>>> shuffled = np.random.permutation(animals)
>>> train = shuffled[:80]
>>> test = shuffled[80:]

# Fit the model to your data, passing the labels in as the second argument
>>> forest.fit(train[:, :-1], train[:, -1])

# Test the accuracy with the testing set
>>> forest.score(test[:, :-1], test[:, -1])
0.85
```

Problem 8. The file `parkinsons.csv` contains annotated speech data from people with and without Parkinson's Disease. The first column is the subject ID, columns 2-27 are various features, and the last column is the label indicating whether or not the subject has Parkinson's. You will need to remove the first column so your forest doesn't use participant ID to predict class labels. Feature names are contained in the file `parkinsons_features.csv`.

Write a function to compare your forest implementation to the package from scikit-learn. Because of the size of this data set, we will only use a small portion of the samples and build a very simple forest. Randomly select 130 samples. Use 100 in training your forest and 30 more in testing it. Include 5 trees in the forest and use `min_samples_leaf=15`. Time how long it takes to train and analyze your forest.

Repeat this with scikit-learn's package, using the same 100 training samples and 30 test samples. Set `n_estimators=5` and `min_samples_leaf=15`.

Next, using scikit-learn's package, run the whole data set, using the default parameters. Use 80% of the data to train the forest and the other 20% to test it.

Return three tuples, where each tuple contains the accuracy and time for each variation.

13 Apache Spark

Lab Objective: Dealing with massive amounts of data often requires parallelization and cluster computing; Apache Spark is an industry standard for doing just that. In this lab we introduce the basics of PySpark, Spark's Python API, including data structures, syntax, and use cases. Finally, we conclude with a brief introduction to the Spark Machine Learning Package.

Apache Spark

Apache Spark is an open-source, general-purpose distributed computing system used for big data analytics. Spark is able to complete jobs substantially faster than previous big data tools (i.e. Apache Hadoop) because of its in-memory caching, and optimized query execution. Spark provides development APIs in Python, Java, Scala, and R. On top of the main computing framework, Spark provides machine learning, SQL, graph analysis, and streaming libraries.

Spark's Python API can be accessed through the PySpark package. You must install Spark, along with the supporting tools like Java, on your local machine for PySpark to work. This will include ensuring that both Java and Spark are included in the environment variable PATH.¹ Installation of PySpark for local execution or remote connection to an existing cluster can be done with conda or pip commands.²

```
# install Java
$ sudo apt-get install openjdk-8-jdk
# check the version, it may not be exactly the same
$ java -version
openjdk version "1.8.0_242"
OpenJDK Runtime Environment (build 1.8.0_242-b09)
OpenJDK 64-Bit Server VM (build 25.242-b09, mixed mode)

# Install Spark by following instructions in the footnote
# Following these steps, you must configure your PATH environment variable

# CHOOSE ONE
```

¹See the Apache Spark configuration instructions for detailed installation instructions

²You may also use the script provided with the spec file that will completely install Spark and its requirements. Note however that this script is provided AS IS and is not the recommended method.

```
# PySpark installation with conda
$ conda install -c conda-forge pyspark

# PySpark installation with pip
$ pip install pyspark
```

If you use `python3` in your terminal, you will need to set the `PYSPARK_PYTHON` environment variable to `python3`. When using an IDE, you must call it from the terminal or set the variables inside the editor so that PySpark can be found.

PySpark

One major benefit of using PySpark is the ability to run it in an interactive environment. One such option is the interactive Spark shell that comes prepackaged with PySpark. To use the shell, simply run `pyspark` in the terminal. In the Spark shell you can run code one line at a time without the need to have a fully written program. This is a great way to get a feel for Spark. To get help with a function use `help(function)`; to exit the shell simply run `quit()`.

In the interactive shell, the `SparkSession` object - the main entrypoint to all Spark functionality - is available by default as `spark`. When running Spark in a standard Python script (or in IPython) you need to define this object explicitly. The code box below outlines how to do this. It is standard practice to name your `SparkSession` object `spark`.

Achtung!

It is important that when you are finished with a `SparkSession` you should end it by calling `spark.stop()`.

Note

While the interactive shell is very robust, it may be easier to learn Spark in an environment that you are more familiar with (like IPython). To do so, just use the code given below. Help can be accessed in the usual way for your environment. Just remember to `stop()` the `SparkSession`!

```
>>> from pyspark.sql import SparkSession

# instantiate your SparkSession object
>>> spark = SparkSession\
...         .builder\
...         .appName("app_name")\
...         .getOrCreate()

# stop your SparkSession
>>> spark.stop()
```

Note

The syntax

```
>>> spark = SparkSession\  
...     .builder\  
...     .appName("app_name")\  
...     .getOrCreate()
```

is somewhat unusual. While this code can be written on a single line, it is often more readable to break it up when dealing with many chained operations; this is standard styling for Spark. Note that you cannot write a comment after a line continuation character '`\`'.

Resilient Distributed Datasets

The most fundamental data structure used in Apache Spark is the Resilient Distributed Dataset (RDD). RDDs are immutable distributed collections of objects. They are resilient because performing an operation on one RDD produces a new RDD without altering the original; if something goes wrong, you can always go back to your original RDD and restart. They are distributed because the data resides in logical partitions across multiple machines. While RDDs can be difficult to work with, they offer the most granular control of all the Spark data structures.

There are two main ways of creating RDDs. The first is reading a file directly into Spark and the second is parallelizing an existing collection (list, numpy array, pandas dataframe, etc.). We will use the Titanic dataset³ in most of the examples throughout this lab. The example below shows various ways to load the Titanic dataset as an RDD.

```
# initialize your SparkSession object  
>>> spark = SparkSession\  
...     .builder\  
...     .appName("app_name")\  
...     .getOrCreate()  
  
# load the data directly into an RDD  
>>> titanic = spark.sparkContext.textFile('titanic.csv')  
  
# the file is of the format  
# Pclass,Survived,Name,Sex,Age,Sibsp,Parch,Ticket,Fare  
# Survived | Class | Name | Sex | Age | Siblings/Spouses Aboard | Parents/←  
# Children Aboard | Fare  
  
>>> titanic.take(2)  
['0,3,Mr. Owen Harris Braund,male,22,1,0,7.25',  
 '1,1,Mrs. John Bradley (Florence Briggs Thayer) Cumings,female,38,1,0,71.283']  
  
# note that each element is a single string - not particularly useful  
# one option is to first load the data into a numpy array
```

³<https://web.stanford.edu/class/archive/cs/cs109/cs109.1166/problem12.html>

```
>>> np_titanic = np.loadtxt('titanic.csv', delimiter=',', dtype=list)

# use sparkContext to parallelize the data into 4 partitions
>>> titanic_parallelize = spark.sparkContext.parallelize(np_titanic, 4)

>>> titanic_parallelize.take(2)
[array(['0', '3', ..., 'male', '22', '1', '0', '7.25'], dtype=object),
 array(['1', '1', ..., 'female', '38', '1', '0', '71.2833'], dtype=object)]

# end SparkSession
>>> spark.stop()
```

Achtung!

Because Apache Spark partitions and distributes data, calling for the first `n` objects using the same code (such as `take(n)`) may yield different results on different computers (or even each time you run it on one computer). This is not something you should worry about; it is the result of variation in partitioning and will not affect data analysis.

RDD Operations

Transformations

There are two types of operations you can perform on RDDs: transformations and actions. Transformations are functions that produce new RDDs from existing ones. Transformations are also lazy; they are not executed until an action is performed. This allows Spark to boost performance by optimizing how a sequence of transformations is executed at runtime.

One of the most commonly used transformations is the `map(func)`, which creates a new RDD by applying `func` to each element of the current RDD. This function, `func`, can be any callable python function, though it is often implemented as a `lambda` function. Similarly, `flatMap(func)` creates an RDD with the flattened results of `map(func)`.

```
# initialize your SparkSession object
>>> spark = SparkSession\
...     .builder\
...     .appName("app_name")\
...     .getOrCreate()

# use map() to format the data
>>> titanic = spark.sparkContext.textFile('titanic.csv')
>>> titanic.take(2)
['0,3,Mr. Owen Harris Braund,male,22,1,0,7.25',
 '1,1,Mrs. John Bradley (Florence Briggs Thayer) Cumings,female,38,1,0,71.283']

# apply split(',') to each element of the RDD with map()
>>> titanic.map(lambda row: row.split(','))\
```

```

...           .take(2)
[['0', '3', 'Mr. Owen Harris Braund', 'male', '22', '1', '0', '7.25'],
 ['1', '1', ..., 'female', '38', '1', '0', '71.283']]]

# compare to flatMap(), which flattens the results of each row
>>> titanic.flatMap(lambda row: row.split(','))\
...           .take(2)
['0', '3']

```

The `filter(func)` transformation returns a new RDD containing only the elements that satisfy `func`. In this case, `func` should be a callable python function that returns a Boolean. The elements of the RDD that evaluate to `True` are included in the new RDD while those that evaluate to `False` are excluded.

```

# create a new RDD containing only the female passengers
>>> titanic = titanic.map(lambda row: row.split(','))
>>> titanic_f = titanic.filter(lambda row: row[3] == 'female')
>>> titanic_f.take(3)
[['1', '1', ..., 'female', '38', '1', '0', '71.283'],
 ['1', '3', ..., 'female', '26', '0', '0', '7.925'],
 ['1', '1', ..., 'female', '35', '1', '0', '53.1']]

```

Note

A great transformation to help validate or explore your dataset is `distinct()`. This will return a new RDD containing only the distinct elements of the original. In the case of the Titanic dataset, if you did not know how many classes there were, you could do the following:

```

>>> titanic.map(lambda row: row[1])\
...           .distinct()\ \
...           .collect()
['1', '3', '2']

```

Spark Command	Transformation
<code>map(f)</code>	Returns a new RDD by applying <code>f</code> to each element of this RDD
<code>flatmap(f)</code>	Same as <code>map(f)</code> , except the results are flattened
<code>filter(f)</code>	Returns a new RDD containing only the elements that satisfy <code>f</code>
<code>distinct()</code>	Returns a new RDD containing the distinct elements of the original
<code>reduceByKey(f)</code>	Takes an RDD of <code>(key, val)</code> pairs and merges the values for each <code>key</code> using an associative and commutative reduce function <code>f</code>
<code>sortBy(f)</code>	Sorts this RDD by the given function <code>f</code>
<code>sortByKey(f)</code>	Sorts an RDD assumed to consist of <code>(key, val)</code> pairs by the given function <code>f</code>
<code>groupByKey(f)</code>	Returns a new RDD of groups of items based on <code>f</code>
<code>groupByKey()</code>	Takes an RDD of <code>(key, val)</code> pairs and returns a new RDD with <code>(key, (val1, val2, ...))</code> pairs

```
# the following counts the number of passengers in each class
# note that this isn't necessarily the best way to do this

# create a new RDD of (pclass, 1) elements to count occurrences
>>> pclass = titanic.map(lambda row: (row[1], 1))
>>> pclass.take(5)
[('3', 1), ('1', 1), ('3', 1), ('1', 1), ('3', 1)]

# count the members of each class
>>> pclass = pclass.reduceByKey(lambda x, y: x + y)
>>> pclass.collect()
[('3', 487), ('1', 216), ('2', 184)]

# sort by number of passengers in each class, ascending order
>>> pclass.sortBy(lambda row: row[1]).collect()
[('2', 184), ('1', 216), ('3', 487)]

# end SparkSession
>>> spark.stop()
```

Achtung!

Note that you must use `.collect()` to extract data from an RDD. Using `.collect()` will return an array.

Problem 1. Write a function that accepts the name of a text file with default `filename=huck_finn.txt`.^a Load the file as a PySpark RDD, and count the number of occurrences of each word. Sort the words by count, in descending order, and return a list of the `(word, count)` pairs for the 20 most used words. The data does not need to be cleaned.

^a<https://www.gutenberg.org/files/76/76-0.txt>

Actions

Actions are operations that return non-RDD objects. Two of the most common actions, `take(n)` and `collect()`, have already been seen above. The key difference between the two is that `take(n)` returns the first `n` elements from one (or more) partition(s) while `collect()` returns the contents of the entire RDD. When working with small datasets this may not be an issue, but for larger datasets running `collect()` can be very expensive.

Another important action is `reduce(func)`. Generally, `reduce()` combines (reduces) the data in each row of the RDD using `func` to produce some useful output. Note that `func` must be an associative and commutative binary operation; otherwise the results will vary depending on partitioning.

```
# create an RDD with the first million integers in 4 partitions
>>> ints = spark.sparkContext.parallelize(range(1, 1000001), 4)
# [1, 2, 3, 4, 5, ..., 1000000]
# sum the first one million integers
>>> ints.reduce(lambda x, y: x + y)
500000500000

# create a new RDD containing only survival data
>>> survived = titanic.map(lambda row: int(row[0]))
>>> survived.take(5)
[0, 1, 1, 1, 0]

# find total number of survivors
>>> survived.reduce(lambda x, y: x + y)
500
```

Spark Command	Action
<code>take(n)</code>	returns the first n elements of an RDD
<code>collect()</code>	returns the entire contents of an RDD
<code>reduce(f)</code>	merges the values of an RDD using an associative and commutative operator f
<code>count()</code>	returns the number of elements in the RDD
<code>min(); max(); mean()</code>	returns the minimum, maximum, or mean of the RDD, respectively
<code>sum()</code>	adds the elements in the RDD and returns the result
<code>saveAsTextFile(path)</code>	saves the RDD as a collection of text files (one for each partition) in the directory specified
<code>foreach(f)</code>	immediately applies f to each element of the RDD; not to be confused with <code>map()</code> , <code>foreach()</code> is useful for saving data somewhere not natively supported by PySpark

Problem 2. Since the area of a circle of radius r is $A = \pi r^2$, one way to estimate π is to estimate the area of the unit circle. A Monte Carlo approach to this problem is to uniformly sample points in the square $[-1, 1] \times [-1, 1]$ and then count the percentage of points that land within the unit circle. The percentage of points within the circle approximates the percentage of the area occupied by the circle. Multiplying this percentage by 4 (the area of the square $[-1, 1] \times [-1, 1]$) gives an estimate for the area of the circle.^a

Write a function that uses Monte Carlo methods to estimate the value of π . Your function should accept two keyword arguments: `n=10**5` and `parts=6`. Use `n*parts` sample points and partition your RDD with `parts` partitions. Return your estimate.

^aSee Example 7.1.1 in the Volume 2 textbook

DataFrames

While RDDs offer granular control, they can be slower than their Scala and Java counterparts when implemented in Python. The solution to this was the creation of a new data structure: Spark DataFrames. Just like RDDs, DataFrames are immutable distributed collections of objects; however, unlike RDDs, DataFrames are organized into named (and typed) columns. In this way they are conceptually similar to a relational database (or a pandas DataFrame).

The most important difference between a relational database and Spark DataFrames is in the execution of transformations and actions. When working with DataFrames, Spark's Catalyst Optimizer creates and optimizes a logical execution plan before sending any instructions to the drivers. After the logical plan has been formed, an optimal physical plan is created and executed. This provides significant performance boosts, especially when working with massive amounts of data. Since the Catalyst Optimizer functions the same across all language APIs, DataFrames bring performance parity to all of Spark's APIs.

Spark SQL and DataFrames

Creating a DataFrame from an existing text, csv, or JSON file is generally easier than creating an RDD. The DataFrame API also has arguments to deal with file headers or to automatically infer the schema.

```
# note that you should initialize your spark object first
# load the titanic dataset using default settings
>>> titanic = spark.read.csv('titanic.csv')
>>> titanic.show(2)
+---+---+-----+---+---+---+---+
|_c0|_c1|      _c2| _c3|_c4|_c5|_c6|      _c7|
+---+---+-----+---+---+---+---+-----+
|  0|  3|Mr. Owen Harris B...| male| 22|   1|   0|    7.25|
|  1|  1|Mrs. John Bradley...|female| 38|   1|   0| 71.2833|
+---+---+-----+---+---+---+---+-----+
only showing top 2 rows

# spark.read.csv('titanic.csv', inferSchema=True) will try to infer
# data types for each column

# load the titanic dataset specifying the schema
>>> schema = ('survived INT, pclass INT, name STRING, sex STRING, '
...             'age FLOAT, sibsp INT, parch INT, fare FLOAT'
...         )
>>> titanic = spark.read.csv('titanic.csv', schema=schema)
>>> titanic.show(2)
+-----+-----+-----+-----+-----+-----+
|survived|pclass|          name| sex|age|sibsp|parch|  fare|
+-----+-----+-----+-----+-----+-----+
|       0|     3|Mr. Owen Harris B...| male| 22|   1|   0|    7.25|
|       1|     1|Mrs. John Bradley...|female| 38|   1|   0| 71.2833|
+-----+-----+-----+-----+-----+-----+
only showing top 2 rows

# for files with headers, the following is convenient
spark.read.csv('my_file.csv', header=True, inferSchema=True)
```

Note

To convert a DataFrame to an RDD use `my_df.rdd`; to convert an RDD to a DataFrame use `spark.createDataFrame(my_rdd)`. You can also use `spark.createDataFrame()` on numpy arrays and pandas DataFrames.

DataFrames can be easily updated, queried, and analyzed using SQL operations. Spark allows you to run queries directly on DataFrames similar to how you perform transformations on RDDs. Additionally, the `pyspark.sql.functions` module contains many additional functions to further analysis. Below are many examples of basic DataFrame operations; further examples involving the `pyspark.sql.functions` module can be found in the additional materials section. Full documentation can be found at <https://spark.apache.org/docs/latest/api/python/pyspark.sql.html>.

```
# select data from the survived column
>>> titanic.select(titanic.survived).show(3) # or titanic.select("survived")
+-----+
|survived|
+-----+
|     0|
|     1|
|     1|
+-----+
only showing top 3 rows

# find all distinct ages of passengers (great for data exploration)
>>> titanic.select("age")\
...      .distinct()\
...      .show(3)
+---+
| age|
+---+
|18.0|
|64.0|
|0.42|
+---+
only showing top 3 rows

# filter the DataFrame for passengers between 20-30 years old (inclusive)
>>> titanic.filter(titanic.age.between(20, 30)).show(3)
+-----+-----+-----+-----+-----+-----+
|survived|pclass|          name| sex| age|sibsp|parch| fare|
+-----+-----+-----+-----+-----+-----+
|     0|     3|Mr. Owen Harris B...| male|22.0|   1|  0| 7.25|
|     1|     3|Miss. Laina Heikk...|female|26.0|   0|  0| 7.925|
|     0|     3| Mr. James Moran| male|27.0|   0|  0|8.4583|
+-----+-----+-----+-----+-----+-----+
only showing top 3 rows

# find total fare by pclass (or use .avg('fare') for an average)
>>> titanic.groupBy('pclass')\
...      .sum('fare')\
...      .show()
+-----+-----+
|pclass|sum(fare)|
+-----+-----+
```

```

|      1| 18177.41|
|      3|  6675.65|
|      2|  3801.84|
+-----+


# group and count by age and survival; order age/survival descending
>>> titanic.groupBy("age", "survived").count()\
...         .sort("age", "survived", ascending=False) \
...         .show(2)
+---+-----+
|age|survived|count|
+---+-----+
| 80|        1|     1|
| 74|        0|     1|
+---+-----+
only showing top 2 rows

# join two DataFrames on a specified column (or list of columns)
>>> titanic_cabins.show(3)
+-----+-----+
|           name|  cabin|
+-----+-----+
|Miss. Elisabeth W...|      B5|
|Master. Hudson Tr...|C22 C26|
|Miss. Helen Lorai...|C22 C26|
+-----+-----+
only showing top 3 rows

>>> titanic.join(titanic_cabins, on='name').show(3)
+-----+-----+-----+-----+-----+-----+-----+
|           name|survived|pclass|   sex|  age|sibsp|parch|  fare|  cabin|
+-----+-----+-----+-----+-----+-----+-----+
|Miss. Elisabeth W...|      0|     3| male|22.0|     1|     0|  7.25|      B5|
|Master. Hudson Tr...|      1|     3|female|26.0|     0|     0| 7.925|C22 C26|
|Miss. Helen Lorai...|      0|     3| male|27.0|     0|     0| 8.4583|C22 C26|
+-----+-----+-----+-----+-----+-----+-----+
only showing top 3 rows

```

Note

If you prefer to use traditional SQL syntax you can use `spark.sql("SQL QUERY")`. Note that this requires you to first create a temporary view of the DataFrame.

```

# create the temporary view so we can access the table through SQL
>>> titanic.createOrReplaceTempView("titanic")

```

```
# query using SQL syntax
>>> spark.sql("SELECT age, COUNT(*) AS count\
...           FROM titanic\
...           GROUP BY age\
...           ORDER BY age DESC").show(3)
+-----+
|age|count|
+---+---+
| 80|    1|
| 74|    1|
| 71|    2|
+---+---+
only showing top 3 rows
```

Spark SQL Command	SQLite Command
select(*cols)	SELECT
groupBy(*cols)	GROUP BY
sort(*cols, **kwargs)	ORDER BY
filter(condition)	WHERE
when(condition, value)	WHEN
between(lowerBound, upperBound)	BETWEEN
drop(*cols)	DROP
join(other, on=None, how=None)	JOIN (join type specified by how)
count()	COUNT()
sum(*cols)	SUM()
avg(*cols) or mean(*cols)	AVG()
collect()	fetchall()

Problem 3. Write a function with keyword argument `filename='titanic.csv'`. Load the file into a PySpark DataFrame and find (1) the number of women on-board, (2) the number of men on-board, (3) the survival rate of women, and (4) the survival rate of men. Return these four values in the order given as a tuple of floats.

Problem 4. In this problem, you will be using the `london_income_by_borough.csv` and the `london_crime_by_lsoa.csv` files to visualize the relationship between income and the frequency of crime.^a The former contains estimated mean and median income data for each London borough, averaged over 2008-2016; the first line of the file is a header with columns `borough`, `mean-08-16`, and `median-08-16`. The latter contains over 13 million lines of crime data, organized by borough and LSOA (Lower Super Output Area) code, for London between 2008 and 2016; the first line of the file is a header, containing the following seven columns:

`lsoa_code`: LSOA code (think area code) where the crime was committed
`borough`: London borough where the crime was committed
`major_category`: major or general category of the crime
`minor_category`: minor or specific category of the crime
`value`: number of occurrences of this crime in the given `lsoa_code`, `month`, and `year`
`year`: year the crime was committed
`month`: month the crime was committed

Write a function that accepts three keyword arguments:

`crimefile='london_crime_by_lsoa.csv'`, `incomefile='london_income_by_borough.csv'`, and `major_cat='Robbery'`. Load the two files as PySpark DataFrames. Use them to create a new DataFrame. The new DataFrame will contain a row for each borough and have columns for borough, total number of crimes for the given major category (`major_cat`), and median income. Order the DataFrame by the total number of crimes for `major_cat`, descending. The final DataFrame should have three columns: `borough`, `major_cat_total_crime`, and `median` -08-16 (column names may be different).

Convert the DataFrame to a numpy array using `np.array(df.collect())`, and create a scatter plot of the number of murders by the median income for each borough. Return the numpy array.

^adata.london.gov.uk

Machine Learning with Apache Spark

Apache Spark includes a vast and expanding ecosystem to perform machine learning. PySpark's primary machine learning API, `pyspark.ml`, is DataFrame-based.

Here we give a start to finish example using Spark ML to tackle the classic Titanic classification problem.

```
# prepare data
# convert the 'sex' column to binary categorical variable
>>> from pyspark.ml.feature import StringIndexer, OneHotEncoder
>>> sex_binary = StringIndexer(inputCol='sex', outputCol='sex_binary')

# one-hot-encode pclass (Spark automatically drops a column)
>>> onehot = OneHotEncoder(inputCols=['pclass'],
...                           outputCols=['pclass_onehot'])

# create single features column
from pyspark.ml.feature import VectorAssembler
features = ['sex_binary', 'pclass_onehot', 'age', 'sibsp', 'parch', 'fare']
features_col = VectorAssembler(inputCols=features, outputCol='features')

# now we create a transformation pipeline to apply the operations above
# this is very similar to the pipeline ecosystem in sklearn
>>> from pyspark.ml import Pipeline
>>> pipeline = Pipeline(stages=[sex_binary, onehot, features_col])
>>> titanic = pipeline.fit(titanic).transform(titanic)
```

```

# drop unnecessary columns for cleaner display (note the new columns)
>>> titanic = titanic.drop('pclass', 'name', 'sex')
>>> titanic.show(2)
+-----+-----+-----+-----+-----+-----+
|survived| age|sibsp|parch|fare|sex_binary|pclass_onehot|  features|
+-----+-----+-----+-----+-----+-----+-----+
|      0|22.0|     1|    0|7.25|      0.0| (3, [], [])|(8,[4,5...|
|      1|38.0|     1|    0|71.3|      1.0| (3, [1], ...|[0.0,1....|
+-----+-----+-----+-----+-----+-----+-----+

# split into train/test sets (75/25)
>>> train, test = titanic.randomSplit([0.75, 0.25], seed=11)

# initialize logistic regression
>>> from pyspark.ml.classification import LogisticRegression
>>> lr = LogisticRegression(labelCol='survived', featuresCol='features')

# run a train-validation-split to fit best elastic net param
# ParamGridBuilder constructs a grid of parameters to search over.
>>> from pyspark.ml.tuning import ParamGridBuilder, TrainValidationSplit
>>> from pyspark.ml.evaluation import MulticlassClassificationEvaluator as MCE
>>> paramGrid = ParamGridBuilder() \
...           .addGrid(lr.elasticNetParam, [0, 0.5, 1]).build()
# TrainValidationSplit will try all combinations and determine best model using
# the evaluator (see also CrossValidator)
>>> tvs = TrainValidationSplit(estimator=lr,
...                             estimatorParamMaps=paramGrid,
...                             evaluator=MCE(labelCol='survived'),
...                             trainRatio=0.75,
...                             seed=11)

# we train the classifier by fitting our tvs object to the training data
>>> clf = tvs.fit(train)

# use the best fit model to evaluate the test data
>>> results = clf.bestModel.evaluate(test)
>>> results.predictions.select(['survived', 'prediction']).show(5)
+-----+-----+
|survived|prediction|
+-----+-----+
|      0|      1.0|
|      0|      1.0|
|      0|      1.0|
|      0|      1.0|
|      0|      0.0|
+-----+-----+

# performance information is stored in various attributes of "results"

```

```

>>> results.accuracy
0.7527272727272727

>>> results.weightedRecall
0.7527272727272727

>>> results.weightedPrecision
0.751035147726004

# many classifiers do not have this object-oriented interface (yet)
# it isn't much more effort to generate the same statistics for a ←
    DecisionTreeClassifier, for example
>>> dt_clf = dt_tvs.fit(train) # same process, except for a different paramGrid

# generate predictions - this returns a new DataFrame
>>> preds = clf.bestModel.transform(test)
>>> preds.select('survived', 'probability', 'prediction').show(5)
+-----+-----+-----+
|survived|probability|prediction|
+-----+-----+-----+
|      0|   [1.0,0.0]|      0.0|
|      0|   [1.0,0.0]|      0.0|
|      0|   [1.0,0.0]|      0.0|
|      0|   [0.0,1.0]|      1.0|
+-----+-----+-----+

# initialize evaluator object
>>> dt_eval = MCE(labelCol='survived')
>>> dt_eval.evaluate(preds, {dt_eval.metricName: 'accuracy'})
0.8433179723502304

```

Below is a broad overview of the `pyspark.ml` ecosystem. It should help give you a starting point when looking for a specific functionality.

PySpark ML Module	Module Purpose
<code>pyspark.ml.feature</code>	provides functions to transform data into feature vectors
<code>pyspark.ml.tuning</code>	grid search, cross validation, and train/validation split functions
<code>pyspark.ml.evaluation</code>	tools to compute prediction metrics (accuracy, f1, etc.)
<code>pyspark.ml.classification</code>	classification models (logistic regression, SVM, etc.)
<code>pyspark.ml.clustering</code>	clustering models (k-means, Gaussian mixture, etc.)
<code>pyspark.ml.regression</code>	regression models (linear regression, decision tree regressor, etc.)

Problem 5. Write a function with keyword argument `filename='titanic.csv'`. Load the file into a PySpark DataFrame, and use the `pyspark.ml` package to train a classifier that outperforms the logistic regression each of the three metrics from the example above (`accuracy`, `weightedRecall`, `weightedPrecision`).

Some of Spark's available classifiers are listed below. For complete documentation, visit <https://spark.apache.org/docs/latest/api/python/pyspark.ml.html>.

```
# from pyspark.ml.classification import LinearSVC
#                                         DecisionTreeClassifier
#                                         GBTClassifier
#                                         MultilayerPerceptronClassifier
#                                         NaiveBayes
#                                         RandomForestClassifier
```

Use `randomSplit([0.75, 0.25], seed=11)` to split your data into train and test sets before fitting the model. Return the `accuracy`, `weightedRecall`, and `weightedPrecision` for your model, in the given order as a tuple.

Hint: to calculate the accuracy of a classifier in PySpark, use `accuracy = MCE(labelCol = 'survived', metricName='accuracy').evaluate(predictions)`.

Additional Material

Further DataFrame Operations

There are a few other functions built directly on top of DataFrames to further analysis. Additionally, the `pyspark.sql.functions` module expands the available functions significantly.⁴

```
# some immediately accessible functions
# covariance between pclass and fare
>>> titanic.cov('pclass', 'fare')
-22.86289824115662

# summary of statistics for selected columns
>>> titanic.select("pclass", "age", "fare")\
...     .summary().show()
+-----+-----+-----+
|summary|      pclass|        age|       fare|
+-----+-----+-----+
|  count|      887|      887|      887|
|  mean| 2.305524239007892|29.471443066501564|32.305420253026846|
| stddev|0.8366620036697728|14.121908405492908| 49.78204096767521|
|   min|          1|      0.42|      0.0|
|  25%|          2|      20.0|      7.925|
|  50%|          3|      28.0|      14.4542|
|  75%|          3|      38.0|      31.275|
|   max|          3|      80.0|      512.3292|
+-----+-----+-----+

# additional functions from the functions module
>>> from pyspark.sql import functions as sqlf

# finding the mean of a column without grouping requires sqlf.avg()
# alias(new_name) allows us to rename the column on the fly
>>> titanic.select(sqlf.avg("age").alias("Average Age")).show()
+-----+
|    Average Age|
+-----+
|29.471443066516347|
+-----+

# use .agg([dict]) on GroupedData to specify [multiple] aggregate
# functions, including those from pyspark.sql.functions
>>> titanic.groupBy('pclass')\
...     .agg({'fare': 'var_samp', 'age': 'stddev'})\
...     .show(3)
+-----+-----+
|pclass| var_samp(fare)| stddev(age)|
+-----+-----+
```

⁴<https://spark.apache.org/docs/latest/api/python/pyspark.sql.html>

```

| 1| 6143.483042924841|14.183632587264817|
| 3|139.64879027298073|12.095083834183779|
| 2|180.02658999396826|13.756191206499766|
+-----+-----+
# perform multiple aggregate actions on the same column
>>> titanic.groupBy('pclass')\
...     .agg(sqlf.sum('fare'), sqlf.stddev('fare'))\
...     .show()
+-----+-----+
|pclass|      sum(fare)| stddev_samp(fare) |
+-----+-----+
| 1|18177.412506103516| 78.38037409278448|
| 3| 6675.653553009033| 11.81730892686574|
| 2|3801.8417053222656|13.417398778972332|
+-----+-----+

```

pyspark.sql.functions	Operation
ceil(col)	computes the ceiling of each element in col
floor(col)	computes the floor of each element in col
min(col), max(col)	returns the minimum/maximum value of col
mean(col)	returns the average of the values of col
stddev(col)	returns the unbiased sample standard deviation of col
var_samp(col)	returns the unbiased variance of the values in col
rand(seed=None)	generates a random column with i.i.d. samples from [0, 1]
randn(seed=None)	generates a random column with i.i.d. samples from the standard normal distribution
exp(col)	computes the exponential of col
log(arg1, arg2=None)	returns arg1-based logarithm of arg2; if there is only one argument, then it returns the natural logarithm
cos(col), sin(col), etc.	computes the given trigonometric or inverse trigonometric (asin(col), etc.) function of col

14 Web Crawling

Lab Objective: Gathering data from the internet often requires information from several web pages. In this lab, we present two methods for crawling through multiple web pages without violating copyright laws or straining the load on a server. We also demonstrate how to scrape data from asynchronously loaded web pages and how to interact programmatically with web pages when needed.

Scraping Etiquette

There are two main ways that web scraping can be problematic for a website owner.

1. The scraper doesn't respect the website's terms and conditions or gathers private or proprietary data.
2. The scraper imposes too much extra server load by making requests too often or in quick succession.

These are extremely important considerations in any web scraping program. Scraping copyrighted information without the consent of the copyright owner can have severe legal consequences. Many websites, in their terms and conditions, prohibit scraping parts or all of the site. Websites that do allow scraping usually have a file called `robots.txt` (for example, www.google.com/robots.txt) that specifies which parts of the website are off-limits, and how often requests can be made according to the robots exclusion standard.¹

Achtung!

Be careful and considerate when doing any sort of scraping, and take care when writing and testing code to avoid unintended behavior. It is up to the programmer to create a scraper that respects the rules found in the terms and conditions and in `robots.txt`. Make sure to scrape websites legally.

¹See www.robotstxt.org/orig.html and en.wikipedia.org/wiki/Robots_exclusion_standard.

Recall that consecutive requests without pauses can strain a website's server and provoke retaliation. Most servers are designed to identify such scrapers, block their access, and sometimes even blacklist the user. This is especially common in smaller websites that aren't built to handle enormous amounts of traffic. To briefly pause the program between requests, use `time.sleep()`.

```
>>> import time
>>> time.sleep(3)      # Pause execution for 3 seconds.
```

The amount of necessary wait time depends on the website. Sometimes, `robots.txt` contains a `Crawl-delay` directive which gives a number of seconds to wait between successive requests. If this doesn't exist, pausing for a half-second to a second between requests is typically sufficient. An email to the site's webmaster is always the safest approach and may be necessary for large scraping operations.

Python provides a parsing library called `urllib.robotparser` for reading `robot.txt` files. Below is an example of using this library to check where robots are allowed on arxiv.org. A website's `robots.txt` file will often include different instructions for specific crawlers. These crawlers are identified by a `User-agent` string. For example, Google's webcrawler, `User-agent` Googlebot, may be directed to index only the pages the website wants to have listed on a Google search. We will use the default `User-agent`, `"*"`.

```
>>> from urllib import robotparser
>>> rp = robotparser.RobotFileParser()
# Set the URL for the robots.txt file. Note that the URL contains `robots.txt'
>>> rp.set_url("https://arxiv.org/robots.txt")
>>> rp.read()
# Request the crawl-delay time for the default User-agent
>>> rp.crawl_delay("*")
15
# Check if User-agent "*" can access the page
>>> rp.can_fetch("*", "https://arxiv.org/archive/math/")
True
>>> rp.can_fetch("*", "https://arxiv.org/IgnoreMe/")
False
```

Problem 1. Write a program that accepts a web address defaulting to the site `http://example.webscraping.com` and a list of pages defaulting to `["/", "/trap", "/places/default/search"]`. For each page, check if the `robots.txt` file permits access. Return a list of boolean values corresponding to each page. Also return the crawl delay time.

Crawling Through Multiple Pages

While web scraping refers to the actual gathering of web-based data, web crawling refers to the navigation of a program between webpages. Web crawling allows a program to gather related data from multiple web pages and websites.

Consider `books.toscrape.com`, a site to practice web scraping that mimics a bookstore. The page `http://books.toscrape.com/catalogue/category/books/mystery_3/index.html` lists mystery books with overall ratings and review. More mystery books can be accessed by clicking on the `next` link. The following example demonstrates how to navigate between webpages to collect all of the mystery book titles.

```
def scrape_books(start_page = "index.html"):
    """ Crawl through http://books.toscrape.com and extract mystery titles"""

    # Initialize variables, including a regex for finding the 'next' link.
    base_url="http://books.toscrape.com/catalogue/category/books/mystery_3/"
    titles = []
    page = base_url + start_page                      # Complete page URL.
    next_page_finder = re.compile(r"next")             # We need this button.

    current = None

    for _ in range(4):
        while current == None: # Try downloading until it works.
            # Download the page source and PAUSE before continuing.
            page_source = requests.get(page).text
            time.sleep(1)          # PAUSE before continuing.
            soup = BeautifulSoup(page_source, "html.parser")
            current = soup.find_all(class_="product_pod")

        # Navigate to the correct tag and extract title
        for book in current:
            titles.append(book.h3.a["title"])

        # Find the URL for the page with the next data.
        if "page-4" not in page:
            new_page = soup.find(string=next_page_finder).parent["href"]
            page = base_url + new_page      # New complete page URL.
            current = None
    return titles
```

In this example, the `for` loop cycles through the pages of books, and the `while` loop ensures that each website page loads properly: if the downloaded `page_source` doesn't have a tag whose class is `product_pod`, the request is sent again. After recording all of the titles, the function locates the link to the next page. This link is stored in the HTML as a relative website path (`page-2.html`); the complete URL to the next day's page is the concatenation of the base URL `http://books.toscrape.com/catalogue/category/books/mystery_3/` with this relative link.

Problem 2. Modify `scrape_books()` so that it gathers the price for each fiction book and returns the mean price, in £, of a fiction book.

Asynchronously Loaded Content and User Interaction

Web crawling with the methods presented in the previous section fails under a few circumstances. First, many webpages use JavaScript, the standard client-side scripting language for the web, to load portions of their content asynchronously. This means that at least some of the content isn't initially accessible through the page's source code (for example, if you have to scroll down to load more results). Second, some pages require user interaction, such as clicking buttons which aren't links (tags which contain a URL that can be loaded) or entering text into form fields (like search bars).

The Selenium framework provides a solution to both of these problems. Originally developed for writing unit tests for web applications, Selenium allows a program to open a web browser and interact with it in the same way that a human user would, including clicking and typing. It also has BeautifulSoup-esque tools for searching the HTML source of the current page.

Note

Selenium requires an executable driver file for each kind of browser. The following examples use Google Chrome, but Selenium supports Firefox, Internet Explorer, Safari, Opera, and PhantomJS (a special browser without a user interface). See <https://seleniumhq.github.io/selenium/docs/api/py> or <http://selenium-python.readthedocs.io/installation.html> for installation instructions and driver download instructions.

If your program still can't find the driver after you've downloaded it, add the argument `executable_path = "path/to driver/file"` when you call `webdriver`. If this doesn't work, you may need to add the location to your system PATH. On a Mac, open the file `/etc/path` and add the new location. On Linux, add `export PATH="path/to driver/file:$PATH"` to the file `./bashrc`. For Windows, follow a tutorial such as this one: <https://www.architectryan.com/2018/03/17/add-to-the-path-on-windows-10/>.

To use Selenium, start up a browser using one of the drivers in `selenium.webdriver`. The browser has a `get()` method for going to different web pages, a `page_source` attribute containing the HTML source of the current page, and a `close()` method to exit the browser.

```
>>> from selenium import webdriver

# Start up a browser and go to example.com.
>>> browser = webdriver.Chrome()
>>> browser.get("https://www.example.com")

# Feed the HTML source code for the page into BeautifulSoup for processing.
>>> soup = BeautifulSoup(browser.page_source, "html.parser")
>>> print(soup.prettify())
<!DOCTYPE html>
<html xmlns="http://www.w3.org/1999/xhtml">
<head>
<title>
    Example Domain
</title>
<meta charset="utf-8"/>
```

```
<meta content="text/html; charset=utf-8" http-equiv="Content-type"/>
# ...

>>> browser.close()                                # Close the browser.
```

Selenium can deliver the HTML page source to BeautifulSoup, but it also has its own tools for finding tags in the HTML.

Method	Returns
<code>find_element_by_tag_name()</code>	The first tag with the given name
<code>find_element_by_name()</code>	The first tag with the specified <code>name</code> attribute
<code>find_element_by_class_name()</code>	The first tag with the given <code>class</code> attribute
<code>find_element_by_id()</code>	The first tag with the given <code>id</code> attribute
<code>find_element_by_link_text()</code>	The first tag with a matching <code>href</code> attribute
<code>find_element_by_partial_link_text()</code>	The first tag with a partially matching <code>href</code> attribute

Table 14.1: Methods of the `selenium.webdriver.Chrome` class.

Each of the `find_element_by_*`() methods returns a single object representing a web element (of type `selenium.webdriver.remote.webelement.WebElement`), much like a BeautifulSoup tag (of type `bs4.element.Tag`). If no such element can be found, a Selenium `NoSuchElementException` is raised. If you want to find more than just the first matching object, each webdriver also has several `find_elements_by_*`() methods (elements, plural) that return a list of all matching elements, or an empty list if there are no matches.

Web element objects have methods that allow the program to interact with them: `click()` sends a click, `send_keys()` enters in text, and `clear()` deletes existing text. This functionality makes it possible for Selenium to interact with a website in the same way that a human would. For example, the following code opens up <https://www.google.com>, types “Python Selenium Docs” into the search bar, and hits enter.

```
>>> from selenium.webdriver.common.keys import Keys
>>> from selenium.common.exceptions import NoSuchElementException

>>> browser = webdriver.Chrome()
>>> try:
...     browser.get("https://www.google.com")
...     try:
...         # Get the search bar, type in some text, and press Enter.
...         search_bar = browser.find_element_by_name('q')
...         search_bar.clear()                      # Clear any pre-set text.
...         search_bar.send_keys("Python Selenium Docs")
...         search_bar.send_keys(Keys.RETURN)      # Press Enter.
...     except NoSuchElementException:
...         print("Could not find the search bar!")
...     raise
... finally:
...     browser.close()
...
```

Problem 3. The website IMDB contains a variety of information on movies. Specifically, information on the top 10 box office movies of the week can be found at <https://www.imdb.com/chart/boxoffice>. Using BeautifulSoup, Selenium, or both, return a list, with each title on a new row, of the top 10 movies of the week and order the list according to the total grossing of the movies, from most money to the least. Break ties using the weekend gross, from most money to the least.

Using CSS Selectors

In addition to the methods listed in Table 14.1, you can also use CSS or XPath selectors to interact more precisely with the page source. Refer to Table 3 from the WebScraping lab for a review of CSS syntax. The following code searches Google for “Python Selenium Docs” and then clicks on the second result.

```
#As before, go to Google and type in the search bar,
# but this time we use CSS selectors

>>> from selenium.webdriver.common.keys import Keys
>>> from selenium.common.exceptions import NoSuchElementException

>>> browser = webdriver.Chrome()
>>> try:
...     browser.get("https://google.com")
...     try:
...         search_bar = browser.find_element_by_css_selector(
...             "input[name='q']")
...         search_bar.clear()
...         search_bar.send_keys("Python Selenium Docs")
...         search_bar.send_keys(Keys.RETURN)
...         try:
...             # Wait a second, then get the second search result
...             time.sleep(1)
...             # "+ div" returns the element's next sibling with a "div" tag
...             second_result = browser.find_element_by_css_selector(
...                 "div[class='g'] + div")
...             try:
...                 # Get the link, which is a child of second_result
...                 link = second_result.find_element_by_css_selector(
...                     "div[class='r']")
...                 link.click()
...                 time.sleep(1)

...             #Remember to handle exceptions
...             except NoSuchElementException:
...                 print("Could not find link")
...             except NoSuchElementException:
...                 print("Could not find second result")
```

```

...     except NoSuchElementException:
...         print("Could not find the search bar")
... finally:
...     browser.close()

```

In the above example, we could have used `find_element_by_class_name()`, but when you need more precision than that, CSS selectors can be very useful. Remember that to view specific HTML associated with an object in Chrome or Firefox, you can right click on the object and click “Inspect.” For Safari, you need to first enable “Show Develop menu” in “Preferences” under “Advanced.” Keep in mind that you can also search through the source code (ctrl+f or cmd+f) to make sure you’re using a unique identifier.

Note

Using Selenium to access a page’s source code is typically much safer, though slower, than using `requests.get()`, since Selenium waits for each web page to load before proceeding. For instance, some websites are somewhat defensive about scrapers, but Selenium can sometimes make it possible to gather info without offending the administrators.

Problem 4. Project Euler (<https://projecteuler.net>) is a collection of mathematical computing problems. Each problem is listed with an ID, a description/title, and the number of users that have solved the problem.

Using Selenium, BeautifulSoup, or both, record the number of people who have solved each of the 700+ problems in the archive at <https://projecteuler.net/archives>. Plot the number of people who have solved each problem against the problem IDs, using a log scale for the y-axis. Display the scatter plot, then state the IDs of which problems have been solved most and least number of times.

Problem 5. The website <http://example.webscraping.com> contains a list of countries of the world. Using Selenium, go to the search page, enter the letters “ca”, and hit `enter`. Remember to use the crawl delay time you found in Problem 1 so you don’t send your requests too fast. Gather the `href` links associated with the `<a>` tags of all 10 displayed results. Print each link on a different line.

15 Web Scraping

Lab Objective: Web Scraping is the process of gathering data from websites on the internet. Since almost everything rendered by an internet browser as a web page uses HTML, the first step in web scraping is being able to extract information from HTML. In this lab, we introduce the requests library for scraping web pages, and BeautifulSoup, Python's canonical tool for efficiently and cleanly navigating and parsing HTML.

HTTP and Requests

HTTP stands for Hypertext Transfer Protocol, which is an application layer networking protocol. It is a higher level protocol than TCP, which we used to build a server in the Web Technologies lab, but uses TCP protocols to manage connections and provide network capabilities. The HTTP protocol is centered around a request and response paradigm, in which a client makes a request to a server and the server replies with a response. There are several methods, or requests, defined for HTTP servers, the three most common of which are GET, POST, and PUT. A GET request asks for information from the server, a POST request modifies the state of the server, and a PUT request adds new pieces of data to the server.

The standard way to get the source code of a website using Python is via the `requests` library.¹ Calling `requests.get()` sends an HTTP GET request to a specified website. The website returns a response code, which indicates whether or not the request was received, understood, and accepted. If the response code is good, typically 200², then the response will also include the website source code as an HTML file.

```
>>> import requests

# Make a request and check the result. A status code of 200 is good.
>>> response = requests.get("http://www.byu.edu")
>>> print(response.status_code, response.ok, response.reason)
200 True OK
```

¹Though `requests` is not part of the standard library, it is recognized as a standard tool in the data science community. See <http://docs.python-requests.org/>.

²See https://en.wikipedia.org/wiki/List_of_HTTP_status_codes for explanation of specific response codes.

```
# The HTML of the website is stored in the 'text' attribute.  
>>> print(response.text)  
<!DOCTYPE html>  
<html lang="en" dir="ltr" prefix="content: http://purl.org/rss/1.0/modules/←  
    content/ dc: http://purl.org/dc/terms/ foaf: http://xmlns.com/foaf/0.1/ ←  
    og: http://ogp.me/ns# rdfs: http://www.w3.org/2000/01/rdf-schema# schema:←  
        http://schema.org/ sioc: http://rdfs.org/sioc/ns# sioct: http://rdfs.org/←  
        /sioc/types# skos: http://www.w3.org/2004/02/skos/core# xsd: http://www.←  
        w3.org/2001/XMLSchema# " class=" ">  
  
<head>  
    <meta charset="utf-8" />  
# ...
```

Note that some websites aren't built to handle large amounts of traffic or many repeated requests. Most are built to identify web scrapers or crawlers that initiate many consecutive GET requests without pauses, and retaliate or block them. When web scraping, always make sure to store the data that you receive in a file and include error checks to prevent retrieving the same data unnecessarily. We won't spend much time on that in this lab, but it's especially important in larger applications.

Problem 1. Use the `requests` library to get the HTML source for the website `http://www.example.com`. Save the source as a file called `example.html`. If the file already exists, make sure not to scrape the website or overwrite the file. You will use this file later in the lab.

Achtung!

Scraping copyrighted information without the consent of the copyright owner can have severe legal consequences. Many websites, in their terms and conditions, prohibit scraping parts or all of the site. Websites that do allow scraping usually have a file called `robots.txt` (for example, `www.google.com/robots.txt`) that specifies which parts of the website are off-limits and how often requests can be made according to the robots exclusion standard.^a

Be careful and considerate when doing any sort of scraping, and take care when writing and testing code to avoid unintended behavior. It is up to the programmer to create a scraper that respects the rules found in the terms and conditions and in `robots.txt`.^b

We will cover this more in the next lab.

^aSee www.robotstxt.org/orig.html and en.wikipedia.org/wiki/Robots_exclusion_standard.

^bPython provides a parsing library called `urllib.robotparser` for reading `robot.txt` files. For more information, see <https://docs.python.org/3/library/urllib.robotparser.html>.

HTML

Hyper Text Markup Language, or HTML, is the standard markup language—a language designed for the processing, definition, and presentation of text—for creating webpages. It structures a document using pairs of tags that surround and define content. Opening tags have a tag name surrounded by angle brackets (`<tag-name>`). The companion closing tag looks the same, but with a forward slash before the tag name (`</tag-name>`). A list of all current HTML tags can be found at <http://htmldog.com/reference/htmltags>.

Most tags can be combined with attributes to include more data about the content, help identify individual tags, and make navigating the document much simpler. In the following example, the `<a>` tag has `id` and `href` attributes.

```
<html>                                <!-- Opening tags -->
  <body>
    <p>
      Click <a id='info' href='http://www.example.com'>here</a>
      for more information.
    </p>                               <!-- Closing tags -->
  </body>
</html>
```

In HTML, `href` stands for hypertext reference, a link to another website. Thus the above example would be rendered by a browser as a single line of text, with `here` being a clickable link to <http://www.example.com>:

Click here for more information.

Unlike Python, HTML does not enforce indentation (or any whitespace rules), though indentation generally makes HTML more readable. The previous example can be written in a single line.

```
<html><body><p>Click <a id='info' href='http://www.example.com/info'>here</a>
  for more information.</p></body></html>
```

Special tags, which don't contain any text or other tags, are written without a closing tag and in a single pair of brackets. A forward slash is included between the name and the closing bracket. Examples of these include `<hr/>`, which describes a horizontal line, and ``, the tag for representing an image.

Note

You can open .html files using a text editor or any web browser. In a browser, you can inspect the source code associated with specific elements. Right click the element and select **Inspect**. If you are using Safari, you may first need to enable “Show Develop menu” in “Preferences” under the “Advanced” tab.

BeautifulSoup

BeautifulSoup (`bs4`) is a package³ that makes it simple to navigate and extract data from HTML documents. See <http://www.crummy.com/software/BeautifulSoup/bs4/doc/index.html> for the full documentation.

The `bs4.BeautifulSoup` class accepts two parameters to its constructor: a string of HTML code and an HTML parser to use under the hood. The HTML parser is technically a keyword argument, but the constructor prints a warning if one is not specified. The standard choice for the parser is "`html.parser`", which means the object uses the standard library's `html.parser` module as the engine behind the scenes.

Note

Depending on project demands, a parser other than "`html.parser`" may be useful. A couple of other options are "`lxml`", an extremely fast parser written in C, and "`html5lib`", a slower parser that treats HTML in much the same way a web browser does, allowing for irregularities. Both must be installed independently; see <https://www.crummy.com/software/BeautifulSoup/bs4/doc/#installing-a-parser> for more information.

A BeautifulSoup object represents an HTML document as a tree. In the tree, each tag is a node with nested tags and strings as its children. The `prettyify()` method returns a string that can be printed to represent the BeautifulSoup object in a readable format that reflects the tree structure.

```
>>> from bs4 import BeautifulSoup

>>> small_example_html = """
<html><body><p>
    Click <a id='info' href='http://www.example.com'>here</a>
    for more information.
</p></body></html>
"""

>>> small_soup = BeautifulSoup(small_example_html, 'html.parser')
>>> print(small_soup.prettify())
<html>
  <body>
    <p>
      Click
      <a href="http://www.example.com" id="info">
        here
      </a>
      for more information.
    </p>
  </body>
</html>
```

³BeautifulSoup is not part of the standard library; install it with `conda install beautifulsoup4` or with `pip install beautifulsoup4`.

Each tag in a `BeautifulSoup` object's HTML code is stored as a `bs4.element.Tag` object, with actual text stored as a `bs4.element.NavigableString` object. Tags are accessible directly through the `BeautifulSoup` object.

```
# Get the <p> tag (and everything inside of it).
>>> small_soup.p
<p>
    Click <a href="http://www.example.com" id="info">here</a>
    for more information.
</p>

# Get the <a> sub-tag of the <p> tag.
>>> a_tag = small_soup.p.a
>>> print(a_tag, type(a_tag), sep='\n')
<a href="http://www.example.com" id="info">here</a>
<class 'bs4.element.Tag'>

# Get just the name, attributes, and text of the <a> tag.
>>> print(a_tag.name, a_tag.attrs, a_tag.string, sep="\n")
a
{'id': 'info', 'href': 'http://www.example.com'}
here
```

Attribute	Description
<code>name</code>	The name of the tag
<code>attrs</code>	A dictionary of the attributes
<code>string</code>	The single string contained in the tag
<code>strings</code>	Generator for strings of children tags
<code>stripped_strings</code>	Generator for strings of children tags, stripping whitespace
<code>text</code>	Concatenation of strings from all children tags

Table 15.1: Data attributes of the `bs4.element.Tag` class.

Problem 2. The `BeautifulSoup` class has a `find_all()` method that, when called with `True` as the only argument, returns a list of all tags in the HTML source code.

Write a function that accepts a string of HTML code as an argument. Use `BeautifulSoup` to return a list of the `names` of the tags in the code.

Navigating the Tree Structure

Not all tags are easily accessible from a `BeautifulSoup` object. Consider the following example.

```
>>> pig_html = """
<html><head><title>Three Little Pigs</title></head>
```

```

<body>
<p class="title"><b>The Three Little Pigs</b></p>
<p class="story">Once upon a time, there were three little pigs named
<a href="http://example.com/larry" class="pig" id="link1">Larry,</a>
<a href="http://example.com/mo" class="pig" id="link2">Mo</a>, and
<a href="http://example.com/curly" class="pig" id="link3">Curly.</a>
<p>The three pigs had an odd fascination with experimental construction.</p>
<p>...</p>
</body></html>
"""

>>> pig_soup = BeautifulSoup(pig_html, "html.parser")
>>> pig_soup.p
<p class="title"><b>The Three Little Pigs</b></p>

>>> pig_soup.a
<a class="pig" href="http://example.com/larry" id="link1">Larry,</a>

```

Since the HTML in this example has several `<p>` and `<a>` tags, only the **first** tag of each name is accessible directly from `pig_soup`. The other tags can be accessed by manually navigating through the HTML tree.

Every HTML tag (except for the topmost tag, which is usually `<html>`) has a parent tag. Each tag also has zero or more sibling and children tags or text. Following a true tree structure, every `bs4.element.Tag` in a soup has multiple attributes for accessing or iterating through parent, sibling, or child tags.

Attribute	Description
<code>parent</code>	The parent tag
<code>parents</code>	Generator for the parent tags up to the top level
<code>next_sibling</code>	The tag immediately after to the current tag
<code>next_siblings</code>	Generator for sibling tags after the current tag
<code>previous_sibling</code>	The tag immediately before the current tag
<code>previous_siblings</code>	Generator for sibling tags before the current tag
<code>contents</code>	A list of the immediate children tags
<code>children</code>	Generator for immediate children tags
<code>descendants</code>	Generator for all children tags (recursively)

Table 15.2: Navigation attributes of the `bs4.element.Tag` class.

```

# Start at the first <a> tag in the soup.
>>> a_tag = pig_soup.a
>>> a_tag
<a class="pig" href="http://example.com/larry" id="link1">Larry,</a>

# Get the names of all of <a>'s parent tags, traveling up to the top.
# The name '[document]' means it is the top of the HTML code.
>>> [par.name for par in a_tag.parents]      # <a>'s parent is <p>, whose
['p', 'body', 'html', '[document]']          # parent is <body>, and so on.

```

```
# Get the next siblings of <a>.
>>> a_tag.next_sibling
'\n'                                         # The first sibling is just text.
>>> a_tag.next_sibling.next_sibling          # The second sibling is a tag.
<a class="pig" href="http://example.com/mo" id="link2">Mo</a>
```

Note carefully that newline characters are considered to be children of a parent tag. Therefore iterating through children or siblings often requires checking which entries are tags and which are just text. In the next example, we use a tag's `attrs` attribute to access specific attributes within the tag (see Table 15.1).

```
# Get to the <p> tag that has class="story" using these commands.
>>> p_tag = pig_soup.body.p.next_sibling.next_sibling
>>> p_tag.attrs["class"]                      # Make sure it's the right tag.
['story']

# Iterate through the child tags of <p> and print hrefs whenever they exist.
>>> for child in p_tag.children:
...     # Skip the children that are not bs4.element.Tag objects
...     # These don't have the attribute "attrs"
...     if hasattr(child, "attrs") and "href" in child.attrs:
...         print(child.attrs["href"])
http://example.com/larry
http://example.com/mo
http://example.com/curly
```

Note that the `"class"` attribute of the `<p>` tag is a list. This is because the `"class"` attribute can take on several values at once; for example, the tag `<p class="story book">` is of class `'story'` and of class `'book'`.

The behavior of the `string` attribute of a `bs4.element.Tag` object depends on the structure of the corresponding HTML tag.

1. If the tag has a string of text and no other child elements, then `string` is just that text.
2. If the tag has exactly one child tag and the child tag has only a string of text, then the tag has the same `string` as its child tag.
3. If the tag has more than one child, then `string` is `None`. In this case, use `strings` to iterate through the child strings. Alternatively, the `get_text()` method returns all text belonging to a tag and to all of its descendants. In other words, it returns anything inside a tag that isn't another tag.

```
>>> pig_soup.head
<head><title>Three Little Pigs</title></head>

# Case 1: the <title> tag's only child is a string.
>>> pig_soup.head.title.string
'Three Little Pigs'
```

```
# Case 2: The <head> tag's only child is the <title> tag.
>>> pig_soup.head.string
'Three Little Pigs'

# Case 3: the <body> tag has several children.
>>> pig_soup.body.string is None
True
>>> print(pig_soup.body.get_text().strip())
The Three Little Pigs
Once upon a time, there were three little pigs named
Larry,
Mo, and
Curly.
The three pigs had an odd fascination with experimental construction.
...
```

Problem 3. Write a function that reads a file of the same format as the output from Problem 1 and loads it into BeautifulSoup. Find the first `<a>` tag, and return its text along with a boolean value indicating whether or not it has a hyperlink (`href` attribute).

Searching for Tags

Navigating the HTML tree manually can be helpful for gathering data out of lists or tables, but these kinds of structures are usually buried deep in the tree. The `find()` and `find_all()` methods of the `BeautifulSoup` class identify tags that have distinctive characteristics, making it much easier to jump straight to a desired location in the HTML code. The `find()` method only returns the `first` tag that matches a given criteria, while `find_all()` returns a list of all matching tags. Tags can be matched by name, attributes, and/or text.

```
# Find the first <b> tag in the soup.
>>> pig_soup.find(name='b')
<b>The Three Little Pigs</b>

# Find all tags with a class attribute of 'pig'.
# Since 'class' is a Python keyword, use 'class_=' as the argument.
>>> pig_soup.find_all(class_="pig")
[<a class="pig" href="http://example.com/larry" id="link1">Larry,</a>,
 <a class="pig" href="http://example.com/mo" id="link2">Mo</a>,
 <a class="pig" href="http://example.com/curly" id="link3">Curly.</a>]

# Find the first tag that matches several attributes.
>>> pig_soup.find(attrs={"class": "pig", "href": "http://example.com/mo"})
<a class="pig" href="http://example.com/mo" id="link2">Mo</a>

# Find the first tag whose text is 'Mo'.
```

```
>>> pig_soup.find(string='Mo')
'Mo'                                     # The result is the actual string,
>>> pig_soup.find(string='Mo').parent    # so go up one level to get the tag.
<a class="pig" href="http://example.com/mo" id="link2">Mo</a>
```

Problem 4. The file `san_diego_weather.html` contains the HTML source for an old page from Weather Underground.^a Write a function that reads the file and loads it into BeautifulSoup.

Return a list of the following tags:

1. The tag containing the date “Thursday, January 1, 2015”.
2. The tags which contain the `links` “Previous Day” and “Next Day.”
3. The tag which contains the number associated with the Actual Max Temperature.

^aSee http://www.wunderground.com/history/airport/KSAN/2015/1/1/DailyHistory.html?req_city=San+Diego&req_state=CA&req_statename=California&reqdb.zip=92101&reqdb.magic=1&reqdb.wmo=99999&MR=1

Advanced Search Techniques: Regular Expressions

Consider the problem of finding the tag that is a link to the URL `http://example.com/curly`.

```
>>> pig_soup.find(href="http://example.com/curly")
<a class="pig" href="http://example.com/curly" id="link3">Curly.</a>
```

This approach works, but it requires entering in the entire URL. To perform generalized searches, the `find()` and `find_all()` method also accept compiled regular expressions from the `re` module. This way, the methods locate tags whose name, attributes, and/or string matches a pattern.

```
>>> import re

# Find the first tag with an href attribute containing 'curly'.
>>> pig_soup.find(href=re.compile(r"curly"))
<a class="pig" href="http://example.com/curly" id="link3">Curly.</a>

# Find the first tag with a string that starts with 'Cu'.
>>> pig_soup.find(string=re.compile(r"^Cu")).parent
<a class="pig" href="http://example.com/curly" id="link3">Curly.</a>

# Find all tags with text containing 'Three'.
>>> [tag.parent for tag in pig_soup.find_all(string=re.compile(r"Three"))]
[<title>Three Little Pigs</title>, <b>The Three Little Pigs</b>]
```

Finally, to find a tag that has a particular attribute, regardless of the actual value of the attribute, use `True` in place of search values.

```
# Find all tags with an 'id' attribute.
>>> pig_soup.find_all(id=True)
[<a class="pig" href="http://example.com/larry" id="link1">Larry,</a>,
 <a class="pig" href="http://example.com/mo" id="link2">Mo</a>,
 <a class="pig" href="http://example.com/curly" id="link3">Curly.</a>]

# Find the names all tags WITHOUT an 'id' attribute.
>>> [tag.name for tag in pig_soup.find_all(id=False)]
['html', 'head', 'title', 'body', 'p', 'b', 'p', 'p']
```

Advanced Search Techniques: CSS Selectors

`BeautifulSoup` also supports the use of CSS selectors. CSS (Cascading Style Sheet) describes the style and layout of a webpage, and CSS selectors provide a useful way to navigate HTML code. Use the method `soup.select()` to find all elements matching an argument. The general format for an argument is `tag-name[attribute-name = 'attribute value']`. The table below lists symbols you can use to more precisely locate various elements.

Symbol	Meaning
=	Matches an attribute value exactly
*=	Partially matches an attribute value
^=	Matches the beginning of an attribute value
\$=	Matches the end of an attribute value
+	Next sibling of matching element
>	Search an element's children

Table 15.3: CSS symbols for use with Selenium

You can do many other useful things with CSS selectors. A helpful guide can be found at https://www.w3schools.com/cssref/css_selectors.asp. The code below gives an example using arguments described above.

```
# Find all <a> tags with id="link1"
>>> pig_soup.select("[id='link1']")
[<a class="pig" href="http://example.com/larry" id="link1">Larry,</a>]

# Find all tags with an href attribute containing 'curly'.
>>> pig_soup.select("[href*='curly']")
[<a class="pig" href="http://example.com/curly" id="link3">Curly.</a>]

# Find all <a> tags with an href attribute
>>> pig_soup.select("a[href]")
[<a class="pig" href="http://example.com/larry" id="link1">Larry,</a>,
 <a class="pig" href="http://example.com/mo" id="link2">Mo</a>,
 <a class="pig" href="http://example.com/curly" id="link3">Curly.</a>]

# Find all <b> tags within a <p> tag with class='title'
```

```
>>> pig_soup.select("p[class='title'] b")
[<b>The Three Little Pigs</b>]

# Use a comma to find elements matching one of two arguments
>>> pig_soup.select("a[href$='mo'],[id='link3']")
[<a class="pig" href="http://example.com/mo" id="link2">Mo</a>,
 <a class="pig" href="http://example.com/curly" id="link3">Curly.</a>]
```

Problem 5. The file `large_banks_index.html` is an index of data about large banks, as recorded by the Federal Reserve.^a Write a function that reads the file and loads the source into BeautifulSoup. Return a list of the tags containing the links to bank data from September 30, 2003 to December 31, 2014, where the dates are in reverse chronological order.

^aSee <https://www.federalreserve.gov/releases/lbr/>.

Problem 6. The file `large_banks_data.html` is one of the pages from the index in Problem 5.^a Write a function that reads the file and loads the source into BeautifulSoup. Create a single figure with two subplots:

1. A sorted bar chart of the seven banks with the most domestic branches.
2. A sorted bar chart of the seven banks with the most foreign branches.

In the case of a tie, sort the banks alphabetically by name.

^aSee <http://www.federalreserve.gov/releases/lbr/20030930/default.htm>.

16

K-Means Clustering

Lab Objective: Clustering is the one of the main tools in unsupervised learning—machine learning problems where the data comes without labels. In this lab we implement the k-means algorithm, a simple and popular clustering method, and apply it to geographic clustering and color quantization.

Jupyter Notebooks

Unlike previous labs where the python file submitted was a normal .py file, this lab among others will be done in a Jupyter Notebook (.ipynb or an iPython Notebook). Jupyter Notebooks is a powerful tool in visualizing data. If you have used Google Colab, this works in a similar manner but it is run on your personal machine.

Once Jupyter Notebook is installed, there are several ways of starting a Jupyter Notebook. The easiest way is to open a new terminal window and navigate to the directory with your .ipynb file, once in the desired directory, type Jupyter notebook. This should automatically open a web browser to the Jupyter Notebook dashboard, from there you can select the .ipynb file and open and edit it.

The Python kernel will keep running in the background until told to stop. So when you are done, to close the Jupyter Notebook, you need to go to file-> Close and Halt, or in the terminal window press ctrl+c (cmd+c for Mac).

Achtung!

Before you push this file to Bitbucket to be graded, be sure to run each cell. When you push a .ipynb file, the current state of the file is pushed. This means what you see is exactly what the graders will see.

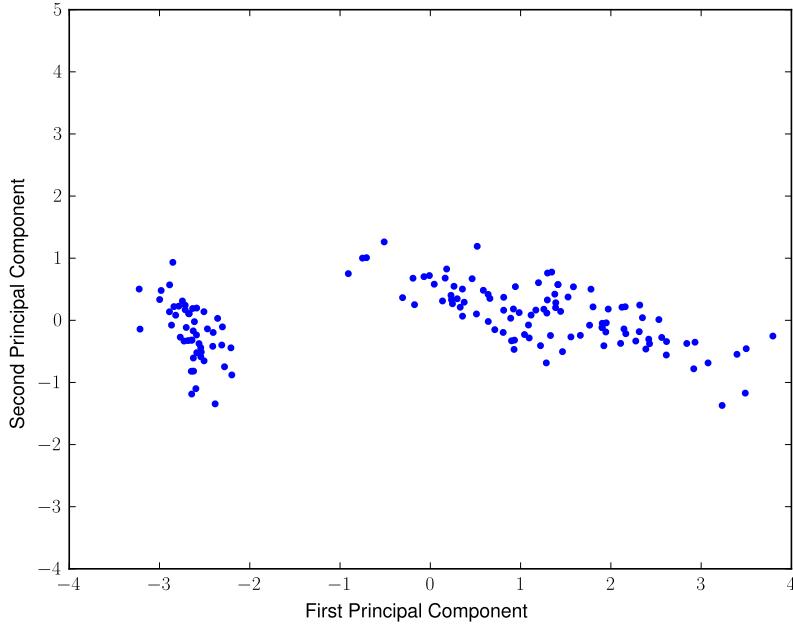


Figure 16.1: The first two principal components of the iris dataset.

Clustering

In this lab, we will analyze a few different datasets from Scikit-Learn's library and use the K-means algorithm. Figure 16.1 is a graph of the iris dataset. As a human, it is easy to identify the two distinct groups of data. Can we create an algorithm to identify these groups without human supervision? This task is called clustering, an instance of unsupervised learning. The K-means algorithm is a simple way of helping computers see the group distinctions.

The objective of clustering is to find a partitions of the data such that points in the same subset will be “close” according to some metric. The metric used will likely depend on the data, but some obvious choices include Euclidean distance and angular distance. Throughout this lab, we will use the metric $d(x, y) = \|x - y\|_2$, the Euclidean distance between x and y , unless we specify a different metric to be used.

More formally, suppose we have a collection of \mathbb{R}^K -valued observations $X = \{x_1, x_2, \dots, x_n\}$. Let $N \in \mathbb{N}$ and let \mathcal{S} be the set of all N -partitions of X , where an N -partition is a partition with exactly N nonempty elements. We can represent a typical partition in \mathcal{S} as $S = \{S_1, S_2, \dots, S_N\}$, where

$$X = \bigcup_{i=1}^N S_i$$

and

$$|S_i| > 0, \quad i = 1, 2, \dots, N.$$

We seek the N -partition S^* that minimizes the within-cluster sum of squares, i.e.

$$S^* = \arg \min_{S \in \mathcal{S}} \sum_{i=1}^N \sum_{x_j \in S_i} \|x_j - \mu_i\|_2^2,$$

where μ_i is the mean of the elements in S_i , i.e.

$$\mu_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j.$$

The K-Means Algorithm

Finding the global minimizing partition S^* is generally intractable since the set of partitions can be very large indeed, but the k-means algorithm is a heuristic approach that can often provide reasonably accurate results.

We begin by specifying an initial cluster mean $\mu_i^{(1)}$ for each $i = 1, \dots, N$. This can be done by random initialization, or according to some heuristic. For each iteration, we adopt the following procedure. Given a current set of cluster means $\mu^{(t)}$, we find a partition $S^{(t)}$ of the observations such that

$$S_i^{(t)} = \{x_j : \|x_j - \mu_i^{(t)}\|_2^2 \leq \|x_j - \mu_l^{(t)}\|_2^2, l = 1, \dots, N\}.$$

We then update our cluster means by computing for each $i = 1, \dots, N$. We continue to iterate in this manner until the partition ceases to change.

Figure 16.2 shows two different clusterings of the iris data produced by the k-means algorithm. Note that the quality of the clustering can depend heavily on the initial cluster means. We can use the within-cluster sum of squares as a measure of the quality of a clustering (a lower sum of squares is better). Where possible, it is advisable to run the clustering algorithm several times, each with a different initialization of the means, and keep the best clustering. Note also that it is possible to have very slow convergence. Thus, when implementing the algorithm, it is a good idea to terminate after some specified maximum number of iterations.

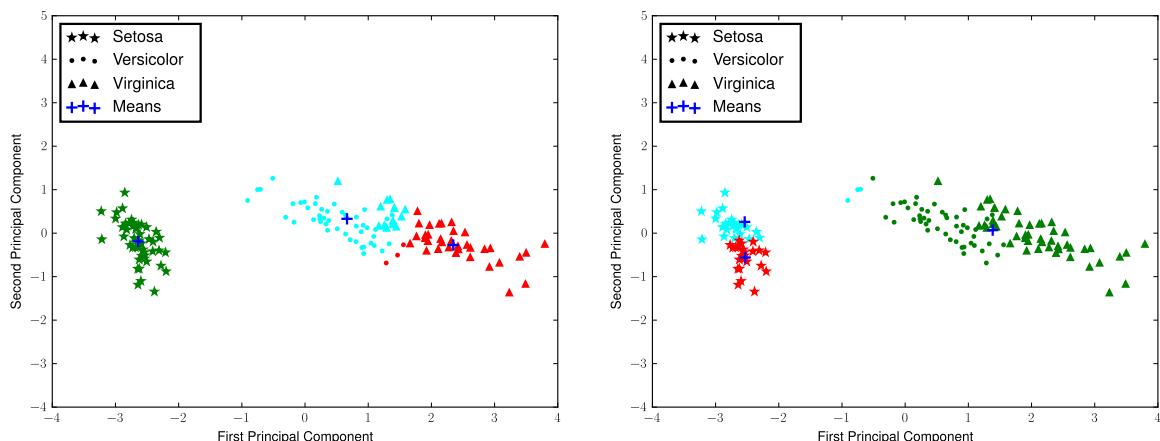


Figure 16.2: Two different K-Means clusterings for the iris dataset. Notice that the clustering on the left predicts the flower species to a high degree of accuracy, while the clustering on the right is less effective.

The algorithm can be summarized as follows.

1. From the data points, choose k initial cluster centers.
2. For $i = 0, \dots, \text{max_iter}$,

- (a) Assign each data point to the cluster center that is closest, forming k clusters.
- (b) Recompute the cluster centers as the means of the new clusters.
- (c) If the old cluster centers and the new cluster centers are sufficiently close, terminate early.

Problem 1. Write a `KMeans` class for doing basic k -means clustering. Implement the following methods, following `sklearn` class conventions.

1. `__init__()`: Accept a number of clusters k , a maximum number of iterations, and a convergence tolerance. Store these as attributes.
2. `fit()`: Accept an $m \times n$ matrix X of m data points with n features. Choose k random rows of X as the initial cluster centers. Run the k -means iteration until consecutive centers are within the convergence tolerance, or until iterating the maximum number of times. Save the cluster centers as attributes.
If a cluster is empty, reassign the cluster center as a random row of X .
3. `predict()`: Accept an $l \times n$ matrix X of data. Return an array of l integers where the i th entry indicates which cluster center the i th row of X is closest to.

Test your class on the iris data set after reducing the data to two principal components. Plot the data, coloring by cluster.

Fire Station Placement

When urban planners are making plans for a city, there are many city elements to consider. One of which is the locations of the fire stations that will service the city. When choosing a suitable location for the city, urban planners look at the current building locations, the roads nearby each location, prior traffic history and the areas of potential growth. We will simplify this complex problem by only taking into account the distances from each building to the nearest fire station (see Additional Material for a harder version of this problem).

Using another data set from SKLearn, we can get the data from the 1990 US Census for California housing based on the blocks of the residents. This has been saved in `sacramento.npy` and can be accessed by using the `np.load()` function. This file contains demographic data for each block in Sacramento and nearby cities. The eight columns in the file are: median block income, median house age in the block, average number of rooms, average number of bedrooms, average house occupancy, latitude and longitude.

There are couple ways for a fire station to be optimally placed. The stations could be placed to minimize the average distance to each house. Another option is to minimize the distance to the farthest house in each group. For this problem, minimize the distance to the farthest house in each group.

Problem 2. Using the Methods you wrote in Problem 1, add a parameter, p , to your class that denotes the norm and defaults to 2. Save p as an attribute to be used in your `fit()` and `predict()` functions. Using the data in `sacramento.npy` find the optimal placement for the fire stations. Plot the longitude and latitudes, the centers, and color them by cluster. Try different values for p to find the optimal locations for the fire stations. As the initial centers are chosen at random, make sure to run the `predict()` function several times. In a Markdown cell report which norm was the best at keeping the maximum distance small.

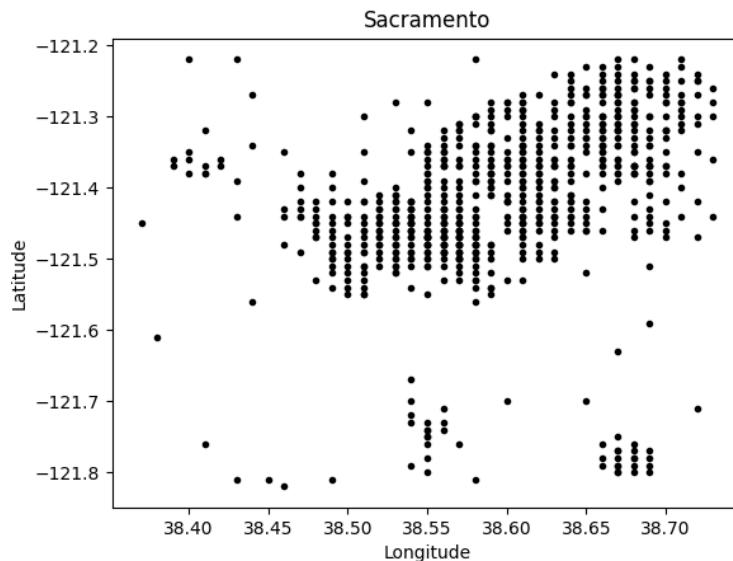


Figure 16.3: Sacramento Housing Data (1990 US Census).

Detecting Active Earthquake Regions

Suppose we are interested in learning about which regions are prone to experience frequent earthquake activity. We could make a map of all earthquakes over a given period of time and examine it ourselves, but this, as an unsupervised learning problem, can be solved using our k -means clustering tool.

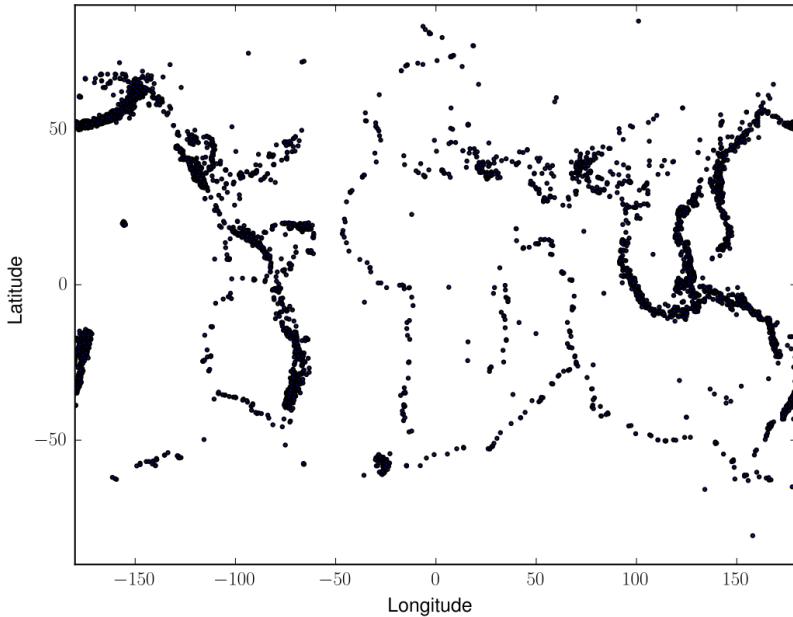


Figure 16.4: Earthquake epicenters over a 6 month period.

The file `earthquake_coordinates.npy` contains earthquake data throughout the world from January 2010 through June 2010. Each row represents a different earthquake; the columns are scaled longitude and latitude measurements. We want to cluster this data into active earthquake regions. For this task, we might think that we can regard any epicenter as a point in \mathbb{R}^2 with coordinates being their latitude and longitude. This, however, would be incorrect, because the earth is not flat. Instead, latitude and longitude should be viewed in spherical coordinates in \mathbb{R}^3 , which could then be clustered.

A simple way to accomplish this transformation is to first transform the latitude and longitude values to spherical coordinates, and then to Euclidean coordinates. Recall that a spherical coordinate in \mathbb{R}^3 is a triple (r, θ, φ) , where r is the distance from the origin, θ is the radial angle in the xy -plane from the x -axis, and φ is the angle from the z -axis. In our earthquake data, once the longitude is converted to radians it is an appropriate θ value; the latitude needs to be offset by 90° degrees, then converted to radians to obtain φ . For simplicity, we can take $r = 1$, since the earth is roughly a sphere. We can then transform to Euclidean coordinates using the following relationships.

$$\theta = \frac{\pi}{180} (\text{longitude}) \quad \varphi = \frac{\pi}{180} (90 - \text{latitude})$$

$$\begin{aligned} r &= \sqrt{x^2 + y^2 + z^2} & x &= r \sin \varphi \cos \theta \\ \varphi &= \arccos \frac{z}{r} & y &= r \sin \varphi \sin \theta \\ \theta &= \arctan \frac{y}{x} & z &= r \cos \varphi \end{aligned}$$

There is one last issue to solve before clustering. Each earthquake data point has norm 1 in Euclidean coordinates, since it lies on the surface of a sphere of radius 1. Therefore, the cluster centers should also have norm 1. Otherwise, the means can't be interpreted as locations on the surface of the earth, and the k-means algorithm will struggle to find good clusters. A solution to this problem is to normalize the mean vectors at each iteration, so that they are always unit vectors.

Problem 3. Add a keyword argument `normalize=False` to your `KMeans` constructor. Modify `fit()` so that if `normalize` is `True`, the cluster centers are normalized at each iteration.

Cluster the earthquake data in three dimensions by converting the data from raw data to spherical coordinates to euclidean coordinates on the sphere.

1. Convert longitude and latitude to radians, then to spherical coordinates.
(Hint: `np.deg2rad()` may be helpful.)
2. Convert the spherical coordinates to euclidean coordinates in \mathbb{R}^3 .
3. Use your `KMeans` class with normalization to cluster the euclidean coordinates.
4. Translate the cluster center coordinates back to spherical coordinates, then to degrees.
Transform the cluster means back to latitude and longitude coordinates.
(Hint: use `numpy.arctan2()` for arctan, so that the correct quadrant is chosen).
5. Plot the data, coloring by cluster. Also mark the cluster centers.

With 15 clusters, your plot should resemble the Figure 16.5.

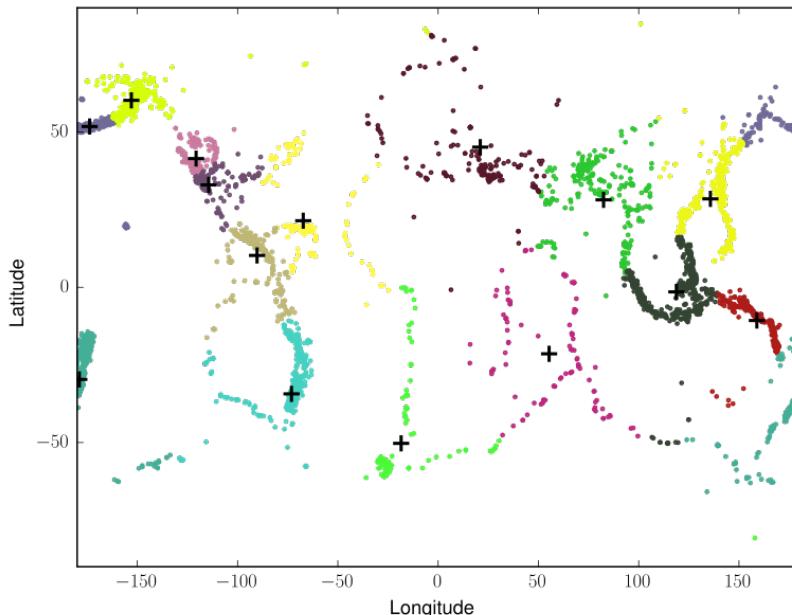


Figure 16.5: Earthquake epicenter clusters with $k = 15$.

Color Quantization

The k -means algorithm uses the euclidean metric, so it is natural to cluster geographic data. However, clustering can be done in any abstract vector space. The following application is one example.

Images are usually represented on computers as 3-dimensional arrays. Each 2-dimensional layer represents the red, green, and blue color values, so each pixel on the image is really a vector in \mathbb{R}^3 . Clustering the pixels in RGB space leads a one kind of image segmentation that facilitate memory reduction.

Reading: https://en.wikipedia.org/wiki/Color_quantization

Problem 4. Write a function that accepts an image array (of shape $(m, n, 3)$), an integer number of clusters k , and an integer number of samples S . Reshape the image so that each row represents a single pixel. Choose S pixels to train a k -means model on with k clusters. Make a copy of the original picture where each pixel has the same color as its cluster center. Return the new image. For this problem, you may use `sklearn.cluster.KMeans` instead of your `KMeans` class from Problem 1.

Test your function on some of the provided NASA images.

Additional Material

Spectral Clustering

We now turn to another method for solving a clustering problem, namely that of Spectral Clustering. As you can see in Figure ???, it can cluster data not just by its location on a graph, but can even separate shapes that overlap others into distinct clusters. It does so by utilizing the spectral properties of a Laplacian matrix. Different types of Laplacian matrices can be used. In order to construct a Laplacian matrix, we first need to create a graph of vertices and edges from our data points. This graph can be represented as a symmetric matrix W where w_{ij} represents the edge from x_i to x_j . In the simplest approach, we can set $w_{ij} = 1$ if there exists an edge and $w_{ij} = 0$ otherwise. However, we are interested in the similarity of points, so we will weight the edges by using a similarity measure. Points that are similar to one another are assigned a high similarity measure value, and dissimilar points a low value. One possible measure is the Gaussian similarity function, which defines the similarity between distinct points x_i and x_j as

$$s(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

for some set value σ .

Note that some similarity functions can yield extremely small values for dissimilar points. We have several options for dealing with this possibility. One is simply to set all values which are less than some ε to be zero, entirely erasing the edge between these two points. Another option is to keep only the T largest-valued edges for each vertex. Whichever method we choose to use, we will end up with a weighted similarity matrix W . Using this we can find the diagonal degree matrix D , which gives the number of edges found at each vertex. If we have the original fully-connected graph, then $D_{ii} = n - 1$ for each i . If we keep the T highest-valued edges, $D_{ii} = T$ for each i .

As mentioned before, we may use different types of Laplacian matrices. Three such possibilities are:

1. The unnormalized Laplacian, $L = D - W$
2. The symmetric normalized Laplacian, $L_{sym} = I - D^{-1/2}WD^{-1/2}$
3. The random walk normalized Laplacian, $L_{rw} = I - D^{-1}W$.

Given a similarity measure, which type of Laplacian to use, and the desired number of clusters k , we can now proceed with the Spectral Clustering algorithm as follows:

- Compute W , D , and the appropriate Laplacian matrix.
- Compute the first k eigenvectors u_1, \dots, u_k of the Laplacian matrix.
- Set $U = [u_1, \dots, u_k]$, and if using L_{sym} or L_{rw} normalize U so that each row is a unit vector in the Euclidean norm.
- Perform k -means clustering on the n rows of U .
- The n labels returned from your `kmeans` function correspond to the label assignments for x_1, \dots, x_n .

As before, we need to run through our k -means function multiple times to find the best measure when we use random initialization. Also, if you normalize the rows of U , then you will need to set the argument `normalize = True`.

Problem 5. Implement the Spectral Clustering Algorithm by calling your `kmeans` function, using the following function declaration:

```
def specClus(measure,Laplacian,args,arg1=None,kiters=10):
    """
    Cluster a dataset using the k-means algorithm.

    Parameters
    -----
    measure : function
        The function used to calculate the similarity measure.
    Laplacian : int in {1,2,3}
        Which Laplacian matrix to use. 1 corresponds to the unnormalized,
        2 to the symmetric normalized, 3 to the random walk normalized.
    args : tuple
        The arguments as they were passed into your k-means function,
        consisting of (data, n_clusters, init, max_iter, normalize). Note
        that you will not pass 'data' into your k-means function.
    arg1 : None, float, or int
        If Laplacian==1, it should remain as None
        If Laplacian==2, the cut-off value, epsilon.
        If Laplacian==3, the number of edges to retain, T.
    kiters : int
        How many times to call your kmeans function to get the best
        measure.

    Returns
    -----
    labels : ndarray of shape (n,)
        The i-th entry is an integer in [0,n_clusters-1] indicating
        which cluster the i-th row of data belongs to.
    """
    pass
```

We now need a way to test our code. The website <http://cs.joensuu.fi/sipu/datasets/> contains many free data sets that will be of use to us. Scroll down to the "Shape sets" heading, and download some of the datasets found there to use for trial datasets.

Problem 6. Create a function that will return the accuracy of your spectral clustering implementation, as follows:

```
def test_specClus(location,measure,Laplacian,args,arg1=None,kiters=10):
    """
    Cluster a dataset using the k-means algorithm.
```

```

Parameters
-----
location : string
    The location of the dataset to be tested.
measure : function
    The function used to calculate the similarity measure.
Laplacian : int in {1,2,3}
    Which Laplacian matrix to use. 1 corresponds to the unnormalized,
    2 to the symmetric normalized, 3 to the random walk normalized.
args : tuple
    The arguments as they were passed into your k-means function,
    consisting of (data, n_clusters, init, max_iter, normalize). Note
    that you will not pass 'data' into your k-means function.
arg1 : None, float, or int
    If Laplacian==1, it should remain as None
    If Laplacian==2, the cut-off value, epsilon.
    If Laplacian==3, the number of edges to retain, T.
kitters : int
    How many times to call your kmeans function to get the best
    measure.

Returns
-----
accuracy : float
    The percent of labels correctly predicted by your spectral
    clustering function with the given arguments (the number
    correctly predicted divided by the total number of points).
"""
pass

```

Fire Station Placement II

In problem 2 we looked at choosing the best location for a fire station. However, because we looked at the city of Sacramento where the geography doesn't role in choosing a location, we didn't need to double check that there is a place for the station. The `sanfrancisco.npy` data is organized the same way as `sacramento.py`, as this also comes from the SKLearn California Housing Module. Doing the same method as before will give us groups of houses, however, the group centers may be in the middle of the bay. When implementing this problem, perform a check on the centers to make sure they are not in water. The file `bayboundary.npy` gives a rough outline of where the bay is. The `bayboundary.npy` has only 2 columns, longitude and latitude. Using the boundaries set, make sure that the chosen centers are on land and not on water.

Problem 7. Import and parse the data from the `bayboundary.npy` and the `sanfrancisco.npy` files. Using either the algorithm that you wrote in problem 1 or the *k*-means algorithm in the SK Learn library, find the optimal locations for the 16 fire stations.

After the algorithm has finished running, check to see if the new coordinates are on land.
Return the graph of the clusters, the centers (the fire station locations) as different colors.

17

Data Augmentation

Lab Objective: Explore different methods of extending data sets to create more robust classifiers.

Data Augmentation

It is not hard to find amusing examples of deep neural networks or other machine learning systems that are brittle and respond poorly to inputs that are only slightly different than the data that the systems were trained on. One way to address brittleness in machine learning systems is to train them on a much wider range of examples. Adult humans, for example, have seen many images of stop signs in a wide variety of settings. If a machine learning system had seen as many different images of stop signs in as many different settings, it would be much more robust. But all this requires large amounts of data, and good labeled data is hard to come by.

One common approach for generating new data is data augmentation, that is, generating new data from old data by applying various transformations that we know should not change the label. For example, an image of a stop sign can be slightly translated, rotated, skewed, or cropped and still be a legitimate image of a stop sign. It may even be blurred or partially obscured, so a classifier for identifying a stop sign must be robust in the face of this sort of interference. Images that don't contain text can often be flipped horizontally, and depending on the image, can also sometimes be flipped vertically. We can use these techniques to generate a larger data set and create more robust classifiers.

As shown in Figure 1 below, these transformations still accurately depict a lion while providing slight modifications to the original image. Image transformations are comprised of three steps. First, create a $2 \times (d_1 * d_2)$ coordinate representation of the image (d_1 and d_2 are the dimensions of the image). For example, if the image is 3 pixels by 3 pixels, the coordinate representation would be

$$C = \begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 2 & 2 & 2 \\ 0 & 1 & 2 & 0 & 1 & 2 & 0 & 1 & 2 \end{bmatrix}$$

where each column represents the coordinate point of a pixel. Next, perform a linear transformation on the coordinate matrix. Note that some transformations will return float values; however, they need to be integers because the matrix represents coordinate points so you will need to round the values. Finally, use the transformed coordinates to create a new image. (The function `np.take()` may be useful for this.)

Visualizing the code below would give you the second image in Figure 1.

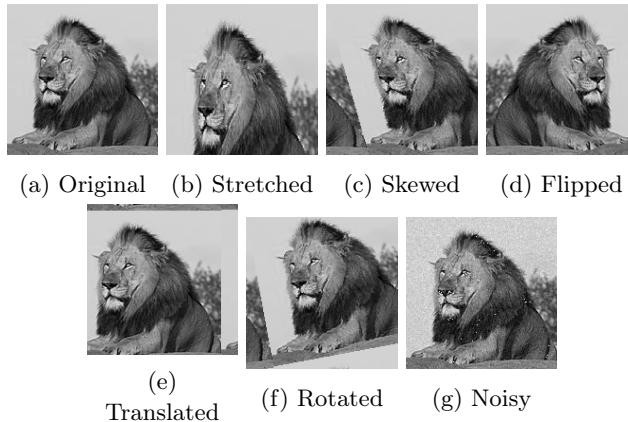


Figure 17.1: Different image transformations

```
>>> import numpy as np
>>> from imageio import imread
>>> import matplotlib.pyplot as plt

>>> lion = imread('sq_lion.png')
>>> d1, d2 = lion.shape
>>> # Get the coordinate points for each pixel in the image
>>> coords = np.mgrid[0:d1, 0:d2].reshape((2,d1*d2))
>>> # Create a linear transformation matrix. (This one will stretch the matrix)
>>> stretch_matrix = np.array([[1, 0], [0, .7]])
>>> # apply the linear transformation to the coordinate matrix
>>> new_coords = stretch_matrix@coords
>>> # some transformations will return entries as floats, but we need them to ←
      be integers because they are coordinates
>>> new_coords = new_coords.astype(int)
>>> # the next two steps apply the transformation to the image
>>> x, y = new_coords.reshape((2, d1, d2), order='F')
>>> stretched_lion = np.take(lion, x+d1*y, mode='wrap').reshape((d1, d2))
```

When augmenting a data set, it is also important to consider what types of augmentation would provide useful samples. For example, if training a dataset to recognize lowercase letters, flipping an image upside down may not be a useful transformation (consider the letters b and p). It is important to ensure that the transformed data is still a reasonable example of the object it is classified as.

Problem 1. Code from scratch the following simple black-and-white image augmenters that take as input the data X (a $d_1 \times d_2$ array that contains an image) and parameters controlling the transformation. It should return the transformed data $f(X)$. Note that the image should receive its own random treatment; for example, if the image is being translated, then the image should be translated by a different (randomly drawn) amount. Your functions should have the the following names and perform the corresponding transform:

1. `translate(X,A,B)`, with parameters A, B . Returns an image translated by a random amount (a, b) , where $a \sim \text{Uniform}(-A, A)$, and $b \sim \text{Uniform}(-B, B)$. The resulting image should be cropped to be of size $d_1 \times d_2$. Note that this translation will leave a border on two sides of the image. Fill the empty border with the parts that were cropped off the opposite sides.
 2. `rotate(X,T)`, with parameter Θ . Returns the image rotated by a random amount $\theta \sim \text{Uniform}(-\Theta, \Theta)$. The resulting image should be cropped to be the same size as the original, and any blank parts should be filled with one of the parts cropped off the other side. HINT: The rotation matrix is:
- $$\begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}.$$
3. `skew(X,A)`, with parameter A . Returns the image with the linear transformation $\begin{bmatrix} 1 & a \\ 0 & 1 \end{bmatrix}$ applied, where $a \sim \text{Uniform}(0, A)$. Crop parts that go outside the image boundaries and fill missing areas with the appropriate cropped piece.
 4. `flip_horizontal(X)`. Returns a horizontally-flipped version of the image.
 5. `gauss_noise(X,s)`, with parameter σ^2 . For the image draw $d_1 \times d_2$ random noise values from $\text{Normal}(0, \sigma^2)$ and add those to the original image.

Show that each transformation works by displaying a transformed version of lion.png.

Problem 2. Create a function called `image_augment` that will augment your data set using each of the transformations created in problem 1. This function should accept the parameters X (the images), Y (labels), and a list of the parameters for each transformation. The function should return an augmented data set with 6 times the number of images N and an array containing the appropriate label for each image.

Take the sklearn digits dataset (see the code box below for importing instructions), make an 80-20 train-test split, and then apply each of your transformations to the entire training set using `image_augment`. You must decide good values of each of the parameters to use. This should give you a larger (augmented) training set with roughly 8,600 training points. Fit a random forest to the augmented training set and to the original training set and score it using the test set. Return the average score for both classifiers (be sure to label the scores).

Your augmented score should be better than your original score.

```
from sklearn import datasets

digits = datasets.load_digits()
X, y = digits.data, digits.target
```

Audio Augmentation

For audio data, adding various forms of noise is still a reasonable augmentation choice, but many of the other augmentation methods used for images aren't really suitable. Some useful methods include dropping data at certain time steps, blocking certain frequencies, and changing the pitch or speed.

When adding noise, it may be useful to consider the types of noise most likely to be encountered when the method is in use. For example, if the method will be used to identify voice commands in an outdoor environment, then adding in typical outdoor noises would probably be more useful than adding white noise, which may not occur much in everyday life. Be thoughtful in choosing how to augment your data, as some types of manipulations may change or obscure the data to the point where it is no longer recognizable.

Audio Packages

There are many different python packages that provide different analysis and audio manipulation tools. Some common packages include `pyAudioAnalysis`, `PYO`, and `ffmpeg-python`. For the purposes of this lab, we will be using `LibROSA`. `LibROSA` is a python package that allows users to read in, write, analyze, and alter .wav files. It can be used to augment an audio dataset and extract features that can be used to distinguish between different sounds or types of music. (`LibROSA` is not included in Anaconda, so you may need to install it with `pip install librosa`.)

```
>>> import matplotlib.pyplot as plt
>>> import librosa
>>> import librosa.display #this needs to be imported separately,
                           #but is included in the librosa package
>>> import numpy as np
# load the audio time series and sampling rate
>>> chopin, sample_rate = librosa.load('chopin.wav', sr = 22050)
>>> plt.figure(figsize = (15,5))
>>> librosa.display.waveplot(chopin,sample_rate)           #generate plot
>>> plt.show()
```

The `LibROSA` library heavily relies on NumPy arrays. If you already have a NumPy array and its sample rate, you can skip the `librosa.load()`.

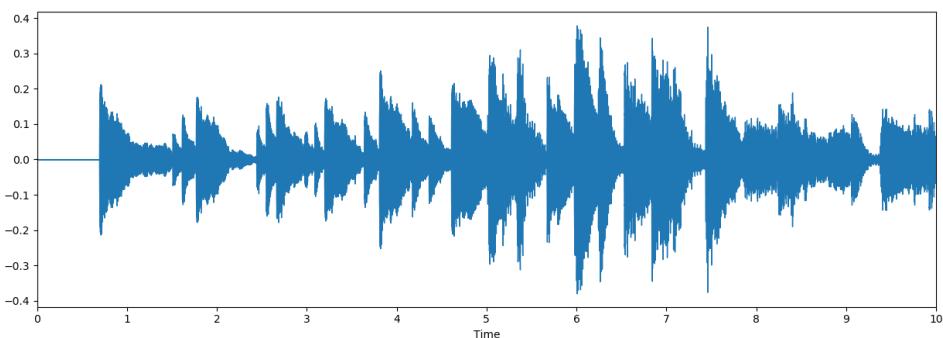


Figure 17.2: Visualization of an audio file

When you load in a .wav file, LibROSA returns the audio as an ndarray and a sample rate.

Depending on the type of audio you are analyzing, different functions may provide better distinguishing characteristics than others. It is important to take these characteristics into account when augmenting data. One possible feature that could be used for classifying music is the "predominant local pulse estimation" or PLP, as shown below. (PLP essentially takes the pulse of the music, just like you can take your own pulse in your wrist.)

```
>>> pulse = librosa.beat.plp(chopin)
>>> plt.figure(figsize = (15,5))
>>> plt.plot(np.linspace(0,10,len(pulse)),pulse)
>>> plt.show()
```

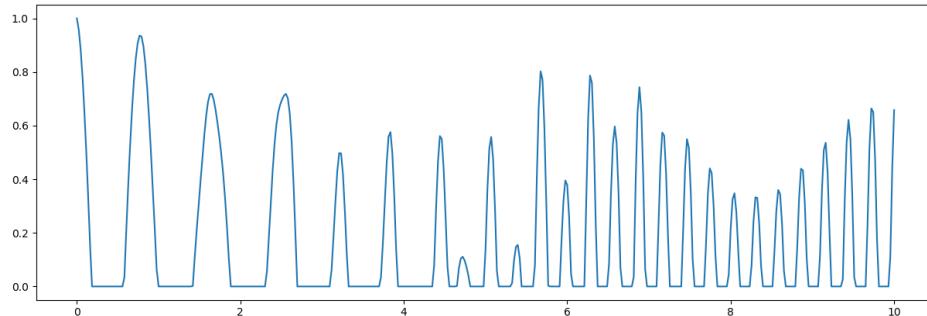


Figure 17.3: Predominant local pulse estimation of audio from figure 17.2

The LibROSA package contains several different functions that can be used to manipulate audio data, several of which are described in the table below.

Function	Returns
<code>time_stretch()</code>	slows or speeds up audio series by a fixed rate
<code>pitch_shift()</code>	shifts the pitch by n_steps semitones
<code>harmonic()</code>	extracts the harmonic elements from an audio time-series
<code>percussive()</code>	extracts percussive elements from an audio time-series
<code>split()</code>	splits an interval into non-silent intervals
<code>remix()</code>	re-orders time intervals

Table 17.1: These descriptions were taken directly from librosa.org/librosa/effects.html

Problem 3. The file `music.npy` contains the audio time series data of 10 second clips from 150 different songs, with `styles.npy` describing the associated style of ballroom dance. The styles included are Chacha, Foxtrot, Jive, Samba, Rumba, and Waltz. Use `train_test_split` from `sklearn.model_selection` with `test_size=.5` to create train and test sets.

Create two training sets by augmenting this original training set. Each new augmented training set will include the original data and the augmented data. For the first, add ambient noise from the file `restaurant-ambience.wav`. For the second, use `time_stretch`.

HINT: Since the ambient noise clip is much longer than the other music clips, you will have to select a sample of the ambient noise to add to the other clips. It may also benefit you to randomize which ambient noise sample you add to each clip, you can do this by choosing a random index to start from, and sampling starting at that index.

Problem 4. Do the following steps 5 times:

- Use the original data set and the augmented data sets to fit three RandomForestClassifiers, one only on the original data, one on the original data and the data with ambient noise added, and one on the original data and the time stretched data.
- Score each classifier.

Print the mean score for each of the classifiers and print the standard deviation for the scores.

HINT: Use the PLP as a feature you use to fit and classify. This example may be helpful for printing your results nicely.

```
print('\t\t Mean \t STD')
print('Original', '\t', np.round(orig.mean(), 3), '\t', np.round(orig.std(), 3))
print('Ambient Noise', '\t', np.round(amb.mean(), 3), '\t', np.round(amb.std()←
    , 3))
print('Time Stretch:', '\t', np.round(time.mean(), 3), '\t', np.round(time.std←
    (), 3))
```

Synthetic Minority Oversampling

Another situation where generating data can be helpful is in a classification problem where one class is rare compared to the others. For example the problem of identifying glioblastoma (a rare malignant brain tumor) is difficult in part because this cancer only occurs in 3 out of every 100,000 people. A classifier that predicts “no cancer” in every case performs very well (99.997% accurate).

If the training set has 100,000 total cases, only three of which are positive, then undersampling (taking the same number of negatives as positives) gives a dataset with only six total instances, and this is not enough to make a good classifier. Naïve oversampling (repeatedly drawing, with replacement, from the three positive cases to get the same number of positives as negatives) works better than undersampling the negatives, but does not perform very well, because the oversampled dataset just has (roughly) 33,332 repeated instances of each of the three positive instances, and the resulting classifier is likely to be overfit on those three instances.

In the special case where the features are all continuous, we can partially address this class-imbalance problem by synthetically generating new positive instances from the minority class samples (in this case the three positive cases). The synthetic minority oversampling technique (SMOTE)¹ works by randomly choosing points along the line segments connecting this point to each (or some) of its k nearest minority neighbors.

SMOTE tends to work better on low-dimensional data than on high-dimensional data. For example if the minority class training examples are images (one dimension per pixel, so, high dimensional) of the subject (say possible tumor cells) that are not centered and not uniformized to be of similar size, then many points along the line connecting two of these images could look nothing like the two endpoints. In such cases SMOTE is not usually very helpful.

The Algorithm

The purpose of this section is to provide a high-level understanding of how Synthetic Minority Over-sampling works and will heavily reference the SMOTE paper mentioned earlier.

The goal in creating synthetic observations is to increase the accuracy of the classifier. This means that the synthetic data generated needs to be similar to the minority class data, while creating moderate variation. To do this, select a member of the minority class and find its k nearest neighbors. Randomly select one. Now, for each feature select a random point on the line between the points.

We will demonstrate this using data with two features, represented by x and y coordinates. Consider the points $(0, 0)$, $(1, 3)$, $(2, 1)$, and $(3, 2)$, as shown in figure 17.4.

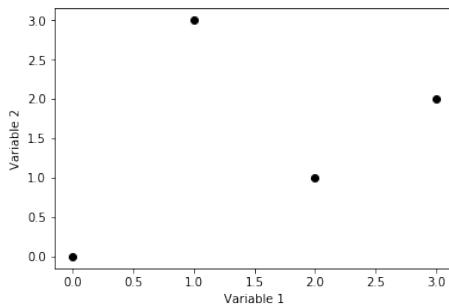


Figure 17.4: Data before SMOTE

To keep things simple, we will use $k = 1$. The nearest neighbor for $(0, 0)$ is $(2, 1)$. Choose a random point between the x values (shown in red) and a random point between the y values (shown in blue). For data with n features, a random point between the two feature values would be chosen for each feature. The intersection of these lines gives us the coordinate for the synthetic point (purple).

¹See <https://jair.org/index.php/jair/article/view/10302>

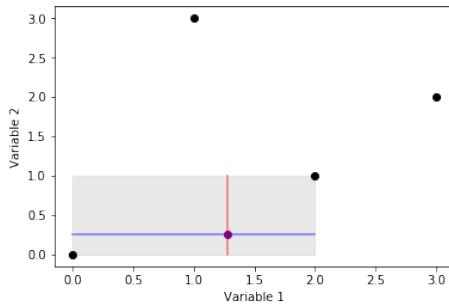


Figure 17.5: SMOTE process

Running this algorithm 500 times per original point, with $k = 1$, returns a graph like figure 17.6a. Running the algorithm 500 times per original point and increasing k to 2, returns the graph like figure 17.6b.

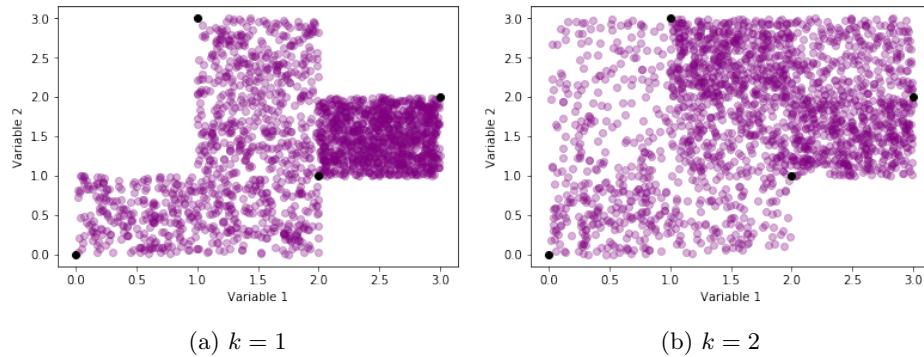


Figure 17.6: After SMOTE

Problem 5. Write a function that uses the synthetic minority oversampling technique to augment an imbalanced data set. Your function should:

- Take arguments: X a matrix of minority class samples, N the number of samples to generate per original point, and k the number of nearest neighbors.
- For each original point in the sample, randomly pick one of the k nearest neighbors and randomly generate a new point that lies between the two original values. You may use `sklearn.neighbors.KDTree` to find the k nearest neighbors.
- Return an array containing the synthetic samples.

Problem 6. The dataset found in `creditcard.npy` contains information about credit card purchases made over a two day period. Of the approximately 285,000 observations, 492 are fraudulent purchases. The last column indicates if the purchase was valid (0) or fraudulent (1).

Do the following steps 10 times:

- Create a training and test set from the data using `train_test_split` from `sklearn.model_selection` with `test_size=.7`.
- Use `smote` with $N = 500$ and $k = 2$ to augment the training set.
- Create two Gaussian Naïve Bayes classifiers (from `sklearn.naive_bayes.GaussianNB`), one which wil be trained on only the original data and the other on the SMOTE augmented data and the original data.
- Fit each classifier and find the recall and accuracy of each model.

Print the mean recall and mean accuracy of each model and and describe the findings.

HINT: Recall = $\frac{tp}{tp+fn}$. This example may be helpful for printing your results nicely.

```
>>> print('\t\t Recall \t Accuracy')
>>> print('Original', '\t', np.round(mean_orig_recall,5), '\t', np.round(←
    mean_orig_score,5))
>>> print('SMOTE', '\t\t', np.round(mean_smote_recall,5), '\t', np.round(←
    mean_smote_score,5))
```


18

Metropolis Algorithm

Lab Objective: Understand the basic principles of the Metropolis algorithm and apply these ideas to the Ising Model.

The Metropolis Algorithm

Sampling from a given probability distribution is an important task in many different applications found throughout the sciences. When these distributions are complicated, as is often the case when modeling real-world problems, direct sampling methods can become difficult, as they might involve computing high-dimensional integrals. The Metropolis algorithm is an effective method to sample from many distributions, requiring only that we be able to evaluate the probability density function up to a constant of proportionality. In particular, the Metropolis algorithm does not require us to compute difficult high-dimensional integrals, such as those that are found in the denominator of Bayesian posterior distributions.

The Metropolis algorithm is an MCMC sampling method which generates a sequence of random variables, similar to Gibbs sampling. These random variables form a Markov Chain whose invariant distribution is equal to the distribution from which we wish to sample. Suppose that $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is the probability density function of distribution, and suppose that $f(\boldsymbol{\theta}) = c \cdot h(\boldsymbol{\theta})$ for some nonzero constant c (in practice, we assume that f is an easy function to evaluate, while h is difficult). Let $Q : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ be a symmetric proposal function (so that $Q(\cdot, \mathbf{y})$ is a probability density function for all $\mathbf{y} \in \mathbb{R}^n$, and $Q(\mathbf{x}, \mathbf{y}) = Q(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$) and let $A : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ be an acceptance function defined by

$$A(\mathbf{x}, \mathbf{y}) = \min\left(1, \frac{f(\mathbf{x})}{f(\mathbf{y})}\right).$$

We can combine these functions in such a way so as to sample from the aforementioned Markov Chain by following Algorithm 18.1. The Metropolis algorithm can be interpreted as follows: given our current state \mathbf{y} , we propose a new state according to the distribution $Q(\cdot, \mathbf{y})$. We then accept or reject it according to A . We continue by repeating the process. So long as Q defines an irreducible, aperiodic, and non-null recurrent Markov chain, we will have a Markov chain whose unique invariant distribution will have density h . Furthermore, given any initial state, the chain will converge to this invariant distribution. Note that for numerical reasons, it is often wise to make calculations of the acceptance functions in log space:

$$\log A(\mathbf{x}, \mathbf{y}) = \min(0, \log f(\mathbf{x}) - \log f(\mathbf{y})).$$

Algorithm 18.1 Metropolis Algorithm

```

1: procedure Metropolis Algorithm
2:   Choose initial point  $\mathbf{y}_0$ .
3:   for  $t = 1, 2, \dots$  do
4:     Draw  $\mathbf{x} \sim Q(\cdot, \mathbf{y}_{t-1})$ 
5:     Draw  $a \sim \text{unif}(0, 1)$ 
6:     if  $a \leq A(\mathbf{x}, \mathbf{y}_{t-1})$  then
7:        $\mathbf{y}_t = \mathbf{x}$ 
8:     else
9:        $\mathbf{y}_t = \mathbf{y}_{t-1}$ 
10:    Return  $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3, \dots$ 
```

Let's apply the Metropolis algorithm to a simple example of Bayesian analysis. Consider the problem of computing the posterior distribution over the mean μ and variance σ^2 of a normal distribution for which we have n data points y_1, \dots, y_n . For concreteness, we use the data in `examscores.csv` and we assume the prior distributions

$$\begin{aligned}\mu &\sim \mathcal{N}(m = 80, s^2 = 16) \\ \sigma^2 &\sim IG(\alpha = 3, \beta = 50).\end{aligned}$$

In this situation, we wish to sample from the posterior distribution

$$p(\mu, \sigma^2 | y_1, \dots, y_N) = \frac{p(\mu)p(\sigma^2) \prod_{i=1}^n \mathcal{N}(y_i | \mu, \sigma^2)}{\int_{-\infty}^{\infty} \int_0^{\infty} p(\mu)p(\sigma^2) \prod_{i=1}^n \mathcal{N}(y_i | \mu, \sigma^2) d\sigma^2 d\mu}.$$

However, we can conveniently calculate only the numerator of this expression. Since the denominator is simply a constant with respect to μ and σ^2 , the numerator can serve as the function f in the Metropolis algorithm, and the denominator can serve as the constant c .

We choose our proposal function to be based on a bivariate Normal distribution:

$$Q(x, y) = \mathcal{N}(x | y, sI),$$

where I is the 2×2 identity matrix and s is some positive scalar.

```

>>> def proposal(y, s):
...     """The proposal function Q(x,y) = N(x|y,sI)."""
...     return stats.multivariate_normal.rvs(mean=y, cov=s*np.eye(len(y)))
...
>>> def propLogDensity(x):
...     """Calculate the log of the proportional density."""
...     logprob = muprior.logpdf(x[0]) + sig2prior.logpdf(x[1])
...     logprob += stats.norm.logpdf(scores, loc=x[0], scale=sqrt(x[1])).sum()
...     return logprob    # ^this is where the scores are used.
...
>>> def acceptance(x, y):
...     return min(0, propLogDensity(x) - propLogDensity(y))
```

We are now ready to code up the Metropolis algorithm using these functions. We will keep track of the samples generated by the algorithm, along with the proportional log densities of the samples and the proportion of proposed samples that were accepted.

We can evaluate the quality of our results by plotting the log probabilities, the μ samples, the σ^2 samples, and kernel density estimators for the marginal posterior distributions of μ and σ^2 . The kernel density estimator is the posterior distribution for a parameter. It measures the frequency of each draw. In this example, the kernel density estimator for μ should be approximately normal, and the kernel density estimator for σ^2 should be approximately an inverse gamma.

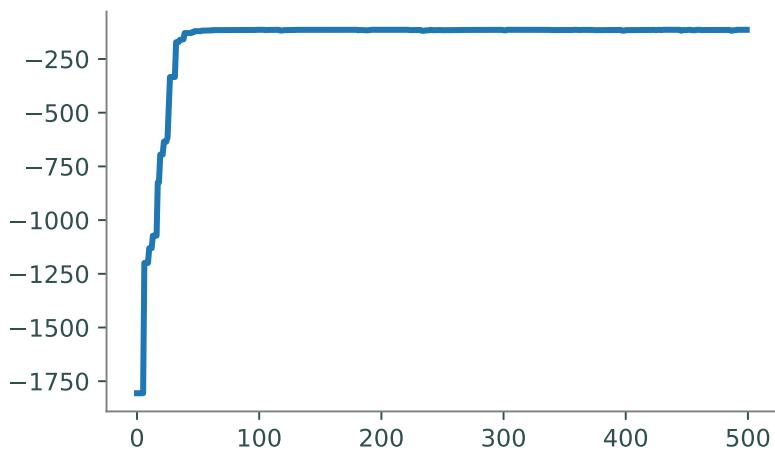


Figure 18.1: Log densities of the first 500 Metropolis samples.

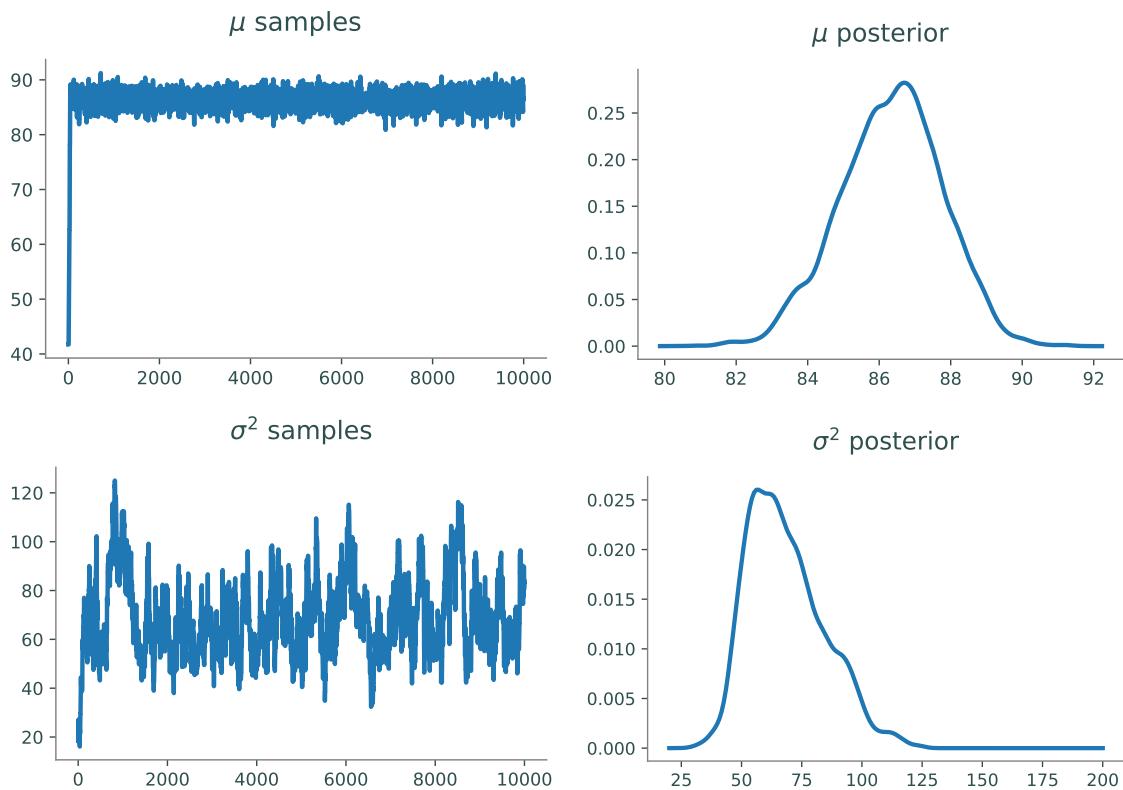


Figure 18.2: Metropolis samples and KDEs for the marginal posterior distribution of μ (top row) and σ^2 (bottom row).

Problem 1. Write a function that uses the Metropolis Hastings algorithm to draw from the posterior distribution over the mean μ and variance σ^2 . Use the given functions and algorithm 18.1 to complete the problem.

Your function should return an array of draws, an array of the log probabilities, and an acceptance rate. Use the following code to check your work. Using the seaborn.kdeplot function, plot the first 500 log probabilities, the μ samples and posterior distribution, and the σ^2 samples and posterior distribution. The results should be similar to Figures 18.1 and 18.2.

When comparing a to the acceptance, remember to use $\log(a)$ as we are in log space.

```
# Load in the data and initialize hyperparameters.
>>> scores = np.load("examscores.npy")

# Prior sigma^2 ~ IG(alpha, beta)
>>> alpha = 3
>>> beta = 50

#Prior mu ~ N(m, s)
>>> m = 80
```

```
>>> s = 4

# Initialize the prior distributions.
>>> muprior = stats.norm(loc=m, scale=sqrt(s**2))
>>> sig2prior = stats.invgamma(alpha, scale=beta)
```

The Ising Model

In statistical mechanics, the Ising model describes how atoms interact in ferromagnetic material. Assume we have some lattice Λ of sites. We say $i \sim j$ if i and j are adjacent sites. Each site i in our lattice is assigned an associated spin $\sigma_i \in \{\pm 1\}$. A state in our Ising model is a particular spin configuration $\sigma = (\sigma_k)_{k \in \Lambda}$. If $L = |\Lambda|$, then there are 2^L possible states in our model. If L is large, the state space becomes huge, which is why MCMC sampling methods (in particular the Metropolis algorithm) are so useful in calculating model estimations.

With any spin configuration σ , there is an associated energy

$$H(\sigma) = -J \sum_{i \sim j} \sigma_i \sigma_j$$

where $J > 0$ for ferromagnetic materials, and $J < 0$ for antiferromagnetic materials. Throughout this lab, we will assume $J = 1$, leaving the energy equation to be $H(\sigma) = -\sum_{i \sim j} \sigma_i \sigma_j$ where the interaction from each pair is added only once.

We will consider a lattice that is a 100×100 square grid. The adjacent sites for a given site are those directly above, below, to the left, and to the right of the site, so to speak. For sites on the edge of the grid, we assume it wraps around. In other words, a site at the farthest left side of the grid is adjacent to the corresponding site on the farthest right side. Thus, a single spin configuration can be represented as a 100×100 array, with entries of ± 1 .

The following code will construct a random spin configuration of size n :

```
def random_lattice(n):
    """Constructs a random spin configuration for an nxn lattice."""
    random_spin = np.zeros((n,n))
    for k in range(n):
        random_spin[k,:] = 2*np.random.binomial(1,.5, n) - 1
    return random_spin
```

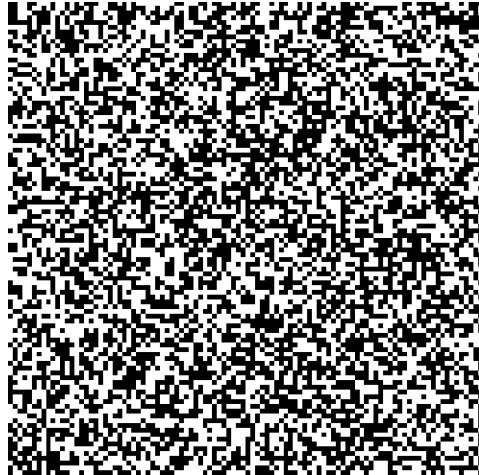


Figure 18.3: Spin configuration from random initialization.

Problem 2. Write a function that accepts a spin configuration σ for a lattice as a NumPy array. Compute the energy $H(\sigma)$ of the spin configuration. Be careful to not double count site pair interactions!

(Hint: `np.roll()` may be helpful.)

Different spin configurations occur with different probabilities, depending on the energy of the spin configuration and $\beta > 0$, a quantity inversely proportional to the temperature. More specifically, for a given β , we have

$$\mathbb{P}_\beta(\sigma) = \frac{e^{-\beta H(\sigma)}}{Z_\beta}$$

where $Z_\beta = \sum_\sigma e^{-\beta H(\sigma)}$. Because there are $2^{100 \cdot 100} = 2^{10000}$ possible spin configurations for our particular lattice, computing this sum is infeasible. However, the numerator is quite simple, provided we can efficiently compute the energy $H(\sigma)$ of a spin configuration. Thus the ratio of the probability densities of two spin configurations is simple:

$$\frac{\mathbb{P}_\beta(\sigma^*)}{\mathbb{P}_\beta(\sigma)} = \frac{e^{-\beta H(\sigma^*)}}{e^{-\beta H(\sigma)}} = e^{\beta(H(\sigma)-H(\sigma^*))}$$

The simplicity of this ratio should lead us to think that a Metropolis algorithm might be an appropriate way by which to sample from the spin configuration probability distribution, in which case the acceptance probability would be

$$A(\sigma^*, \sigma) = \begin{cases} 1 & \text{if } H(\sigma^*) < H(\sigma) \\ e^{\beta(H(\sigma)-H(\sigma^*))} & \text{otherwise.} \end{cases} \quad (18.1)$$

By choosing our transition matrix Q cleverly, we can also make it easy to compute the energy for any proposed spin configuration. We restrict our possible proposals to only those spin configurations in which we have flipped the spin at exactly one lattice site, i.e. we choose a lattice site i and flip its spin. Thus, there are only L possible proposal spin configurations σ^* given σ , each being proposed with probability $\frac{1}{L}$, and such that $\sigma_j^* = \sigma_j$ for all $j \neq i$, and $\sigma_i^* = -\sigma_i$. Note that we would never actually write out this matrix (it would be $2^{10000} \times 2^{10000}$). Computing the proposed site's energy is simple: if the spin flip site is i , then we have

$$H(\sigma^*) = H(\sigma) + 2 \sum_{j:j \sim i} \sigma_i \sigma_j. \quad (18.2)$$

Problem 3. Write a function that accepts an integer n and chooses a pair of indices (i, j) where $0 \leq i, j \leq n - 1$. Each possible pair should have an equal probability $\frac{1}{n^2}$ of being chosen.

Problem 4. Write a function that accepts a spin configuration σ , its energy $H(\sigma)$, and integer indices i and j . Use (18.2) to compute the energy of the new spin configuration σ^* , which is σ but with the spin flipped at the (i, j) th entry of the corresponding lattice. Do not explicitly construct the new lattice for σ^* .

Problem 5. Write a function that accepts a float β and spin configuration energies $H(\sigma)$ and $H(\sigma^*)$. Using (18.1), calculate whether or not the new spin configuration σ^* should be accepted (return `True` or `False`). Consider doing the calculations in log space. (Hint: `np.random.binomial()` might be useful)

To track the convergence of the Markov chain, we would like to look at the probabilities of each sample at each time. However, this would require us to compute the denominator Z_β , which is generally the reason we have to use a Metropolis algorithm to begin with. We can get away with examining only $-\beta H(\sigma)$. We should see this value increase as the algorithm proceeds, and it should converge once we are sampling from the correct distribution. Note that we don't expect these values to converge to a specific value, but rather to a restricted range of values.

Problem 6. Write a function that accepts a float $\beta > 0$ and integers n , `n_samples`, and `burn_in`. Initialize an $n \times n$ lattice for a spin configuration σ using Problem 2. Use the Metropolis algorithm to (potentially) update the lattice `burn_in` times.

1. Use Problem 3 to choose a site for possibly flipping the spin, thus defining a potential new configuration σ^* .
2. Use Problem 4 to calculate the energy $H(\sigma^*)$ of the proposed configuration.
3. Use Problem 5 to accept or reject the proposed configuration. If it is accepted, set $\sigma = \sigma^*$ by flipping the spin at the indicated site.

4. Track $-\beta H(\sigma)$ at each iteration (independent of acceptance).

After the burn-in period, continue the iteration `n_samples` times, also recording every 100th sample (to prevent memory failure). The acceptance rate is counted after the burn-in period. Return the samples, the sequence of weighted energies $-\beta H(\sigma)$, and the acceptance rate.

Test your sampler on a 100×100 grid with 200000 total iterations, with `n_samples` large enough so that you will keep 50 samples, for $\beta = 0.2, 0.4, 1$. Plot the proportional log probabilities, as well as a late sample from each test. How does the ferromagnetic material behave differently with differing temperatures? Recall that β is an inverse function of temperature. You should see more structure with lower temperature, as illustrated in Figure 18.4.

To show the spin configuration, use `plt.imshow(L, cmap='gray')`.

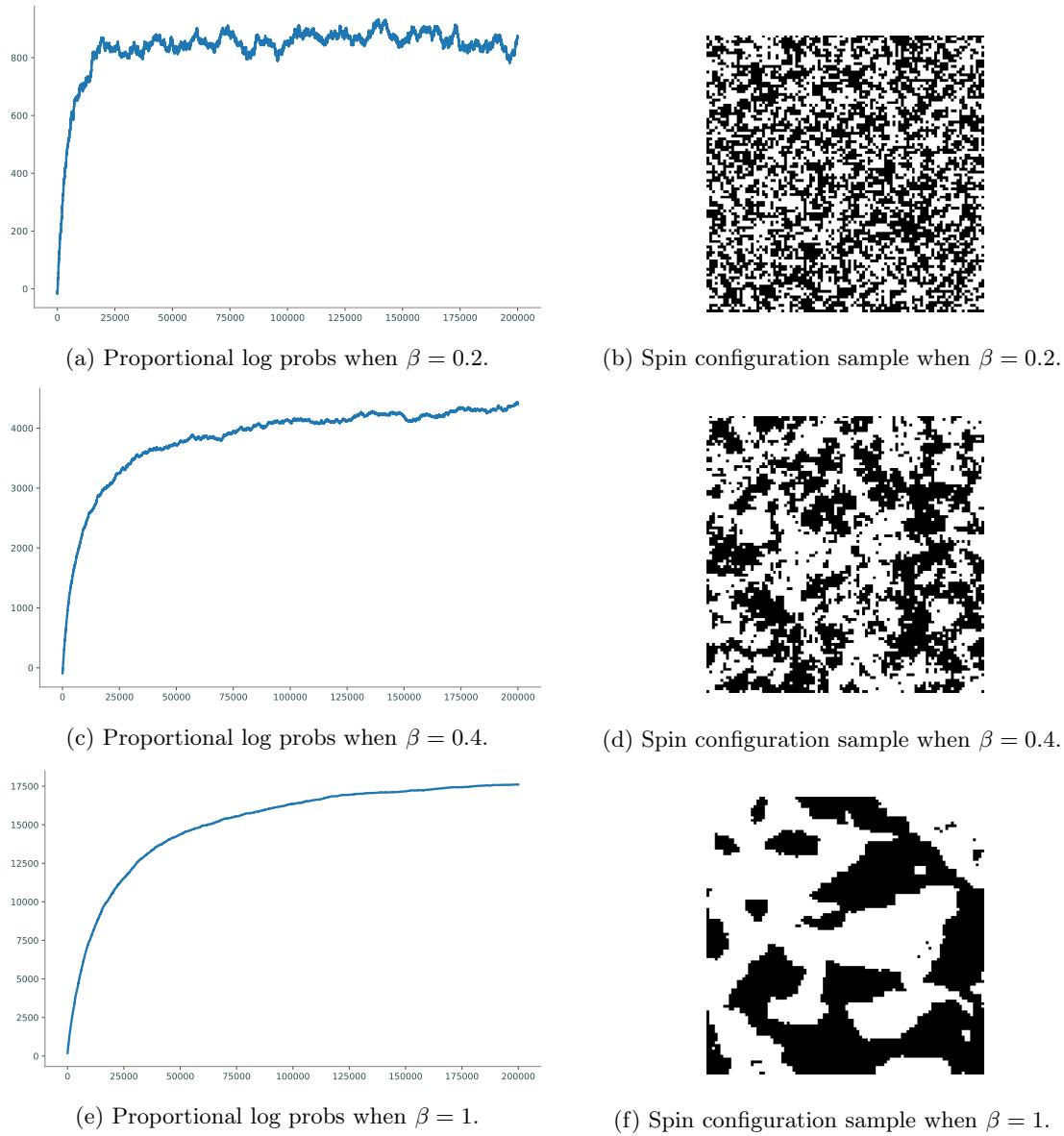


Figure 18.4

19

Gibbs Sampling and LDA

Lab Objective: Understand the basic principles of implementing a Gibbs sampler. Apply this to Latent Dirichlet Allocation.

Gibbs Sampling

Gibbs sampling is an MCMC sampling method in which we construct a Markov chain which is used to sample from a desired joint (conditional) distribution

$$\mathbb{P}(x_1, \dots, x_n | \mathbf{y}).$$

Often it is difficult to sample from this high-dimensional joint distribution, while it may be easy to sample from the one-dimensional conditional distributions

$$\mathbb{P}(x_i | \mathbf{x}_{-i}, \mathbf{y})$$

where $\mathbf{x}_{-i} = x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n$.

Algorithm 19.1 Basic Gibbs Sampling Process.

```

1: procedure Gibbs Sampler
2:   Randomly initialize  $x_1, x_2, \dots, x_n$ .
3:   for  $k = 1, 2, 3, \dots$  do
4:     for  $i = 1, 2, \dots, n$  do
5:       Draw  $x \sim \mathbb{P}(x_i | \mathbf{x}_{-i}, \mathbf{y})$ 
6:       Fix  $x_i = x$ 
7:      $\mathbf{x}^{(k)} = (x_1, x_2, \dots, x_n)$ 
```

A Gibbs sampler proceeds according to Algorithm 19.1. Each iteration of the outer for loop is a sweep of the Gibbs sampler, and the value of $\mathbf{x}^{(k)}$ after a sweep is a sample. This creates an irreducible, non-null recurrent, aperiodic Markov chain over the state space consisting of all possible \mathbf{x} . The unique invariant distribution for the chain is the desired joint distribution

$$\mathbb{P}(x_1, \dots, x_n | \mathbf{y}).$$

Thus, after a burn-in period, our samples $\mathbf{x}^{(k)}$ are effectively samples from the desired distribution.

Consider the dataset of N scores from a calculus exam in the file `examscores.npy`. We believe that the spread of these exam scores can be modeled with a normal distribution of mean μ and variance σ^2 . Because we are unsure of the true value of μ and σ^2 , we take a Bayesian approach and place priors on each parameter to quantify this uncertainty:

$$\begin{aligned}\mu &\sim N(\nu, \tau^2) && \text{(a normal distribution)} \\ \sigma^2 &\sim IG(\alpha, \beta) && \text{(an inverse gamma distribution)}\end{aligned}$$

Letting $\mathbf{y} = (y_1, \dots, y_N)$ be the set of exam scores, we would like to update our beliefs of μ and σ^2 by sampling from the posterior distribution

$$\mathbb{P}(\mu, \sigma^2 | \mathbf{y}, \nu, \tau^2, \alpha, \beta).$$

Sampling directly can be difficult. However, we can easily sample from the following conditional distributions:

$$\begin{aligned}\mathbb{P}(\mu | \sigma^2, \mathbf{y}, \nu, \tau^2, \alpha, \beta) &= \mathbb{P}(\mu | \sigma^2, \mathbf{y}, \nu, \tau^2) \\ \mathbb{P}(\sigma^2 | \mu, \mathbf{y}, \nu, \tau^2, \alpha, \beta) &= \mathbb{P}(\sigma^2 | \mu, \mathbf{y}, \alpha, \beta)\end{aligned}$$

The reason for this is that these conditional distributions are conjugate to the prior distributions, and hence are part of the same distributional families as the priors. In particular, we have

$$\begin{aligned}\mathbb{P}(\mu | \sigma^2, \mathbf{y}, \nu, \tau^2) &= N(\mu^*, (\sigma^*)^2) \\ \mathbb{P}(\sigma^2 | \mu, \mathbf{y}, \alpha, \beta) &= IG(\alpha^*, \beta^*),\end{aligned}$$

where

$$\begin{aligned}(\sigma^*)^2 &= \left(\frac{1}{\tau^2} + \frac{N}{\sigma^2} \right)^{-1} \\ \mu^* &= (\sigma^*)^2 \left(\frac{\nu}{\tau^2} + \frac{1}{\sigma^2} \sum_{i=1}^N y_i \right) \\ \alpha^* &= \alpha + \frac{N}{2} \\ \beta^* &= \beta + \frac{1}{2} \sum_{i=1}^N (y_i - \mu)^2\end{aligned}$$

We have thus set this up as a Gibbs sampling problem, where we have only to alternate between sampling μ and sampling σ^2 . We can sample from a normal distribution and an inverse gamma distribution as follows:

```
>>> from math import sqrt
>>> from scipy.stats import norm
>>> from scipy.stats import invgamma
>>> mu = 0. # the mean
>>> sigma2 = 9. # the variance
>>> normal_sample = norm.rvs(mu, scale=sqrt(sigma))
>>> alpha = 2.
>>> beta = 15.
>>> invgamma_sample = invgamma.rvs(alpha, scale=beta)
```

Note that when sampling from the normal distribution, we need to set the `scale` parameter to the standard deviation, not the variance.

Problem 1. Write a function that accepts data \mathbf{y} , prior parameters ν , τ^2 , α , and β , and an integer n . Use Gibbs sampling to generate n samples of μ and σ^2 for the exam scores problem.

Test your sampler with priors $\nu = 80$, $\tau^2 = 16$, $\alpha = 3$, and $\beta = 50$, collecting 1000 samples. Plot your samples of μ and your samples of σ^2 . They should both converge quickly, so that both plots look like “fuzzy caterpillars”.

We'd like to look at the posterior marginal distributions for μ and σ^2 . To plot these from the samples, use a kernel density estimator from `scipy.stats`. If our samples of μ are called `mu_samples`, then we can do this with the following code.

```
>>> import numpy as np
>>> from matplotlib import pyplot as plt
>>> from scipy.stats import gaussian_kde

>>> mu_kernel = gaussian_kde(mu_samples)
>>> x = np.linspace(min(mu_samples) - 1, max(mu_samples) + 1, 200)
>>> plt.plot(x, mu_kernel(x))
>>> plt.show()
```

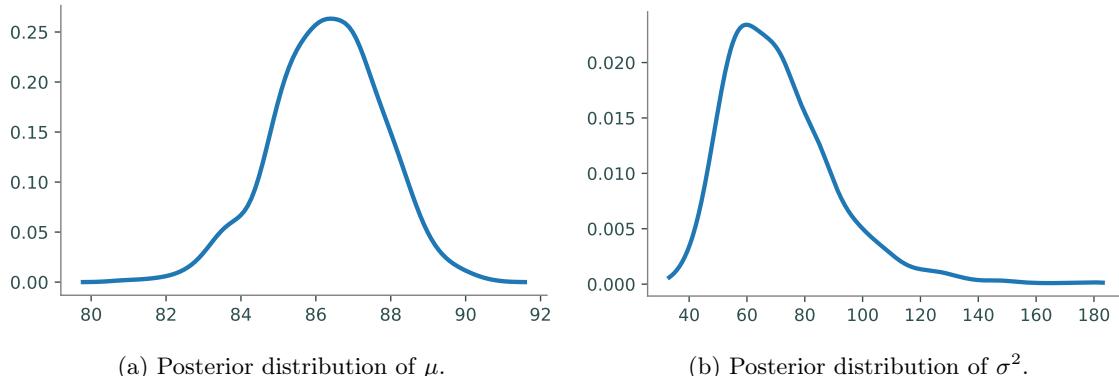


Figure 19.1: Posterior marginal probability densities for μ and σ^2 .

Keep in mind that the plots above are of the posterior distributions of the parameters, not of the scores. If we would like to compute the posterior distribution of a new exam score \tilde{y} given our data \mathbf{y} and prior parameters, we compute what is known as the posterior predictive distribution:

$$\mathbb{P}(\tilde{y}|\mathbf{y}, \lambda) = \int_{\Theta} \mathbb{P}(\tilde{y}|\Theta) \mathbb{P}(\Theta|\mathbf{y}, \lambda) d\Theta$$

where Θ denotes our parameters (in our case μ and σ^2) and λ denotes our prior parameters (in our case ν , τ^2 , α , and β).

Rather than actually computing this integral for each possible \tilde{y} , we can do this by sampling scores from our parameter samples. In other words, sample

$$\tilde{y}_{(t)} \sim N(\mu_{(t)}, \sigma_{(t)}^2)$$

for each sample pair $\mu_{(t)}, \sigma^2_{(t)}$. Now we have essentially drawn samples from our posterior predictive distribution, and we can use a kernel density estimator to plot this distribution from the samples.

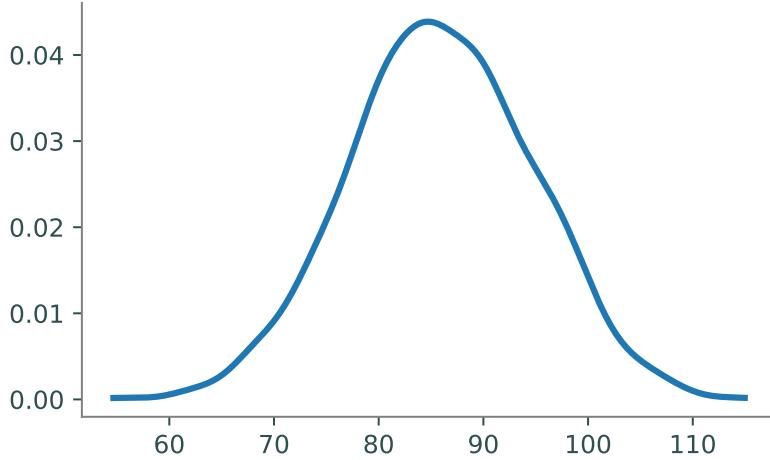


Figure 19.2: Predictive posterior distribution of exam scores.

Problem 2. Plot the kernel density estimators for the posterior distributions of μ and σ^2 . You should get plots similar to those in Figure 19.1.

Next, use your samples of μ and σ^2 to draw samples from the posterior predictive distribution. Plot the kernel density estimator of your sampled scores. Compare your plot to Figure 19.2.

Latent Dirichlet Allocation

Gibbs sampling can be applied to an interesting problem in natural language processing (NLP): determining which topics are prevalent in a document. Latent Dirichlet Allocation (LDA) is a generative model for a collection of text documents. It supposes that there is some fixed vocabulary (composed of V distinct terms) and K different topics, each represented as a probability distribution ϕ_k over the vocabulary, each with a Dirichlet prior β . This means $\phi_{k,v}$ is the probability that topic k is represented by vocabulary term v .

With the vocabulary and topics chosen, the LDA model assumes that we have a set of M documents (each “document” may be a paragraph or other section of the text, rather than a “full” document). The m -th document consists of N_m words, and a probability distribution θ_m over the topics is drawn from a Dirichlet distribution with parameter α . Thus $\theta_{m,k}$ is the probability that document m is assigned the label k . If $\phi_{k,v}$ and $\theta_{m,k}$ are viewed as matrices, their rows sum to one.

We will now iterate through each document in the same manner. Assume we are working on document m , which you will recall contains N_m words. For word n , we first draw a topic assignment $z_{m,n}$ from the categorical distribution θ_m , and then we draw a word $w_{m,n}$ from the categorical distribution $\phi_{z_{m,n}}$. Throughout this implementation, we assume α and β are scalars¹. In summary, we have

1. Draw $\phi_k \sim \text{Dir}(\beta)$ for $1 \leq k \leq K$.
2. For $1 \leq m \leq M$:
 - (a) Draw $\theta_m \sim \text{Dir}(\alpha)$.
 - (b) Draw $z_{m,n} \sim \text{Cat}(\theta_m)$ for $1 \leq n \leq N_m$.
 - (c) Draw $w_{m,n} \sim \text{Cat}(\phi_{z_{m,n}})$ for $1 \leq n \leq N_m$.

We end up with n words which represent document m . Note that these words are not necessarily distinct from one another; indeed, we are most interested in the words that have been repeated the most.

This is typically depicted with graphical plate notation as in Figure 19.3.

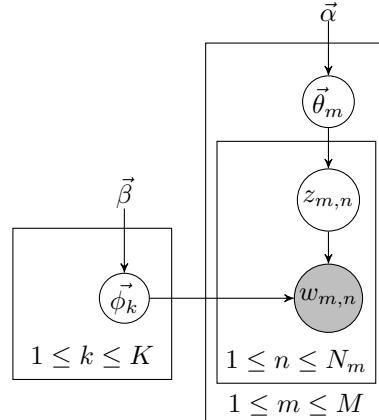


Figure 19.3: Graphical plate notation for LDA text generation.

In the plate model, only the variables $w_{m,n}$ are shaded, signifying that these are the only observations visible to us; the rest are latent variables. Our goal is to estimate each ϕ_k and each θ_m . This will allow us to understand what each topic is, as well as understand how each document is distributed over the K topics. In other words, we want to predict the topic of each document, and also which words best represent this topic. We can estimate these well if we know $z_{m,n}$ for each m, n , collectively referred to as \mathbf{z} . Thus, we need to sample \mathbf{z} from the posterior distribution $\mathbb{P}(\mathbf{z}|\mathbf{w}, \alpha, \beta)$, where \mathbf{w} is the collection words in the text corpus. Unsurprisingly, it is intractable to sample directly from the joint posterior distribution. However, letting $\mathbf{z}_{-(m,n)} = \mathbf{z} \setminus \{z_{m,n}\}$, the conditional posterior distributions

$$\mathbb{P}(z_{m,n} = k | \mathbf{z}_{-(m,n)}, \mathbf{w}, \alpha, \beta)$$

have nice, closed form solutions, making them easy to sample from.

¹The Dirichlet distribution $\text{Dir}(x_1, \dots, x_s, \alpha_1, \dots, \alpha_s)$ usually requires the parameter α to be a vector of length s , but when α is a scalar, it is called the “concentration parameter” and behaves like a vector of length s whose entries are all equal to α .

These conditional distributions have the following form:

$$\mathbb{P}(z_{m,n} = k | \mathbf{z}_{-(m,n)}, \mathbf{w}, \alpha, \beta) \propto \frac{(n_{(k,m,\cdot)}^{-(m,n)} + \alpha)(n_{(k,\cdot,w_{m,n})}^{-(m,n)} + \beta)}{n_{(k,\cdot,\cdot)}^{-(m,n)} + V\beta}$$

where

- $n_{(k,m,\cdot)}$ = the number of words in document m assigned to topic k
- $n_{(k,\cdot,v)}$ = the number of times term $v = w_{m,n}$ is assigned to topic k
- $n_{(k,\cdot,\cdot)}$ = the number of times topic k is assigned in the corpus
- $n_{(k,m,\cdot)}^{-(m,n)} = n_{(k,m,\cdot)} - \mathbf{1}_{z_{m,n}=k}$
- $n_{(k,\cdot,v)}^{-(m,n)} = n_{(k,\cdot,v)} - \mathbf{1}_{z_{m,n}=k}$
- $n_{(k,\cdot,\cdot)}^{-(m,n)} = n_{(k,\cdot,\cdot)} - \mathbf{1}_{z_{m,n}=k}$

Thus, if we simply keep track of these count matrices, then we can easily create a Gibbs sampler over the topic assignments. This is actually a particular class of samplers known as collapsed Gibbs samplers, because we have collapsed the sampler by integrating out θ and ϕ .

We have provided for you the structure of a Python object `LDACGS` with several methods, listed at the end of this lab. The object defines attributes `n_topics`, `alpha`, and `beta` upon initialization. The method `buildCorpus()` then defines attributes `vocab` and `documents`, where `vocab` is a list of strings (terms), and `documents` is a list of dictionaries (a dictionary for each document). For dictionary m in `documents`, each entry is of the form $n : w$, where w is the index in `vocab` of the n^{th} word in document m .

The remainder of this lab will guide you through writing several more methods in order to implement the Gibbs sampler. The first step is to initialize the assignments and create count matrices $n_{(k,m,\cdot)}$, $n_{(k,\cdot,v)}$ and vector $n_{(k,\cdot,\cdot)}$.

Problem 3. Complete the method `initialize()` to initialize as attributes `n_words`, `n_docs`, the count matrices, and the topic assignment dictionary `topics`.

To do this, you will need to initialize `nkm`, `nkv`, and `nk` to be zero arrays of the correct size. Matrix `nkm` corresponds to $n_{(k,m,\cdot)}$, `nkv` to $n_{(k,\cdot,v)}$, and `nk` to $n_{(k,\cdot,\cdot)}$. You will then iterate through each word found in each document. In the second of these for-loops (for each word), you will randomly assign z as an integer from the correct range of topics. Then, you will increment each of the count matrices by 1, given the values for z , m , and w , where w is the index in `vocab` of the n^{th} word in document m . Finally, assign `topics` as given.

The next method fully outlines a sweep of the Gibbs sampler.

Problem 4. Complete the method `_sweep()`.

To do this, iterate through each word of each document. The first part of this method will undo what `initialize()` did by decrementing each of the count matrices by 1. Then, call the method `_conditional()` to use the conditional distribution (instead of the uniform distribution used previously) to pick a more accurate topic assignment z . Finally, repeat what `initialize()` did by incrementing each of the count matrices by 1, but this time using the more accurate topic assignment.

You are now prepared to write the full Gibbs sampler.

Problem 5. Complete the method `sample()`. The argument filename is the name and location of a .txt file, where each line is considered a document. The corpus is built by method `buildCorpus()`, and stopwords are removed (if argument stopwords is provided).

Initialize attributes `total_nkm`, `total_nkv`, and `logprobs` as zero arrays. `total_nkm` and `total_nkv` will be the sums of every `sample_rateth` `nkm` and `nkv` matrix respectively. `logprobs` is of length `burnin + sample_rate * n_samples` and will store each log-likelihood after each sweep of the sampler.

Burn-in the Gibbs sampler. After the burn-in, iterate further for `n_samples` iterations, adding `nkm` and `nkv` to `total_nkm` and `total_nkv` respectively, but only for every `sample_rateth` iteration. Also, compute and save the log-likelihood at each iteration in `logprobs` using the method `_loglikelihood()`.

You should now have a working Gibbs sampler to perform LDA inference on a corpus. Let's test it out on one of Ronald Reagan's State of the Union addresses, found in `reagan.txt`.

Problem 6. Create an `LDACGS` object with 20 topics, letting α and β be the default values. Run the Gibbs sampler, with a burn-in of 100 iterations, accumulating 10 samples, only keeping the results of every 10th sweep. Use `stopwords.txt` as the stopwords file.

Plot the log-likelihoods. How many iterations did it take to burn-in?

We can estimate the values of each ϕ_k and each θ_m as follows:

$$\hat{\phi}_{k,v} = \frac{n_{(k,\cdot,v)} + \beta}{V \cdot \beta + \sum_{v=1}^V n_{(k,\cdot,v)}}$$

$$\hat{\theta}_{m,k} = \frac{n_{(k,m,\cdot)} + \alpha}{K \cdot \alpha + \sum_{k=1}^K n_{(k,m,\cdot)}}$$

We have provided methods `phi` and `theta` that do this for you. We often examine the topic-term distributions ϕ_k by looking at the n terms with the highest probability, where n is small (say 10 or 20). We have provided a method `topterms` which does this for you.

Problem 7. Using the method `topterms()`, examine the topics for Reagan's addresses. If $n_topics = 20$ and $n_samples = 10$, you should get the top 10 words that represent each of the 20 topics. For each topic, decide what these ten words jointly represent, and come up with a label for them.

We can use $\hat{\theta}$ to find the paragraphs in Reagan's addresses that focus the most on each topic. The documents with the highest values of $\hat{\theta}_k$ are those most heavily focused on topic k . For example, if you chose the topic label for topic p to be the Cold War, you can find the five highest values in $\hat{\theta}_p$, which will tell you which five paragraphs are most centered on the Cold War.

Let's take a moment to see what our Gibbs sampler has accomplished. By simply feeding in a group of documents, and with no human input, we have found the most common topics discussed, which are represented by the words most frequently used in relation to that particular topic. The only work that the user has done is to assign topic labels, saying what the words in each group have in common. As you may have noticed, however, these topics may or may not be relevant topics. You might have noticed that some of the most common topics were simply English particles (words such as a, the, an) and conjunctions (and, so, but). Industrial grade packages can effectively remove such topics so that they are not included in the results.

Additional Material

LDACGS Source Code

```

class LDACGS:
    """ Do LDA with Gibbs Sampling. """

    def __init__(self, n_topics, alpha=0.1, beta=0.1):
        """ Initializes attributes n_topics, alpha, and beta. """
        self.n_topics = n_topics
        self.alpha = alpha
        self.beta = beta

    def buildCorpus(self, filename, stopwords_file=None):
        """ Reads the given filename, and using any provided stopwords,
            initializes attributes vocab and documents.

            Vocab is a list of terms found in filename.

            Documents is a list of dictionaries (a dictionary for each
            document); for dictionary m in documents, each entry is of
            the form n:w, where w is the index in vocab of the nth word
            in document m.
        """
        with open(filename, 'r') as infile: # create vocab
            doclines = [line.rstrip().lower().split(' ') for line in infile]
        n_docs = len(doclines)
        self.vocab = list({v for doc in doclines for v in doc})

        if stopwords_file: # if there are stopwords, remove them from vocab
            with open(stopwords_file, 'r') as stopfile:
                stops = stopfile.read().split()
            self.vocab = [x for x in self.vocab if x not in stops]
            self.vocab.sort()

        self.documents = [] # create documents
        for i in range(n_docs):
            self.documents.append({})
            for j in range(len(doclines[i])):
                if doclines[i][j] in self.vocab:
                    self.documents[i][j] = self.vocab.index(doclines[i][j])

    def initialize(self):
        """ Initializes attributes n_words, n_docs, the three count matrices,
            and topics.

            Note that
            n_topics = K, the number of possible topics
        """

```

```

n_docs    = M, the number of documents being analyzed
n_words   = V, the number of words in the vocabulary

To do this, you will need to initialize nkm, nk, and nk
to be zero arrays of the correct size.
Matrix nkm corresponds to n_(k,m,.)
Matrix nk corresponds to n_(k,,v)
Matrix nk corresponds to n_(k,,,.)

You will then iterate through each word found in each document.
In the second of these for-loops (for each word), you will
randomly assign z as an integer from the range of topics.
Then, you will increment each of the count matrices by 1,
given the values for z, m, and w, where w is the index in
vocab of the nth word in document m.
Finally, assign topics as given.

"""
self.n_words = len(self.vocab)
self.n_docs = len(self.documents)

# Initialize the three count matrices.
# The (k, m) entry of self.nkm is the number of words in document m ←
# assigned to topic k.
self.nkm = np.zeros((self.n_topics, self.n_docs))
# The (k, v) entry of self.nkv is the number of times term v is ←
# assigned to topic k.
self.nkv = np.zeros((self.n_topics, self.n_words))
# The (k)-th entry of self.nk is the number of times topic k is ←
# assigned in the corpus.
self.nk = np.zeros(self.n_topics)

# Initialize the topic assignment dictionary.
self.topics = {} # key-value pairs of form (m,i):z

random_distribution = np.ones(self.n_topics) / self.n_topics
for m in range(self.n_docs):
    for i in self.documents[m]:
        # Get random topic assignment, i.e. z = ...
        # Increment count matrices
        # Store topic assignment, i.e. self.topics[(m,i)]=z
        raise NotImplementedError("Problem 3 Incomplete")

def _sweep(self):
    """ Iterates through each word of each document, giving a better
    topic assignment for each word.

    To do this, iterate through each word of each document.
    The first part of this method will undo what initialize() did
    by decrementing each of the count matrices by 1.

```

Then, call the method `_conditional()` to use the conditional distribution (instead of the uniform distribution used previously) to pick a more accurate topic assignment z . Finally, repeat what `initialize()` did by incrementing each of the count matrices by 1, but this time using the more accurate topic assignment.

```
"""
for m in range(self.n_docs):
    for i in self.documents[m]:
        # Retrieve vocab index for i-th word in document m.
        # Retrieve topic assignment for i-th word in document m.
        # Decrement count matrices.
        # Get conditional distribution.
        # Sample new topic assignment.
        # Increment count matrices.
        # Store new topic assignment.
    raise NotImplementedError("Problem 4 Incomplete")
```

`def sample(self, filename, burnin=100, sample_rate=10, n_samples=10, ↵
stopwords_file=None):`

""" Runs the Gibbs sampler on the given filename.

The argument `filename` is the name and location of a `.txt` file, where each line is considered a document.
The corpus is built by method `buildCorpus()`, and stopwords are removed (if argument `stopwords` is provided).

Initialize attributes `total_nkm`, `total_nkv`, and `logprobs` as zero arrays.
`total_nkm` and `total_nkv` will be the sums of every `sample_rate`-th `nkm` and `nkv` matrix respectively.
`logprobs` is of length `burnin + sample_rate * n_samples` and will store each log-likelihood after each sweep of the sampler.

Burn-in the Gibbs sampler.

After the burn-in, iterate further for `n_samples` iterations, adding `nkm` and `nkv` to `total_nkm` and `total_nkv` respectively, but only for every `sample_rate`-th iteration.

Also, compute and save the log-likelihood at each iteration in `logprobs` using the method `_loglikelihood()`.

```
"""
self.buildCorpus(filename, stopwords_file)
self.initialize()

self.total_nzw = np.zeros((self.n_topics, self.n_words))
```

```

self.total_nnz = np.zeros((self.n_docs, self.n_topics))
self.logprobs = np.zeros(burnin + sample_rate * n_samples)

for i in range(burnin):
    # Sweep and store log likelihood.
    raise NotImplementedError("Problem 5 Incomplete")
for i in range(sample_rate * n_samples):
    # Sweep and store log likelihood
    raise NotImplementedError("Problem 5 Incomplete")
    if not i % sample_rate:
        # accumulate counts
        raise NotImplementedError("Problem 5 Incomplete")

def _conditional(self, m, w):
    """ Returns the conditional distribution given m and w.
    Called by _sweep(). """
    dist = (self.nkm[:,m] + self.alpha) * (self.nkv[:,w] + self.beta) / (←
        self.nk + self.beta * self.n_words)
    return dist / np.sum(dist)

def _loglikelihood(self):
    """ Computes and returns the log-likelihood. Called by sample(). """
    lik = 0

    for z in range(self.n_topics):
        lik += np.sum(gammaln(self.nkv[z,:] + self.beta)) - gammaln(np.sum(←
            self.nkv[z,:] + self.beta))
        lik -= self.n_words * gammaln(self.beta) - gammaln(self.n_words * ←
            self.beta)

    for m in range(self.n_docs):
        lik += np.sum(gammaln(self.nkm[:,m] + self.alpha)) - gammaln(np.sum(←
            self.nkm[:,m] + self.alpha))
        lik -= self.n_topics * gammaln(self.alpha) - gammaln(self.n_topics * ←
            self.alpha)

    return lik

def phi(self):
    """ Initializes attribute _phi. Called by topterm(). """
    phi = self.total_nkv + self.beta
    self._phi = phi / np.sum(phi, axis=1)[:,np.newaxis]

def theta(self):
    """ Initializes attribute _theta. Called by topterm(). """
    theta = self.total_nkm + self.alpha
    self._theta = theta / np.sum(theta, axis=1)[:,np.newaxis]

```


20 Gaussian Mixture Models

Lab Objective: Understand the formulation of Gaussian Mixture Models (GMMs) and use the Expectation Maximization algorithm to estimate GMM parameters.

Mixture models are a useful way to combine distributions together that allows us to describe much more complicated distributions than using just the standard list of named distributions. The essential idea of a mixture model is in its name: it is a mixture of several different models, or probability distributions. Each of these model is called a component. Each component has a certain probability associated with it, called its weight, that describes how likely it is for a sample from the model to come from that component. We denote the weight of the i -th component as w_i .

In this lab, we focus on Gaussian Mixture Models, or GMMs for short. In a GMM, each component is a multivariate Gaussian (normal) distribution. Each of these is parameterized by a mean μ_i and a covariance matrix Σ_i .

A GMM with K components thus has parameters $\theta = (w_1, \dots, w_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K)$. We can use the law of total probability to evaluate the density of a GMM, which is given by

$$P(z|\theta) = \sum_{k=1}^K w_k \mathcal{N}(z|\mu_k, \Sigma_k)$$

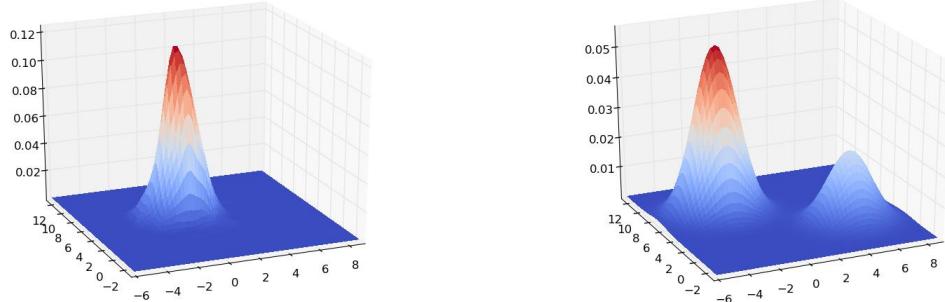
where

$$\mathcal{N}(z|\mu, \Sigma) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{1}{2} (z - \mu)^T \Sigma^{-1} (z - \mu)\right)$$

is the density function of a multivariate normal distribution.

It is important to keep in mind that a GMM does not arise from adding weighted multivariate normal random variables, but rather from weighting the responsibility of each multivariate normal random variable. The first case simply results in a different multivariate normal distribution. Refer to Figure 20.1 for a visualization of these two cases.

Problem 1. Throughout this lab, we will build a GMM class with various methods. Write the `__init__` method for this class. It should accept a parameter for the number of components and optional parameters for the weights, means, and covariance matrices which define the GMM, and store these.^a



(a) Sum of weighted multivariate normal random variables. (b) Weighted mixture of multivariate normal random variables.

Figure 20.1

If we have K components and d dimensions, then the weights should have shape `(K,)`, the means `(K,d)`, and the covariances `(K,d,d)`. The parameters for the k -th component can be found as `weights[k]`, `means[k]`, `covars[k]`.

^aIf we don't have a good guess for the parameters of the GMM to pass into the class, it makes more sense to initialize these from the dataset we are training on, which we will do later in the `fit` method; hence, we let the parameters be optional here.

Problem 2. Write a method `component_logpdf` for your class that accepts a component k and a point \mathbf{z} and computes

$$\log w_k + \log \mathcal{N}(z|\mu_k, \Sigma_k),$$

the logarithm of the contribution of the k -th component of the pdf. Also write a method `pdf` that accepts a point \mathbf{z} and returns the probability density of the whole GMM at that point.

Hint: `scipy.stats.multivariate_normal.pdf` and `scipy.stats.multivariate_normal.logpdf` can be used to efficiently evaluate the multivariate normal pdf.

To test your functions, create the following GMM:

```
gmm = GMM(n_components = 2,
           weights = np.array([0.6, 0.4]),
           means = np.array([[[-0.5, -4.0], [0.5, 0.5]]]),
           covars = np.array([
               [[1, 0], [0, 1]],
               [[0.25, -1], [-1, 8]],
           ]))
```

Your functions should give the following output:

```
>>> gmm.pdf(np.array([1.0, -3.5]))
0.05077912539363083
# Component 0
>>> gmm.component_logpdf(0, np.array([1.0, -3.5]))
-3.598702690175336
# Component 1
>>> gmm.component_logpdf(1, np.array([1.0, -3.5]))
-3.7541677982835004
```

Note that since this GMM is 2-dimensional, the input point must be an array of length 2.

In order to draw a value from a mixture model, we must first draw a variable $X \sim \text{Cat}(w_1, \dots, w_K)$ that represents which component the sample comes from. We can then draw the sample $Z \sim \mathcal{N}(\mu_X, \sigma_X)$.

Problem 3. Write a method `draw` for the GMM class that randomly draws from the model. If m points are drawn and the GMM is d -dimensional, the returned array should have shape (m, d) .

Draw a sample of 10,000 points from the GMM defined in Problem 2. Plot the pdf of the GMM (using `plt.pcolormesh`) and a hexbin plot of the drawn points. How do the plots compare?

The following code can be used to plot the pdf:

```
## Create the grid to plot on
x = np.linspace(-8,8,100)
y = np.linspace(-8,8,100)
X, Y = np.meshgrid(x, y)
## Calculate the pdf at each point
# If your pdf function uses array broadcasting, you can do the following:
Z = gmm.pdf(np.dstack((X,Y)))
# Otherwise, you need to iterate over each point:
Z = np.array([
    gmm.pdf([X[i,j], Y[i,j]]) for j in range(100)
    ] for i in range(100)
)
## Create the plot
plt.pcolormesh(X, Y, Z, shading='auto')
```

We now consider how to estimate the parameters of a GMM given some observed data $Z = z_1, \dots, z_n$. Ordinarily, a good approach would be to try to directly maximize the log-likelihood

$$l(\theta) = \sum_{i=1}^n \log \sum_{j=1}^K w_j \mathcal{N}(z_i | \mu_j, \Sigma_j).$$

However, this expression is very difficult to deal with using standard optimization methods, particularly because of the sum inside of the logarithm. A good alternative in this case is the expectation maximization (EM) algorithm. This is an iterative algorithm, where each step consists of maximizing a function that is designed to approximate the log-likelihood while being much easier to maximize.

Each iteration consists of two steps, the E-step and the M-step. Suppose our estimated parameters at the t -th iteration are $\theta^t = (w_1^t, \dots, w_K^t, \mu_1^t, \dots, \mu_K^t, \Sigma_1^t, \dots, \Sigma_K^t)$. Note that t is an index, not an exponent. For each data point $z_i, 1 \leq i \leq n$ and each component $1 \leq k \leq K$, the E-step consists of computing

$$\begin{aligned} q_i^t(k) &= P(X_i = k | z_i, \theta^t) \\ &= \frac{P(z_i | X_i = k, \theta^t)}{P(z_i | \theta^t)} \\ &= \frac{w_k^t \mathcal{N}(z_i | \mu_k^t, \Sigma_k^t)}{\sum_{k'=1}^K w_{k'}^t \mathcal{N}(z_i | \mu_{k'}^t, \Sigma_{k'}^t)} \end{aligned}$$

In order to accurately compute this quantity, however, we need to be more careful. It is possible that due to floating point underflow¹ that each term $w_{k'}^t \mathcal{N}(z_i | \mu_{k'}, \Sigma_{k'})$ in the sum in the denominator becomes zero, which is a major problem. This particularly happens if the exponents in the multivariate normal densities all are large negative numbers. To avoid this problem, we can rescale the numerator and denominator. Let

$$\ell_{i,k} = \log w_k^t + \log \mathcal{N}(z_i | \mu_k^t, \Sigma_k^t),$$

the logarithm of each term in the denominator. For each data point z_i , we can find

$$L_i = \max_{k'} \ell_{i,k'},$$

the largest of these logarithms. Then, we can rewrite the quantity we want to calculate as

$$\begin{aligned} q_i^t(k) &= \frac{w_k^t \mathcal{N}(z_i | \mu_k^t, \Sigma_k^t)}{\sum_{k'=1}^K w_{k'}^t \mathcal{N}(z_i | \mu_{k'}^t, \Sigma_{k'}^t)} \\ &= \frac{e^{\ell_{i,k}}}{\sum_{k'=1}^K e^{\ell_{i,k'}}} \\ &= \frac{e^{\ell_{i,k}} e^{-L_i}}{\sum_{k'=1}^K e^{\ell_{i,k'}} e^{-L_i}} \\ &= \frac{e^{\ell_{i,k} - L_i}}{\sum_{k'=1}^K e^{\ell_{i,k'} - L_i}}. \end{aligned}$$

This rescaling makes the largest term in the denominator equal to 1, so computing $q_i^t(k)$ in this way avoids underflow problems. Note that for the computation of any individual $q_i^t(k)$, the value L_i is a scalar that is the same for all components; however, you will have as many of these values as you have data points.

¹As a refresher, one way that floating point numbers are limited is that they cannot represent positive numbers arbitrarily close to zero; at some point, if the number in a computation becomes too small, the computer is forced to round it to zero, which is called underflow. The threshold is about 10^{-323} for the 64-bit floating point numbers used in python. Even if underflow does not occur, very small floating points have greatly reduced precision, so it is generally good to avoid using them.

Problem 4. Write a method `_compute_e_step` that calculates the $q_i^t(k)$ as given by the E-step, given a collection of observations. Be sure to do the calculation in a way that avoids underflow, and use array broadcasting when possible.

Now that we have the $q_i^t(k)$, we can perform the M-step. This step consists of maximizing the function

$$Q^t(\theta) = \sum_{i=1}^n \sum_{k=1}^K q_i^t(k) \log w_k^t \mathcal{N}(z_i | \mu_k, \Sigma_k)$$

We then set

$$\theta^{t+1} = \operatorname{argmax}_{\theta} Q^t(\theta)$$

and iterate until the method appears to converge. In the case of GMMs, the maximizer θ^{t+1} of $Q^t(\theta)$ is given by

$$\begin{aligned} w_k^{t+1} &= \frac{1}{n} \sum_{i=1}^n q_i^t(k) \\ \mu_k^{t+1} &= \frac{\sum_{i=1}^n q_i^t(k) z_i}{\sum_{i=1}^n q_i^t(k)} \\ \Sigma_k^{t+1} &= \frac{\sum_{i=1}^n q_i^t(k) (z_i - \mu_k^{t+1})(z_i - \mu_k^{t+1})^\top}{\sum_{i=1}^n q_i^t(k)} \end{aligned}$$

For details on the derivation of the maximizer, refer to the Volume 3 textbook.

Problem 5. Write a method `_compute_m_step` for your GMM class that performs a single iteration of the EM algorithm. Return the updated parameters, as given by the M-step. Be sure to use array broadcasting when possible.

Problem 6. Write a `fit` method for your GMM class.

First, if the GMM's parameters are uninitialized (set to `None`), initialize the parameters of the components. We want to do this in a way that the algorithm starts with reasonable values for the dataset. A good way to initialize the means is to randomly select points from the dataset. The covariance matrices can be initialized as diagonal matrices based on the variance of the data. Ensure that the weights you choose add up to 1.

Then, perform the expectation maximization algorithm. Use the functions you created in Problems 4 and 5 to calculate the parameters at each step. Repeat until the parameters converge. Use the following to measure the change in the parameters with each iteration:

```
change = (np.max(np.abs(new_weights - old_weights))
          + np.max(np.abs(new_means - old_means))
          + np.max(np.abs(new_covars - old_covars)))
```

Problem 7. The file `problem7.npy` contains a collection of data drawn from a two-dimensional GMM. Train a GMM on this data with `n_components=3`. Plot the pdf of your trained GMM (in the same way as in Problem 3), as well as a hexbin plot of the data. Your class should take less than 15 seconds to train on this dataset.

Clustering with GMMS

An important use of mixture models is for clustering. The objective of clustering is to take an unlabeled dataset and separate it into some number of clusters, which can then be labeled. This is an instance of unsupervised learning, as it is a machine learning task where the training algorithm does not need the true answers (in this case, the actual clusters).

In order to cluster a dataset using a GMM, we first need to train the GMM on that data. Then, we can assign each point a label by finding which component has the largest contribution to the pdf there. Written symbolically, for a data point z , we have

$$\text{Cluster}(z) = \operatorname{argmax}_k w_k \mathcal{N}(z|\mu_k, \Sigma_k).$$

Note that the number of clusters (components) is a hyperparameter that must be selected before a GMM is trained. In general, cross-validation or some other method must be used to find the right number of clusters.

Problem 8. Write a `predict` method for your class. Given a set of data points, return which cluster has the highest pdf density for each data point.

The file `classification.npz` contains a set of 3-dimensional data points (`X`) and their labels (`y`). Use your class with `n_components=4` to cluster the data. Plot the points with the predicted and actual labels, and compute and return your model's accuracy. Your class should take less than 30 seconds to train on this dataset.

Note that the labels may be permuted; for instance, your model might cluster the points correctly, but swap the labels of clusters 1 and 2 compared to the true labels. The model would still be considered accurate in this case; we only care what the clusters are, not how the model labels them. To resolve this problem, we need to find the permutation of the labels that results in the highest accuracy. The following function does this in a way that is more efficient than directly checking all permutations:

```
from scipy.optimize import linear_sum_assignment
from sklearn.metrics import confusion_matrix

def get_accuracy(pred_y, true_y):
    """
    Helper function to calculate the actually clustering accuracy,
    accounting for the possibility that labels are permuted.
    """
    # Compute confusion matrix
    cm = confusion_matrix(pred_y, true_y)
    # Find the arrangement that maximizes the score
    r_ind, c_ind = linear_sum_assignment(cm, maximize=True)
```

```
return np.sum(cm[r_ind, c_ind]) / np.sum(cm)
```

For convenience, a method `fit_predict` for the class is also included in the specifications file that calls both `fit` and `predict` to make the clustering process simpler.

Clustering with GMMs is closely related to the K-means algorithm. In fact, K-means can be viewed as a special case of GMMs. We now compare the effectiveness of GMMs for classification on this dataset with K-means, as well as comparing to sklearn's implementation.

Problem 9. Again using `classification.npz`, compare your class, sklearn's GMM implementation, and sklearn's K-means implementation for speed of training and for accuracy of the resulting clusters. Print your results. Be sure to check for permuted labels.

You should find that sklearn's GMM is actually faster on this dataset than K-means. This is in part because the dataset is rather low-dimensional. As the dimension of the dataset grows, GMMs suffer computationally from the curse of dimensionality much more than the K-means algorithm.

21

Discrete Hidden Markov Models

Lab Objective: Understand how to use discrete Hidden Markov Models.

In this lab, we explore Hidden Markov Models (HMMs) with discrete state and observation spaces. Assume the state space \mathcal{X} and observation space \mathcal{Z} are finite sets where $|\mathcal{X}| = n$ and $|\mathcal{Z}| = m$. In addition, a discrete state-space HMM has parameters $\theta = (\pi, A, B)$ and an observation sequence z . We would like to answer three questions about an HMM:

1. What is the likelihood that our model generated the observation sequence? In other words, what is $P(z | \theta)$?
2. What is the most likely state sequence x to have generated z , given θ ?
3. How can we choose the parameters θ that maximize $P(z | \theta)$?

The first question is answered using the forward pass algorithm. For the second question, the approach taken in this lab will be to find the state sequence maximizing the expected number of correct states. The third question is an example of unsupervised learning, since we are attempting to learn (or fit) model parameters using data (the observation sequence z) that is devoid of human-provided labels (the corresponding state sequence); the algorithm does not rely on human supervision or input.

In this context $\theta = (\pi, A, B)$, where π is a stochastic vector of length n (the initial state distribution), A is a $n \times n$ column-stochastic matrix (the state transition model), and B is a $m \times n$ column-stochastic matrix (the state observation model). Further, z is a vector of length T with values in the set $\mathcal{Z} = \{0, 1, 2, \dots, m - 1\}$.

Throughout this lab, we will be using the following toy HMM to verify your code.

```
>>> # toy HMM example to be used to check answers
>>> pi = np.array([.6, .4])
>>> A = np.array([[.7, .4], [.3, .6]])
>>> B = np.array([[.1, .7], [.4, .2], [.5, .1]])
>>> z = np.array([0, 1, 0, 2])
```

Problem 1. To start off your implementation of the HMM, define a class object which you should call `HMM`. Then add the initialization method, storing the `self` aspects `pi`, `A`, and `B` as `None` objects. You will be adding methods throughout the remainder of the lab.

The Forward Pass

Our first task is to efficiently compute $P(\mathbf{z} | \boldsymbol{\theta})$. We can do this using the forward pass algorithm.

First, let $\alpha_t(i) = P(z_0, \dots, z_t, x_t = i | \boldsymbol{\theta})$. Then using the law of total probability and $\alpha_t(i)$, we can efficiently compute $P(\mathbf{z} | \boldsymbol{\theta})$ as

$$P(\mathbf{z} | \boldsymbol{\theta}) = \sum_{i \in \mathcal{X}} \alpha_{T-1}(i).$$

Now we must use a rescaled version of the forward pass to prevent the $\alpha_t(i)$'s from becoming too small as t gets large. Let \odot denote the Hadamard (entry-wise) product of arrays. The algorithm is as below.

Algorithm 21.1 Forward Pass Algorithm

```

1: procedure Forward Pass Algorithm
2:   for t=0 do
3:     Set  $\hat{\alpha}_0(i) = \pi_i \cdot B_{z_0 i}$ ,  $\forall i \in \mathcal{X}$ 
4:     Let  $c_0 = 1 / (\sum_{j \in \mathcal{X}} \hat{\alpha}_0(j))$ 
5:     Set  $\hat{\alpha}_0(i) = c_0 \odot \hat{\alpha}_0(i)$ ,  $\forall i \in \mathcal{X}$ 
6:   for t=1, ..., T - 1 do
7:     Compute  $\tilde{\alpha}_t(i) = \sum_{j \in \mathcal{X}} \hat{\alpha}_{t-1}(j) \cdot A_{ij} \cdot B_{z_t i}$ ,  $\forall i \in \mathcal{X}$ 
8:     Compute  $c_t = 1 / (\sum_{j \in \mathcal{X}} \tilde{\alpha}_t(j))$ 
9:     Rescale by setting  $\hat{\alpha}_t(i) = c_t \cdot \tilde{\alpha}_t(i)$ ,  $\forall i \in \mathcal{X}$ 
```

The matrix $\hat{\alpha}$ will be of use when fitting parameters, but we can compute the desired log probability using the scaling factors c_t as follows:

$$\log P(\mathbf{z} | \boldsymbol{\theta}) = - \sum_{t=0}^{T-1} \log c_t.$$

Problem 2. Implement the forward pass by adding the following method to your class:

```

def _forward(self, z):
    """
    Compute the scaled forward probability matrix and scaling factors.

    Parameters
    -----
    z : ndarray of shape (T,)
        The observation sequence

    Returns
    -----
    None
    """

    # Initialize alpha and c
    alpha = np.zeros_like(z)
    c = 1.0

    # Compute initial forward probabilities
    alpha[0] = self.pi * self.B[z[0], :]
    c = 1.0 / np.sum(alpha[0])

    # Compute scaled forward probabilities and scaling factors
    for t in range(1, len(z)):
        alpha[t] = np.dot(alpha[t-1], self.A) * self.B[z[t], :]
        c *= alpha[t].sum()

    # Return the scaled forward probabilities and scaling factor
    return alpha, c
```

```

-----
alpha : ndarray of shape (T, n)
    The scaled forward probability matrix
c : ndarray of shape (T,)
    The scaling factors c = [c_0, c_1, ..., c_{T-1}]
"""
pass

```

To verify that your code works, you should get the following output using the toy HMM:

```

>>> h = HMM()
>>> h.pi = pi
>>> h.A = A
>>> h.B = B
>>> alpha, c = h._forward(z)
>>> print(-np.log(c).sum()) # the log prob of observation
-4.6429135909

```

The Backward Pass

The backward pass produces values that can be used to calculate the most likely state sequence corresponding to an observation sequence.

We compute a scaled backward probability matrix $\hat{\beta}$ of dimension $T \times n$ as follows:

Algorithm 21.2 Backward Pass Algorithm

- 1: **procedure** Backward Pass Algorithm
 - 2: When $t = T - 1$, set $\hat{\beta}_{T-1}(i) = c_{T-1}$, $\forall i \in \mathcal{X}$
 - 3: **for** $t = T - 2, \dots, 0$ **do**
 - 4: Compute $\tilde{\beta}_t(j) = \sum_{i \in \mathcal{X}} A_{ij} \cdot \hat{\beta}_{t+1}(i) \cdot B_{z_{t+1}i}$, $\forall j \in \mathcal{X}$
 - 5: Rescale by setting $\hat{\beta}_t(i) = c_t \cdot \tilde{\beta}_t(i)$, $\forall i \in \mathcal{X}$
-

Problem 3. Implement the backward pass by adding the following method to your class:

```

def _backward(self, z, c):
    """
    Compute the scaled backward probability matrix.

    Parameters
    -----
    z : ndarray of shape (T,)
        The observation sequence
    c : ndarray of shape (T,)
        The scaling factors from the forward pass

```

```

Returns
-----
beta : ndarray of shape (T, n)
    The scaled backward probability matrix
"""
pass

```

Using the same toy example as before, your code should produce the following output:

```

>>> beta = h._backward(z, c)
>>> print(beta)
[[ 3.1361635   2.89939354]
 [ 2.86699344  4.39229044]
 [ 3.898812    2.66760821]
 [ 3.56816483   3.56816483]]

```

Computing the ξ and γ Probabilities

Having implemented both parts of the forward-backward algorithm, we are closing in on the solution to question three, namely that of fitting parameters θ that maximize $P(\mathbf{z} | \theta)$. At this stage, we combine the information accumulated in the forward and backward algorithms to produce a three-dimensional array ξ of shape $(T - 1) \times n \times n$ whose entries are related to $P(\mathbf{x}_t = i, \mathbf{x}_{t+1} = j | \mathbf{z}, \theta)$, as well as a $T \times n$ matrix γ whose entries are related to $P(\mathbf{x}_t = i | \mathbf{z}, \theta)$. The relevant formulae are

$$\xi_t(i, j) = \hat{\alpha}_t(i) A_{j,i} B_{z_{t+1},j} \hat{\beta}_{t+1}(j)$$

for $t = 0, \dots, T - 1$ and $i, j \in \mathcal{X}$,

$$\gamma_t(i) = \hat{\alpha}_t(i) \hat{\beta}_t(i) / c_t$$

for $t = 0, \dots, T - 1$ and $i \in \mathcal{X}$.

Problem 4. Add the following method to your class to compute the ξ and γ probabilities.

```

def _xi(self, z, alpha, beta, c):
    """
    Compute the xi and gamma probabilities.

    Parameters
    -----
    z : ndarray of shape (T,)
        The observation sequence
    alpha : ndarray of shape (T, n)
        The scaled forward probability matrix from the forward pass
    beta : ndarray of shape (T, n)
        The scaled backward probability matrix from the backward pass
    c : ndarray of shape (T,)

```

The scaling factors from the forward pass

Returns

```
xi : ndarray of shape (T-1, n, n)
    The xi probability array
gamma : ndarray of shape (T, n)
    The gamma probability array
"""
pass
```

While writing a triply-nested loop may be the simplest way to convert the formula into code, it is possible to use array broadcasting to eliminate two of the loops, which will speed up your code.

Check your code by making sure it produces the following output, using the same toy example as before.

```
>>> xi, gamma = h._xi(z, alpha, beta, c)
>>> print(xi)
[[[ 0.14166321  0.0465066 ]
 [ 0.37776855  0.43406164]]

 [[ 0.17015868  0.34927307]
 [ 0.05871895  0.4218493 ]]

 [[ 0.21080834  0.01806929]
 [ 0.59317106  0.17795132]]]

>>> print(gamma)
[[ 0.18816981  0.81183019]
 [ 0.51943175  0.48056825]
 [ 0.22887763  0.77112237]
 [ 0.8039794   0.1960206 ]]
```

Choosing Better Parameters

After running the forward-backward algorithm and computing the ξ probabilities, we are now in a position to choose new parameters $\boldsymbol{\theta}' = (\boldsymbol{\pi}', A', B')$ that increase the probability of observing our data, i.e.

$$P(\mathbf{z} | \boldsymbol{\theta}') \geq P(\mathbf{z} | \boldsymbol{\theta}).$$

The update formulae are given by

$$\begin{aligned}\pi' &= \gamma_0(i) \\ A'_{i,j} &= \frac{\sum_{t=0}^{T-2} \xi_t(i, j)}{\sum_{t=0}^{T-2} \gamma_t(j)} \\ B'_{i,j} &= \frac{\sum_{t=0}^{T-1} \gamma_t(j) \delta_{z_t=i}}{\sum_{t=0}^{T-1} \gamma_t(j)}\end{aligned}$$

where $\delta_{z_t=i}$ equals 1 if $z_t = i$, and it equals 0 otherwise.

Problem 5. Implement the parameter update step by adding the following method to your class:

```
def _estimate(self, z, xi, gamma):
    """
    Estimate better parameter values and update self.pi, self.A, and
    self.B in place.

    Parameters
    -----
    z : ndarray of shape (T,)
        The observation sequence
    xi : ndarray of shape (T-1, n, n)
        The xi probability array
    gamma : ndarray of shape (T, n)
        The gamma probability array
    """
    pass
```

Verify that your code produces the following output on the toy HMM from before:

```
h._estimate(z, xi, gamma)
>>> print(h.pi)
[ 0.18816981  0.81183019]
>>> print(h.A)
[[ 0.55807991  0.49898142]
 [ 0.44192009  0.50101858]]
>>> print(h.B)
[[ 0.23961928  0.70056364]
 [ 0.29844534  0.21268397]
 [ 0.46193538  0.08675238]]
```

Fitting the Model

We are now ready to put everything together into a learning algorithm. Given a sequence of observations, a maximum number of iterations K , and a convergence tolerance threshold ε , we fit a HMM model using the following procedure:

Algorithm 21.3 HMM Fitting Algorithm

```

1: procedure HMM Fitting Algorithm
2:   Randomly initialize parameters  $\boldsymbol{\theta} = (\pi, A, B)$ .
3:   Compute  $\log P(\mathbf{z} | \boldsymbol{\theta})$ 
4:   for  $i = 0, 1, \dots, K - 1$  do
5:     Run forward pass
6:     Run backward pass
7:     Compute  $\xi$  and  $\gamma$  probabilities
8:     Update model parameters
9:     Compute  $\log P(\mathbf{z} | \boldsymbol{\theta})$  according to new parameters
10:    if Change in log probabilities is less than  $\varepsilon$  then
11:      break
12:    else
13:      continue

```

The most convenient way to randomly initialize stochastic matrices is to draw from the Dirichlet distribution, which produces vectors with nonnegative entries that sum to 1. The following Python code initializes M , π , A , and B using this technique:

```

>>> # assume N is defined
>>> # define M to be the number of distinct observed states
>>> M = len(set(obs))
>>> pi = np.random.dirichlet(np.ones(N))
>>> A = np.random.dirichlet(np.ones(N), size=N).T
>>> B = np.random.dirichlet(np.ones(M), size=N).T

```

The learning algorithm is essentially an optimization over the parameter space (i.e. the space of tuples of stochastic arrays having the proper dimensions) with respect to the objective function $P(\mathbf{z} | \boldsymbol{\theta})$. The algorithm is guaranteed to increase the objective function at each iteration, so it is sure to converge. However, the objective function is riddled with local maxima, and so the outcome depends heavily on the randomly selected starting values for π , A , and B . Figure 21.1 illustrates the issues involved. The log probability stays approximately constant for the first 100 iterations. This indicates that the algorithm is not exploring the parameter space enough, and the parameters found at the 100-th iteration are virtually the same as those found at the first or second iteration. After the first 100 iterations, however, the algorithm is finally able to explore more of the parameter space and hence make better progress toward increasing the objective function. The moral of the story is that you may need to train the HMM a few times, using different starting values, and then keep the model that has the highest log likelihood.

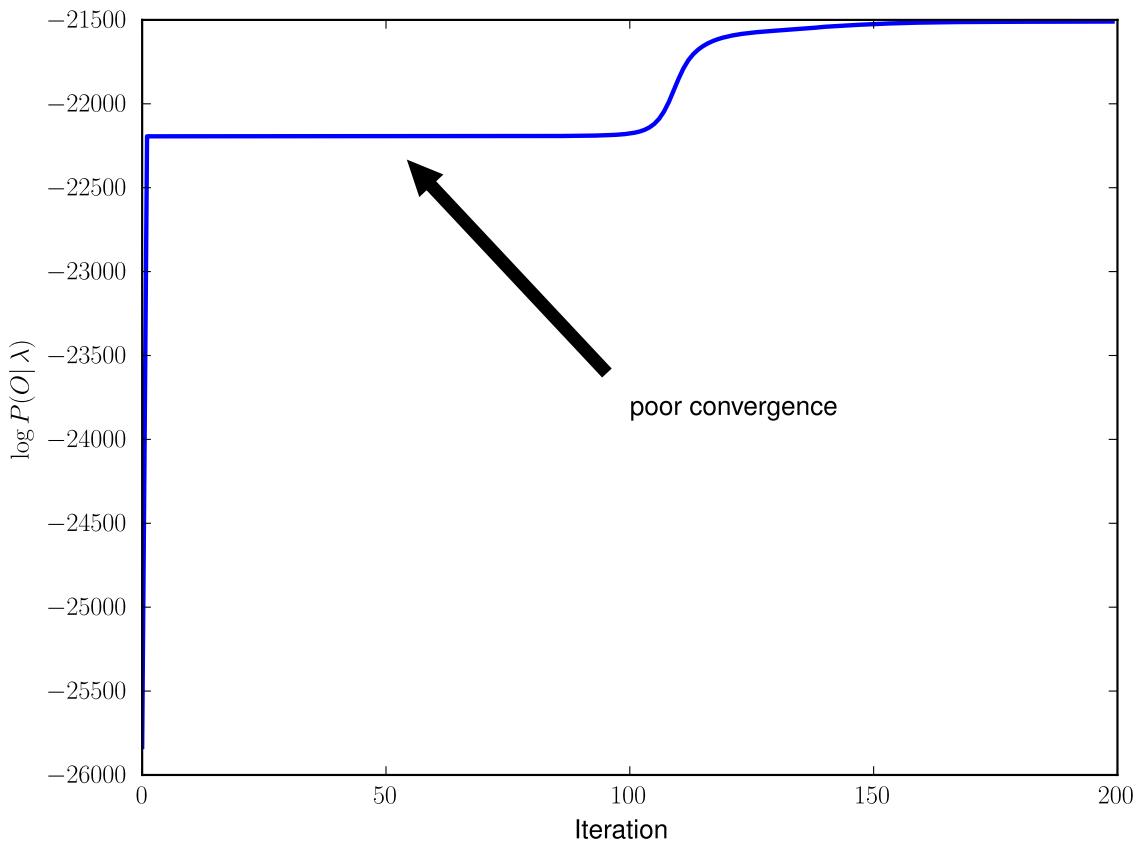


Figure 21.1: The log probabilities for a HMM trained on the Declaration of Independence data with 200 iterations. It takes over 100 iterations for the algorithm to work itself out of a poor local maximum.

Problem 6. Implement the learning algorithm by adding the following method to your class:

```
def fit(self, z, pi, A, B, max_iter=100, tol=1e-4):
    """
    Fit the model parameters to a given observation sequence.

    Parameters
    -----
    z : ndarray of shape (T,)
        Observation sequence on which to train the model.
    pi : Stochastic ndarray of shape (n,)
        Initial state distribution
    A : Stochastic ndarray of shape (n, n)
        Initial state transition matrix
    B : Stochastic ndarray of shape (m, n)
        Initial state observation matrix
```

```

max_iter : int
    The maximum number of iterations to take
tol : float
    The convergence threshold for change in log-probability
"""
# initialize self.pi, self.A, self.B
# run the iteration
pass

```

We now turn to the data found in the file `declaration.txt`. This file contains the text of the Declaration of Independence. We will use the sequence of characters (after stripping out punctuation and converting everything to lower-case) as our observation sequence. In order to convert the raw text into a useable data structure, we need to read in the file, process the string as necessary, and then map the characters to integer values. We provide helper code below to accomplish this task for various files in various languages:

```

>>> import numpy as np
>>> import string
>>> import codecs

>>> def vec_translate(a, my_dict):
>>>     # translate numpy array from symbols to state numbers or vice versa
>>>     return np.vectorize(my_dict.__getitem__)(a)

>>> def prep_data(filename):
>>>     # Get the data as a single string
>>>     with codecs.open(filename, encoding='utf-8') as f:
>>>         data=f.read().lower() # and convert to all lower case

>>>     # remove punctuation and newlines
>>>     remove_punct_map = {ord(char):
>>>                         None for char in string.punctuation+"\n\r"}
>>>     data = data.translate(remove_punct_map)

>>>     # make a list of the symbols in the data
>>>     symbols = sorted(list(set(data)))

>>>     # convert the data to a NumPy array of symbols
>>>     a = np.array(list(data))

>>>     # make a conversion dictionary from symbols to state numbers
>>>     symbols_to_obsstates = {x:i for i,x in enumerate(symbols)}

>>>     # convert the symbols in a to state numbers
>>>     obs_sequence = vec_translate(a,symbols_to_obsstates)

>>>     return symbols, obs_sequence

```

Now apply this helper code to `declaration.txt`.

```
>>> symbols, obs = prep_data('declaration.txt')
```

Problem 7. You are now ready to train a HMM using the Declaration of Independence data. Use $N = 2$ states and $M = \text{len}(\text{set}(\text{obs})) = 27$ observation values (26 lower case characters and 1 whitespace character), and run for 150 iterations with the default value for `tol`. Generally speaking, if you converge to a log probability greater than -21550 , then you have reached an acceptable set of parameters for this dataset.

Once the learning algorithm converges, analyze the state observation matrix B . Note which rows correspond to the largest and smallest probability values in each column of B , and check the corresponding characters. The code below displays typical results for a well-converged HMM. Note that the `u` before the `"` indicates that the string should be unicode, which will be required for languages other than English.

```
>>> for i in range(len(h.B)):
>>>     print(u"{}0", {1:0.4f}, {2:0.4f}"
           .format(symbols[i], h.B[i,0], h.B[i,1]))
    , 0.0051, 0.3324
a, 0.0000, 0.1247
c, 0.0460, 0.0000
b, 0.0237, 0.0000
e, 0.0000, 0.2245
d, 0.0630, 0.0000
g, 0.0325, 0.0000
f, 0.0450, 0.0000
i, 0.0000, 0.1174
h, 0.0806, 0.0070
k, 0.0031, 0.0005
j, 0.0040, 0.0000
m, 0.0360, 0.0000
l, 0.0569, 0.0001
o, 0.0009, 0.1331
n, 0.1207, 0.0000
q, 0.0015, 0.0000
p, 0.0345, 0.0000
s, 0.1195, 0.0000
r, 0.1062, 0.0000
u, 0.0000, 0.0546
t, 0.1600, 0.0000
w, 0.0242, 0.0000
v, 0.0185, 0.0000
y, 0.0147, 0.0058
x, 0.0022, 0.0000
```

```
z, 0.0010, 0.0000
```

What do you notice about the second column of B ? It seems that the HMM has detected a vowel state and a consonant state, without any prior input from an English speaker. Interestingly, the whitespace character is grouped together with the vowels. A HMM can also detect the vowel/consonant distinction in other languages.

Problem 8. Repeat the previous calculation with 3 hidden states and again with 4 hidden states. Interpret/explain your results.

Hint: with 3 hidden states, your print statement will look like the following:

```
>>> print(u"{}0, {}1:0.4f, {}2:0.4f, {}3:0.4f"
       .format(symbols[i], h.B[i,0], h.B[i,1], h.B[i,2]))
```

Now we turn to the Russian file `WarAndPeace.txt`, which is a small subset of the book War and Peace by Tolstoy.

```
>>> symbols, obs = prep_data('WarAndPeace.txt')
```

Problem 9. Repeat the same calculation with `WarAndPeace.txt` for 2 and 3 hidden states. Interpret/explain your results. Which Cyrillic characters appear to be vowels?

22

Speech Recognition using CDHMMs

Lab Objective: Understand how speech recognition via CDHMMs works, and implement a simplified speech recognition system.

22.0.1 Continuous Density Hidden Markov Models

Some of the most powerful applications of Hidden Markov Models, speech and voice recognition, result from allowing the observation space to be continuous instead of discrete. These are called Continuous Density Hidden Markov Models (CDHMMs), and they have two standard formulations: Gaussian HMMs and Gaussian Mixture Model HMMs (GMMHMMs). The former is a special case of the latter, so we will just discuss GMMHMMs in this lab.

In order to understand GMMHMMs, we need to be familiar with a particular continuous, multivariate distribution called a mixture of Gaussians. A mixture of Gaussians is a distribution composed of several Gaussian distributions with corresponding weights. Such a distribution is parameterized by the number of mixture components K (where each component is a Gaussian distribution), the dimension M of the normal distributions involved, a collection of component weights $\{c_1, \dots, c_K\}$ that are nonnegative and sum to 1, and a collection of mean and covariance parameters $\{(\mu_1, \Sigma_1), \dots, (\mu_K, \Sigma_K)\}$ for each Gaussian component. To sample from a mixture of Gaussians, one first chooses the mixture component i according to the probability weights $\{c_1, \dots, c_K\}$, and then one samples from the normal distribution $\mathcal{N}(\mu_k, \Sigma_k)$. The probability density function for a mixture of Gaussians is given by

$$p(\mathbf{z}|\theta) = \sum_{k=1}^K c_k N(\mathbf{z}; \mu_k, \Sigma_k),$$

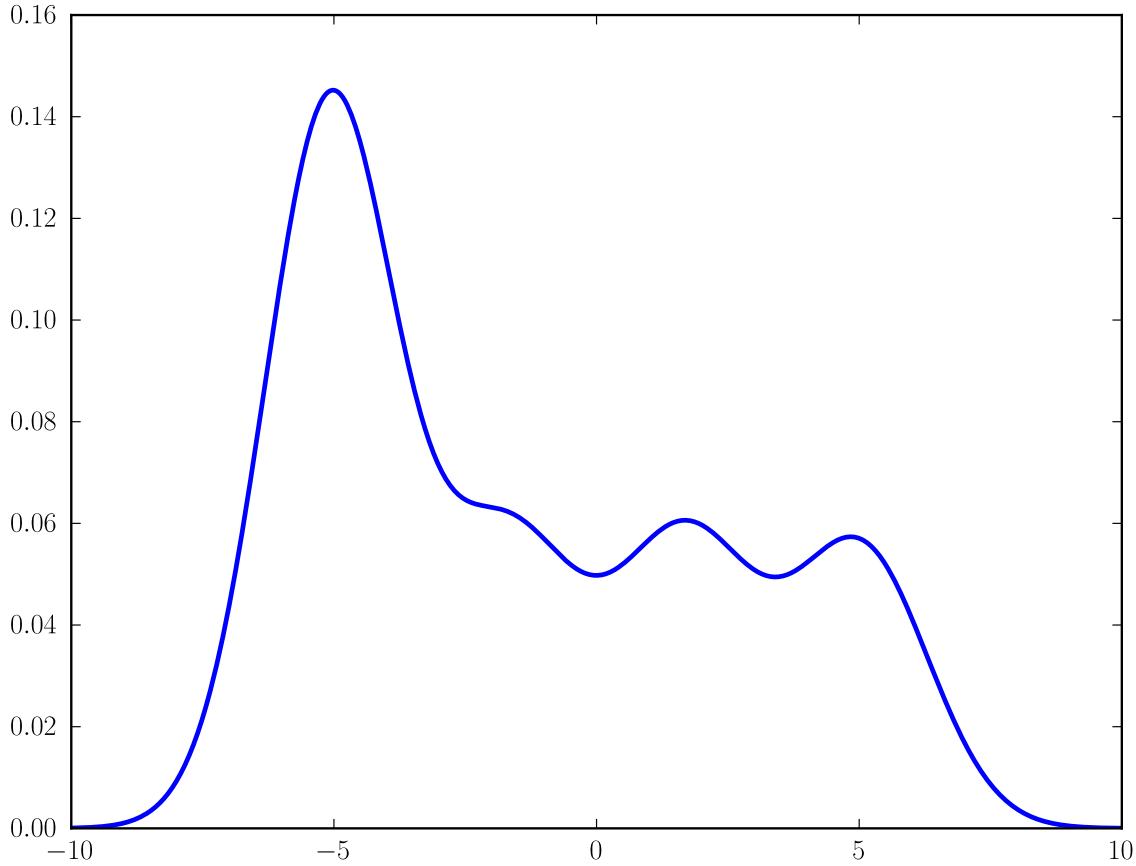


Figure 22.1: The probability density function of a mixture of Gaussians with four components.

where $N(\cdot; \mu_k, \Sigma_k)$ denotes the probability density function for the normal distribution $\mathcal{N}(\mu_k, \Sigma_k)$. See Figure 22.1 for the plot of such a density curve. Note that a mixture of Gaussians with just one mixture component reduces to a simple normal distribution, and so a GMMHMM with just one mixture component is simply a Gaussian HMM.

Similar to discrete HMMs, GMMHMMs seek to model a hidden state sequence $\{X_1, \dots, X_T\}$ and a corresponding observation sequence $\{Z_1, \dots, Z_T\}$ where T is the number of time steps or number of observations. The major difference is that each observation \mathbf{z}_t is a real-valued vector of length M distributed according to a mixture of Gaussians with K components. The parameters for such a model include the initial state distribution π and the state transition matrix A (just as with discrete HMMs). Additionally, for each state $i = 1, \dots, N$, we have component weights $\{c_{i,1}, \dots, c_{i,K}\}$, component means $\{\mu_{i,1}, \dots, \mu_{i,K}\}$, and component covariance matrices $\{\Sigma_{i,1}, \dots, \Sigma_{i,K}\}$.

Let's define a full GMMHMM with $N = 3$ states, components of dimension $M = 2$, and $K = 5$ components.

```
>>> import numpy as np

# 3x3 transition matrix
>>> A = np.array([[.6, .3, .1], [.2, .3, .5], [.7, .1, .2]])
```

```

# 3x5 collection of component weights
>>> weights = np.array([[.5, .1, .25, .09, 0.6], [0, .4, .3, .2, .1], [.1, .3, ←
    .2, .1, .3]])

# 3x5x2 collection of component means
>>> means = np.array([np.floor(np.random.uniform(-20, 20, size = (5, 2))) for i←
    in range(3)])

# 3x5x(2x2) collection of component covariance matrices
>>> covars = np.array([[np.floor(np.random.uniform(1, 20))*np.eye(2) for i in ←
    range(5)] for j in range(3)])

# (3,) ndarray initial state distribution
>>> pi = np.array([.4, .1, .5])

# Save the model parameters
>>> gmm = [A, weights, means, covars, pi]

```

Once we have a GMMHMM, we can randomly choose the first state based on the initial state distribution π . Now we can iteratively sample from our GMMHMM as follows:

- Randomly select a GMM Gaussian component according to the probability weights of the current state.
- Sample from the selected GMM Gaussian component using the corresponding mean and covariance matrix.
- Obtain the next state using the transition matrix A .

```

# choose initial state
>>> state = np.argmax(np.random.multinomial(1, pi))

# steps to randomly sample from GMMHMM
# randomly select a component using the probability weights of the current ←
# state
>>> sample_component = np.argmax(np.random.multinomial(1, weights[state,:]))

# sample an observation from the selected GMM Gaussian component
>>> sample = np.random.multivariate_normal(means[state, sample_component, :], ←
    covars[state, sample_component, :, :])

# obtain the next state using the transition matrix
>>> state = np.argmax(np.random.multinomial(1, A[:, state]))

```

Figure 22.2 shows an observation sequence generated from a GMMHMM with two mixture component and two states.

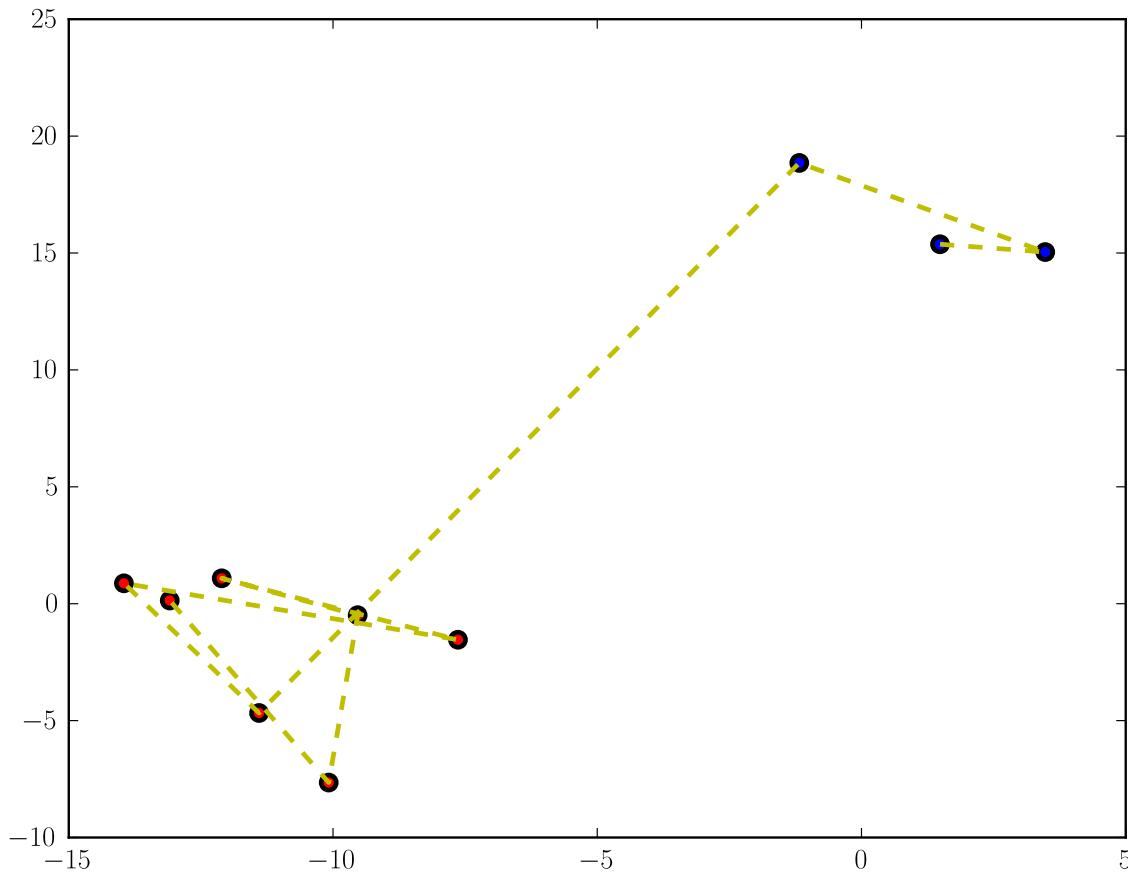


Figure 22.2: An observation sequence generated from a GMMHMM with two mixture components and two states. The observations (points in the plane) are shown as solid dots, the color indicating from which state they were generated. The connecting dotted lines indicate the sequential order of the observations.

Problem 1. Write a function which accepts a GMMHMM in the format above as well as an integer T , and which simulates the GMMHMM process, generating T different observations. Do so by implementing the following function declaration.

```
def sample_gmmhmm(gmmhmm, T):
    """
    Simulate sampling from a GMMHMM.

    Returns
    -----
    states : ndarray of shape (T,)
        The sequence of states
    obs : ndarray of shape (T, M)
        The generated observations
    
```

```
    """  
    pass
```

Test your function by running it on the gmmhmm given in the example, with $T = 900$. Use `sklearn.decomposition.PCA` with 2 components to plot the observations in two-dimensional space. Color the observations by state. How many clusters do you see?

The classic problems for which we normally use discrete observation HMMs can also be solved by using CDHMMs, though with continuous observations it is much more difficult to keep things numerically stable. We will not have you implement any of the three problems for CDHMMs yourself; instead, you will use a stable module we will provide for you. Note, however, that the techniques for solving these problems are still based on the forward-backward algorithm; the implementation may be trickier, but the mathematical ideas are virtually the same as those for discrete HMMs.

Speech Recognition and Hidden Markov Models

Hidden Markov Models are the basis of modern speech recognition systems. However, a fair amount of signal processing must precede the HMM stage, and there are other components of speech recognition, such as language models, that we will not address in this lab.

The basic signal processing and HMM stages of the speech recognition system that we develop in this lab can be summarized as follows: The audio to be processed is divided into small frames of approximately 30 ms. These are short enough that we can treat the signal as being constant over these intervals. We can then take this framed signal and, through a series of transformations, represent it by mel-frequency cepstral coefficients (MFCCs), keeping only the first M (say $M = 10$). Viewing these MFCCs as continuous observations in \mathbb{R}^M , we can train a GMMHMM on sequences of MFCCs for a given word, spoken multiple times. Doing this for several words, we have a collection of GMMHMMs, one for each word. Given a new speech signal, after framing and decomposing it into its MFCC array, we can score the signal against each GMMHMM, returning the word whose GMMHMM scored the highest.

Industrial-grade speech recognition systems do not train a GMMHMM for each word in a vocabulary (that would be ludicrous for a large vocabulary), but rather on phonemes, or distinct sounds. The English language has 44 phonemes, yielding 44 different GMMHMMs. As you could imagine, this greatly facilitates the problem of speech recognition. Each and every word can be represented by some combination of these 44 distinct sounds. By correctly classifying a signal by its phonemes, we can determine what word was spoken. Doing so is beyond the scope of this lab, so we will simply train GMMHMMs on five words/phrases: biology, mathematics, political science, psychology, and statistics.

Problem 2. Samples.zip contains 30 recordings for each of the words/phrases mathematics, biology, political science, psychology, and statistics. Remove the files that end in 00 (eg. Biology00.wav). These audio samples are 2 seconds in duration, recorded at a rate of 44100 samples per second, with samples stored as 16-bit signed integers in WAV format. Load the recordings into Python using `scipy.io.wavfile.read`.

Extract the MFCCs from each sample using code from the file MFCC.py. Store the MFCCs for each word in a separate list. You should have five lists, each containing 30 MFCC arrays, corresponding to each of the five words under consideration.

To load and extract, use the following code:

```
>>> samplerate, data = wavfile.read(file) # load wavfile
>>> model = MFCC.extract(data, show = False) # extract MFCC
```

For a specific word, given enough distinct samples of that word (decomposed into MFCCs), we can train a GMMHMM. Recall, however, that the training procedure does not always produce a very effective model, as it can get stuck in a poor local minimum. To combat this, we will train 10 GMMHMMs for each word (using a different random initialization of the parameters each time) and keep the model with the highest log-likelihood.

For training, we will use the python package called `hmmlearn`, as this is a stable implementation of GMMHMM algorithms. To facilitate random restarts, we need a function to provide initializations for the initial state distribution and the transition matrix.

Let `data` be a list of arrays, where each array is the output of the MFCC extraction for a speech sample. Using a function `initialize()` that returns a random initial state distribution and row-stochastic transition matrix, we can train a GMMHMM with 5 states and 3 mixture components and view its score as follows:

```
>>> import numpy as np # Import packages
>>> from hmmlearn import hmm
>>> startprob, transmat = initialize(5) # Get probabilities and transition ←
    matrices
>>> model = hmm.GMMHMM(n_components=5, covariance_type="diag", init_params = "←
    mc") # Initialize model
>>> model.startprob_ = startprob # Set probabilities and transition matrices
>>> model.transmat_ = transmat
>>> data = train_samples[word] # Reshape data for hmmlearns fit method
>>> lengths = [data[0].shape[0]] * len(data)
>>> data_collected = np.vstack(data)
>>> model.fit(data_collected) # Fit the model
>>> model.monitor_.history[-1] # Check the score
```

Achtung!

The process for problem 3 could take up to a couple of hours. Since you will not want to run this more than once, you may want to save the best model for each word to disk using the pickle module so you can use it later.

```
>>> import pickle
>>> temp = {word: best_model}
>>> with open(word, "wb") as out:
...     pickle.dump(temp, out)
```

Problem 3. Partition each list of MFCCs into a training set of 20 samples, and a test set of the remaining 10 samples. Using the training sets, train a GMMHMM on each of the words from the previous problem with at least 10 random restarts (reinitializing and creating a new model). Use $n_components = 5$ and $initialize(5)$. Keep the best model for each word (the one with the highest log-likelihood).

Given a trained model, we would like to compute the score of a new sample. Letting `obs` be an array of MFCCs for a speech sample we do this as follows:

```
>>> score = model.score(obs)
```

We classify a new speech sample by scoring it against each of the 5 trained GMMHMMs, and returning the word corresponding to the GMMHMM with the highest score.

Problem 4. Classify the 10 test samples for each word.

How does your system perform? Which words are the hardest to correctly classify? Make a dictionary containing the accuracy of the classification of your five testing sets. Specifically, the words/phrases will be the keys, and the values will be the percent accuracy. For example, to find the accuracy for the biology model score (`model.score(sample)`) each model on all 10 samples in the biology test set. The predicted class for each sample is the class of the model with the highest score. The accuracy of the biology model is the number of words in the biology test set that the biology model gave the highest score for over ten, since there were 10 words in the test set.

23 Kalman Filter

Lab Objective: Understand how to implement the standard Kalman Filter. Apply to the problem of projectile tracking.

Measured observations are often prone to significant noise, due to restrictions on measurement accuracy. For example, most commercial GPS devices can provide a good estimate of geolocation, but only within a dozen meters or so. A Kalman filter is an algorithm that takes a sequence of noisy observations made over time and attempts to get rid of the noise, producing more accurate estimates than the original observations. To do this, the algorithm needs information about the system being observed.

Consider the problem of tracking a projectile as it travels through the air. Short-range projectiles approximately trace out parabolas, but a sensor that is recording measurements of the projectile's position over time will likely show a path that is much less smooth. Because we know something about the laws of physics, we can filter out the noise in the measurements using basic Newtonian mechanics, recovering a more accurate estimate of the projectile's trajectory. In this lab, we will simulate measurements of a projectile and implement a Kalman filter to estimate the complete trajectory of the projectile.

Linear Dynamical Systems

The standard Kalman filter assumes that: (1) we have a linear dynamical system, (2) the state of the system evolves over time with some noise, and (3) we receive noisy measurements about the state of the system at each iteration. More formally, letting \mathbf{x}_k denote the state of the system at time k , we have

$$\mathbf{x}_{k+1} = F_k \mathbf{x}_k + G_k \mathbf{u}_k + \mathbf{w}_k \quad (23.1)$$

where F_k is a state-transition model, G_k is a control-input model, \mathbf{u}_k is a control vector, and \mathbf{w}_k is the noise present in state k . This noise is assumed to be drawn from a multivariate Gaussian distribution with zero mean and covariance matrix Q_k . The control-input model and control vector allow the assumption that the state can be additionally influenced by some other factor than the linear state-transition model.

We further assume that the states are “hidden,” and we only get the noisy observations

$$\mathbf{z}_k = H_k \mathbf{x}_k + \mathbf{v}_k \quad (23.2)$$

where H_k is the observation model mapping the state space to the observation space, and \mathbf{v}_k is the observation noise present at iteration k . As with the aforementioned error, we assume that this noise is drawn from a multivariate Gaussian distribution with zero mean and covariance matrix R_k .

The dynamics stated above are all taken to be linear. Thus, for our purposes, the operators F_k , G_k , and H_k are all matrices, and \mathbf{x}_k , \mathbf{u}_k , \mathbf{z}_k , and \mathbf{v}_k are all vectors.

We will assume that the transition and observation models, the control vector, and the noise covariances are constant, i.e. for each k , we will replace F_k , H_k , \mathbf{u}_k , Q_k , and R_k with F , H , \mathbf{u} , Q , and R . We will also assume that $G = I$ is the identity matrix, so it can safely be ignored.

Problem 1. Begin implementing a `KalmanFilter` class by writing an initialization method that stores the transition and observation models, noise covariances, and control vector. We provide an interface below:

```
class KalmanFilter(object):
    def __init__(self,F,Q,H,R,u):
        """
        Initialize the dynamical system models.

        Parameters
        -----
        F : ndarray of shape (n,n)
            The state transition model.
        Q : ndarray of shape (n,n)
            The covariance matrix for the state noise.
        H : ndarray of shape (m,n)
            The observation model.
        R : ndarray of shape (m,m)
            The covariance matrix for observation noise.
        u : ndarray of shape (n,)
            The control vector.
        """
        pass
```

We now derive the linear dynamical system parameters for a projectile traveling through \mathbb{R}^2 undergoing a constant downward gravitational force of 9.8 m/s^2 . The relevant information needed to describe how the projectile moves through space is its position and velocity. Thus, our state vector has the form

$$\mathbf{x} = \begin{pmatrix} s_x \\ s_y \\ V_x \\ V_y \end{pmatrix},$$

where s_x and s_y give the x and y coordinates of the position (in meters), and V_x and V_y give the horizontal and vertical components of the velocity (in meters per second), respectively.

How does the system evolve from one time step to the next? Assuming each time step is 0.1 seconds, it is easy enough to calculate the new position:

$$\begin{aligned}s'_x &= s_x + 0.1V_x \\ s'_y &= s_y + 0.1V_y.\end{aligned}$$

Further, since the only force acting on the projectile is gravity (we are ignoring things like wind resistance), the horizontal velocity remains constant:

$$V'_x = V_x.$$

The vertical velocity, however, does change due to the effects of gravity. From basic Newtonian mechanics, we have

$$V'_y = V_y - 0.1 \cdot 9.8.$$

In summary, over one time step, the state evolves from \mathbf{x} to \mathbf{x}' , where

$$\mathbf{x}' = \begin{pmatrix} s_x + 0.1V_x \\ s_y + 0.1V_y \\ V_x \\ V_y - 0.98 \end{pmatrix}.$$

From this equation, you can extract the state transition model F and the control vector u .

We now turn our attention to the observation model. Imagine that a radar sensor captures (noisy) measurements of the projectile's position as it travels through the air. At each time step, the radar transmits the observation $z = (z_x, z_y)$ given by

$$\begin{aligned}z_x &= s_x + v_x \\ z_y &= s_y + v_y,\end{aligned}$$

where (v_x, v_y) is a noise vector assumed to be drawn from a multivariate Gaussian with mean zero and some known covariance. These equations indicate the appropriate choice of observation model.

Problem 2. Work out the transition and observation models F and H , along with the control vector \mathbf{u} , corresponding to the projectile. Assume that the noise covariances are given by

$$\begin{aligned}Q &= 0.1 \cdot I_4 \\ R &= 5000 \cdot I_2.\end{aligned}$$

Instantiate a `KalmanFilter` object with these values.

We now wish to simulate a sequence of states and observations from the dynamical system. In addition to the system parameters, we need an initial state \mathbf{x}_0 to get started. Computing the subsequent states and observations is simply a matter of following equations 23.1 and 23.2.

Problem 3. Add a method to your `KalmanFilter` class to generate a state and observation sequence by evolving the system from a given initial state (the function `numpy.random.multivariate_normal` will be useful). To do this, implement the following:

```
def evolve(self, x0, N):
    """
```

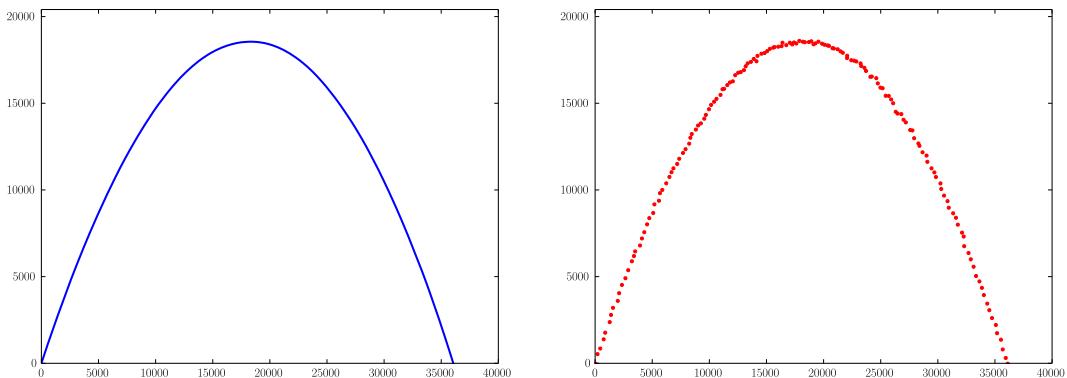


Figure 23.1: State sequence (left) and sampling of observation sequence (right).

Compute the first N states and observations generated by the Kalman system.

Parameters

x0 : ndarray of shape (n,)
The initial state.

N : integer

The number of time steps to evolve.

Returns

states : ndarray of shape (n,N)
States 0 through N-1, given by each column.

obs : ndarray of shape (m,N)
Observations 0 through N-1, given by each column.

....

pass

Simulate the true and observed trajectory of a projectile with initial state

$$\mathbf{x}_0 = \begin{pmatrix} 0 \\ 0 \\ 300 \\ 600 \end{pmatrix}.$$

Approximately 1250 time steps should be sufficient for the projectile to hit the ground (i.e. for the y coordinate to return to 0). Your results should qualitatively match those given in Figure 23.1.

State Estimation with the Kalman Filter

The Kalman filter is a recursive estimator that smooths out the noise in real time, estimating each current state based on the past state estimate and the current measurement. This process is done by repeatedly invoking two steps: Predict and Update. The predict step is used to estimate the current state based on the previous state. The update step then combines this prediction with the current observation, yielding a more robust estimate of the current state.

To describe these steps in detail, we need additional notation. Let

- $\hat{\mathbf{x}}_{n|m}$ be the state estimate at time n given only measurements up through time m ; and
- $P_{n|m}$ be an error covariance matrix, measuring the estimated accuracy of the state at time n given only measurements up through time m .

The elements $\hat{\mathbf{x}}_{k|k}$ and $P_{k|k}$ represent the state of the filter at time k , giving the state estimate and the accuracy of the estimate.

We evolve the filter recursively, as follows:

Predict	$\hat{\mathbf{x}}_{k k-1} = F\hat{\mathbf{x}}_{k-1 k-1} + \mathbf{u}$
	$P_{k k-1} = FP_{k-1 k-1}F^T + Q$
Update	$\tilde{\mathbf{y}}_k = \mathbf{z}_k - H\hat{\mathbf{x}}_{k k-1}$
	$S_k = HP_{k k-1}H^T + R$
	$K_k = P_{k k-1}H^T S_k^{-1}$
	$\hat{\mathbf{x}}_{k k} = \hat{\mathbf{x}}_{k k-1} + K_k\tilde{\mathbf{y}}_k$
	$P_{k k} = (I - K_kH)P_{k k-1}$

The more observations we have, the greater the accuracy of these estimates becomes (i.e the norm of the accuracy matrix converges to 0).

Problem 4. Add code to your `KalmanFilter` class to estimate a state sequence corresponding to a given observation sequence and initial state estimate. Implement the following class method:

```
def estimate(self,x,P,z):
    """
    Compute the state estimates using the Kalman filter.
    If x and P correspond to time step k, then z is a sequence of
    observations starting at time step k+1.

    Parameters
    -----
    x : ndarray of shape (n,)
        The initial state estimate.
    P : ndarray of shape (n,n)
        The initial error covariance matrix.
    z : ndarray of shape(m,N)
        Sequence of N observations (each column is an observation).
    """

    # Your implementation here
    pass
```

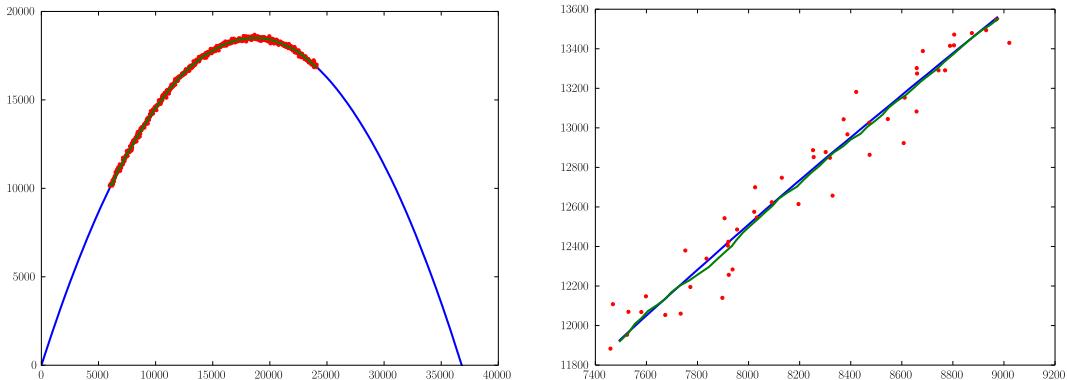


Figure 23.2: State estimates together with observations and true state sequence (detailed view on the right).

```

Returns
-----
out : ndarray of shape (n,N)
      Sequence of state estimates (each column is an estimate).
"""
pass

```

Returning to the projectile example, we now assume that our radar sensor has taken observations from time steps 200 through 800 (take the corresponding slice of the observations produced in Problem 3). Using these observations, we seek to estimate the corresponding true states of the projectile. We must first come up with a state estimate $\hat{\mathbf{x}}_{200}$ for time step 200, and then feed this into the Kalman filter to obtain estimates $\hat{\mathbf{x}}_{201}, \dots, \hat{\mathbf{x}}_{800}$.

Problem 5. Calculate an initial state estimate $\hat{\mathbf{x}}_{200}$ as follows: For the horizontal and vertical positions, simply use the observed position at time 200. For the velocity, compute the average velocity between the observations \mathbf{z}_k and \mathbf{z}_{k+1} for $k = 200, \dots, 208$, then average these 9 values and take this as the initial velocity estimate. (Hint: the NumPy function `diff` is useful here.)

Using the initial state estimate, $P_{200} = 10^6 \cdot Q$, and your Kalman filter, compute the next 600 state estimates, i.e. compute $\hat{\mathbf{x}}_{201}, \dots, \hat{\mathbf{x}}_{800}$. Plot these state estimates as a smooth green curve together with the radar observations (as red dots) and the entire true state sequence (as a blue curve). Zoom in to see how well it follows the true path. Your plots should be similar to Figure 23.2.

In the absence of observations, we can still estimate some information about the state of the system at some future time. We can do this by recognizing that the expected state noise $\mathbb{E}[\boldsymbol{\varepsilon}_k] = 0$ at any time k . Thus, given a current state estimate $\hat{\mathbf{x}}_{n|m}$ using only measurements up through time m , the expected state at time $n + 1$ is

$$\hat{\mathbf{x}}_{n+1|m} = F\hat{\mathbf{x}}_{n|m} + \mathbf{u}$$

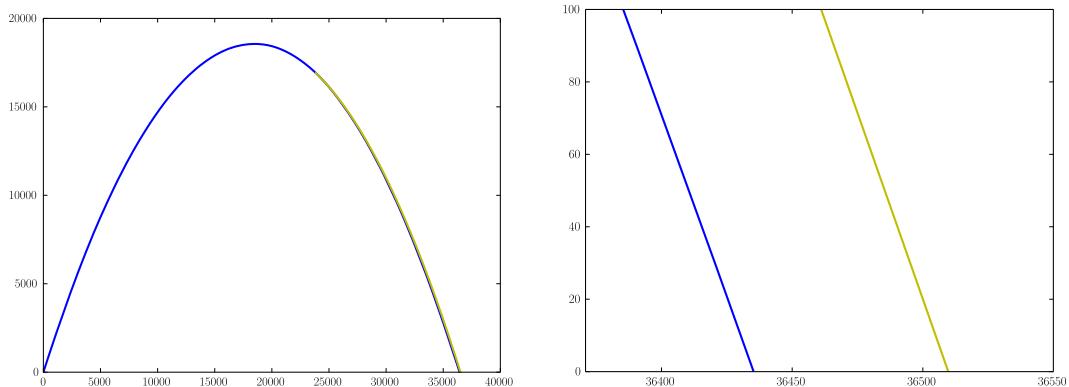


Figure 23.3: Predicted vs. actual point of impact (detailed view on right).

Problem 6. Add a function to your class that predicts the next k states given a current state estimate but in the absence of observations. Do so by implementing the following function:

```
def predict(self,x,k):
    """
    Predict the next k states in the absence of observations.

    Parameters
    -----
    x : ndarray of shape (n,)
        The current state estimate.
    k : integer
        The number of states to predict.

    Returns
    -----
    out : ndarray of shape (n,k)
        The next k predicted states.
    """
    pass
```

We can use this prediction routine to estimate where the projectile will hit the surface.

Problem 7. Using the final state estimate \hat{x}_{800} that you obtained in Problem 5, predict the future states of the projectile until it hits the ground. Predicting approximately the next 450 states should be sufficient.

Plot the actual state sequence together with the predicted state sequence (as a yellow curve), and observe how near the prediction is to the actual point of impact. Your results should be similar to those shown in Figure 23.3.

In the absence of observations, we can also reverse the system and iterate backward in time to infer information about states of the system prior to measured observations. The system is reversed by

$$\mathbf{x}_k = F^{-1}(\mathbf{x}_{k+1} - \mathbf{u} - \boldsymbol{\varepsilon}_{k+1}).$$

Considering again that $\mathbb{E}[\boldsymbol{\varepsilon}_k] = 0$ at any time k , we can ignore this term, simplifying the recursive estimation backward in time.

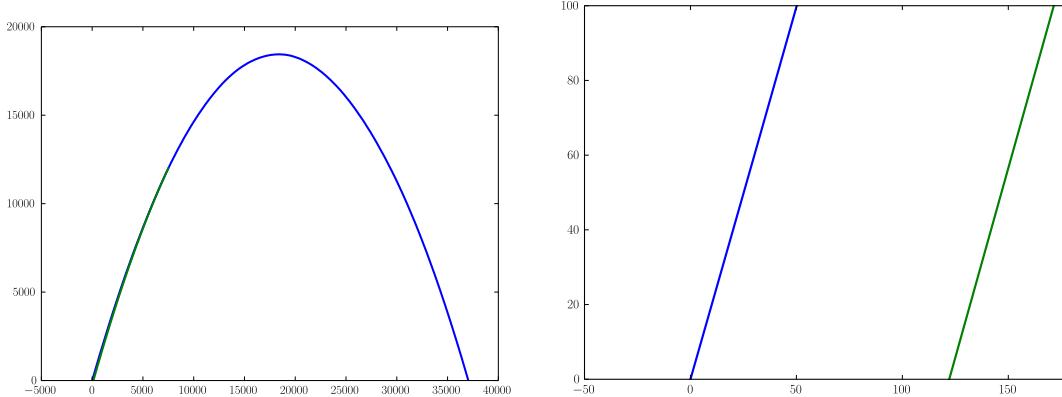


Figure 23.4: Predicted vs. actual point of origin (detailed view on right).

Problem 8. Add a function to your class that rewinds the system from a given state estimate, returning predictions for the previous states. Do so by implementing the following function:

```
def rewind(self, x, k):
    """
    Predict the k states preceding the current state estimate x.

    Parameters
    -----
    x : ndarray of shape (n,)
        The current state estimate.
    k : integer
        The number of preceding states to predict.

    Returns
    -----
    out : ndarray of shape (n,k)
        The k preceding predicted states.
    """
    pass
```

Returning to the projectile example, we can now predict the point of origin.

Problem 9. Using your state estimate $\hat{\mathbf{x}}_{250}$, predict the point of origin of the projectile along with all states leading up to time step 250. Note that you may have to take a few extra time steps to predict the point of origin. (The point of origin is the first point along the trajectory where the y coordinate is 0.) Plot these predicted states (in cyan) together with the original state sequence. Zoom in to see how accurate your prediction is. Your plots should be similar to Figure 23.4.

Repeat the prediction starting with $\hat{\mathbf{x}}_{600}$. Compare to the previous results. Which is better? Why?

24 ARMA Models

Lab Objective: ARMA(p, q) models combine autoregressive and moving-average models in order to forecast future observations using time-series. In this lab, we will build an ARMA(p, q) model to analyze and predict future weather data and then compare this model to statsmodels built-in ARMA package as well as the VARMAX package. Then we will forecast macroeconomic data as well as the future height of the Rio Negro.

Time Series

A time series is any discrete-time stochastic process. In other words, it is a sequence of random variables, $\{Z_t\}_{t=1}^T$, that are determined by their time t . We let the realization of the time series $\{Z_t\}_{t=1}^T$ be denoted by $\{z_t\}_{t=1}^T$. Examples of time series include heart rate readings over time, pollution readings over time, stock prices at the closing of each day, and air temperature. Often when analyzing time series, we want to forecast future data, such as what will the stock price of a company will be in a week and what will the temperature be in 10 days.

ARMA(p, q) Models

One way to forecast a time series is using an ARMA model. The Wold Theorem says that any covariance-stationary time series can be well approximated with an ARMA model. An ARMA(p, q) model combines an autoregressive model of order p and a moving average model of order q on a time series $\{Z_t\}_{t=1}^T$. The model itself is a discrete-time stochastic process $(Z_t)_{t \in \mathbb{Z}}$ satisfying the equation

$$Z_t = \mathbf{c} + \underbrace{\left(\sum_{i=1}^p \Phi_i Z_{t-i} \right)}_{\text{AR}(p)} + \underbrace{\left(\sum_{j=1}^q \Theta_j \varepsilon_{t-j} \right)}_{\text{MA}(q)} + \varepsilon_t \quad (24.1)$$

where each ε_t is an identically-distributed Gaussian variable with mean 0 and constant covariance Σ , $\mathbf{c} \in \mathbb{R}^n$, and Φ_i and Θ_j are in $M_n(\mathbb{R})$.

AR(p) Models

An AR(p) model works similar to a weighted random walk. Recall that in a random walk, the current position depends on the immediate past position. In the autoregressive model, the current data point in the time series depends on the past p data points. However, the importance of each of the past p data points is not uniform. With an error term to represent white noise and a constant term to adjust the model along the y-axis, we can model the stochastic process with the following equation:

$$Z_t = \mathbf{c} + \sum_{i=1}^p \Phi_i Z_{t-i} + \varepsilon_t \quad (24.2)$$

If there is a high correlation between the current and previous values of the time series, then the AR(p) model is a good representation of the data, and thus the ARMA(p, q) model will most likely be a good representation. The coefficients $\{\Phi_i\}_{i=1}^p$ are larger when the correlation is stronger.

In this lab, we will be using weather data from Provo, Utah¹. To check that the data can be represented well, we need to look at the correlation between the current and previous values.

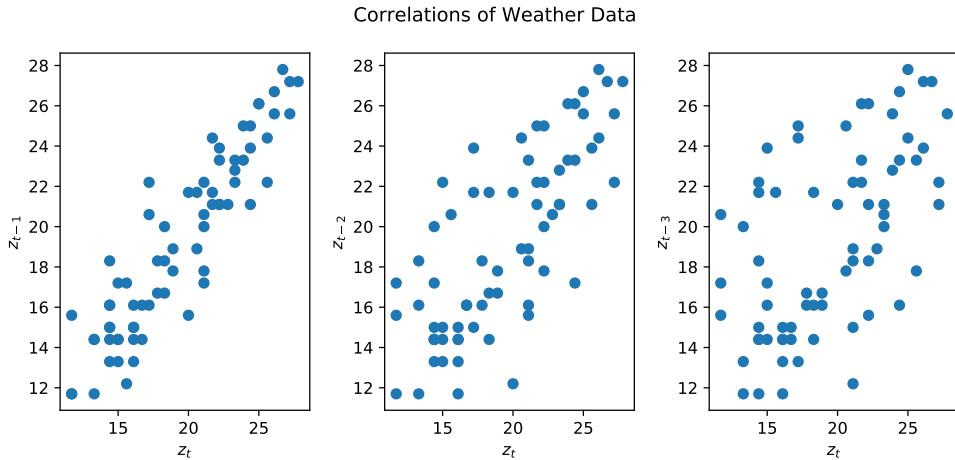


Figure 24.1: These graphs show that the weather data is correlated to its previous values. The correlation is weaker in each graph successively, showing that the further in the past the data is, the less correlated the data becomes.

MA(q) Models

A moving average model of order q is used to factor in the varying error of the time series. This model uses the error of the current data point and the previous data points to predict the next datapoint. Similar to an AR(p) model, this model uses a linear combination (which includes a constant term to adjust along the y-axis..

$$Z_t = \mathbf{c} + \varepsilon_t + \sum_{i=1}^q \Theta_i \varepsilon_{t-i} \quad (24.3)$$

This part of the model simulates shock effects in the time series. Examples of shock effects include volatility in the stock market or sudden cold fronts in the temperature.

¹This data was taken from <https://forecast.weather.gov/data/obhistory/metric/KPVU.html>

Combining both the AR(p) and MA(q) models, we get an ARMA(p, q) model which forecasts based on previous observations and error trends in the data.

ARIMA(p, d, q) Models

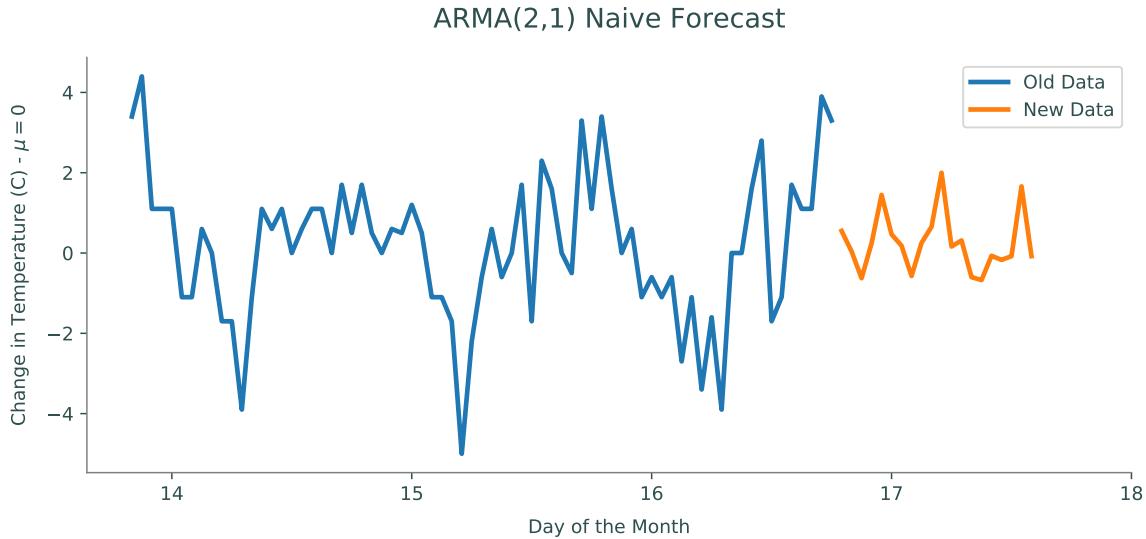
Not all ARMA models are covariance stationary. However, many time series can be made covariance stationary by differencing. Let δZ_t represent the time series $Y_t = Z_t - Z_{t-1}$ obtained by taking a difference of the terms. If the trend is linear a first difference is usually stationary. If the trend is quadratic a second difference may be necessary $\delta^2 Z_t = \delta(\delta Z_t)$. An ARIMA(p, d, q) model is a discrete-time stochastic process $(Z_t)_{t \in \mathbb{Z}}$ satisfying the equation

$$\delta^d Z_t = \mathbf{c} + \underbrace{\left(\sum_{i=1}^p \Phi_i \mathbf{y}_{t-i} \right)}_{\text{AR}(p)} + \underbrace{\left(\sum_{j=1}^q \Theta_j \varepsilon_{t-j} \right)}_{\text{MA}(q)} + \varepsilon_t \quad (24.4)$$

Finding Parameters

One of the most difficult parts of using an ARMA(p, q) model is identifying the proper parameters of the model. For simplicity, at the beginning of this lab we discuss univariate ARMA models with parameters $\{\phi_i\}_{i=1}^p$, $\{\theta_i\}_{i=1}^q$, μ , and σ , where μ and σ are the mean and variance of the error. Note that $\{\phi_i\}_{i=1}^p$ and $\{\theta_i\}_{i=1}^q$ determine the order of the ARMA model.

A naive way to use an ARMA model is to choose p and q based on intuition. Figure 24.1 showed that there is a strong correlation between z_t and z_{t-1} and between z_t and z_{t-2} . The correlation is weaker between z_t and z_{t-3} . Intuition then suggests to choose $p = 2$. By looking at the correlations between the current noise with previous noise, similar to Figure 24.1, it can also be seen that there is a weak correlation between z_t and ε_t and between z_t and ε_{t-1} . Between z_t and ε_{t-2} there is no correlation. For more on how these error correlations were found, see Additional Materials. Intuition from these correlations suggests to choose $q = 1$. Thus, a naive choice for our model is an ARMA(2, 1) model.

Figure 24.2: Naive forecast on `weather.npy`

Problem 1. Write a function `arma_forecast_naive()` that builds an ARMA(p, q) model where the values of $\phi_i = .5$ and $\theta_i = .1$ for all i . Let $\varepsilon_i \sim \mathcal{N}(0, 1)$ for all i .

Use your function to predict the next n values of the time series. The time series that should be used is the first difference of the time series found in the file `weather.npy`, which we denote $\{z_t\}_{t=1}^T$. This is done because we want the time series to be covariance stationary. The function should accept as parameters p , q , and n , where p is the order of the autoregressive model, q is the order of the moving average model, and n is the number of observations to predict. Plot the observed differences $\{z_t\}_{t=1}^T$ followed by your predicted observations of z_t .

Hint: you might find `np.diff()` to be useful.

The file `weather.npy` contains data on the temperature in Provo, Utah from 7:56 PM May 13, 2019 to 6:56 PM May 16, 2019, taken every hour.

Use this file to test your code. For $p = 2$, $q = 1$, and $n = 20$, your plot should look similar to Figure 24.2, however, due to the variance of the error ε_t , the plot will not look exactly like Figure 24.2. The predictions may be higher or lower on the x-axis.

Let $\Theta = \{\phi_i, \theta_j, \mu, \sigma_a^2\}$ be the set of parameters for an ARMA(p, q) model. Suppose we have a set of observations $\{z_t\}_{t=1}^n$. Our goal is to find the p, q , and Θ that maximize the likelihood of the ARMA model given the data. Using the chain rule, we can factorize the likelihood of the model given this data as

$$p(\{z_t\}|\Theta) = \prod_{t=1}^n p(z_t|z_{t-1}, \dots, z_1, \Theta) \quad (24.5)$$

State Space Representation

In a general ARMA(p, q) model, the likelihood is a function of the unobserved error terms ε_t and is not trivial to compute. Simple approximations can be made, but these may be inaccurate under certain circumstances. Explicit derivations of the likelihood are possible, but tedious. However, when the ARMA model is placed in state-space, the Kalman filter affords a straightforward, recursive way to compute the likelihood.

We demonstrate one possible state-space representation of an ARMA(p, q) model. Let $r = \max(p, q + 1)$. Define

$$\hat{\mathbf{x}}_{t|t-1} = [x_{t-1} \ x_{t-2} \ \cdots \ x_{t-r}]^T \quad (24.6)$$

$$F = \begin{bmatrix} \phi_1 & \phi_2 & \cdots & \phi_{r-1} & \phi_r \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix} \quad (24.7)$$

$$H = [1 \ \theta_1 \ \theta_2 \ \cdots \ \theta_{r-1}] \quad (24.8)$$

$$Q = \begin{bmatrix} \sigma_a^2 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix} \quad (24.9)$$

$$w_t \sim \text{MVN}(0, Q), \quad (24.10)$$

where $\phi_i = 0$ for $i > p$, and $\theta_j = 0$ for $j > q$. Note that Equation 24.2 gives

$$F\hat{\mathbf{x}}_{t-1|t-2} + w_t = \begin{bmatrix} \sum_{i=1}^r \phi_i x_{t-i} \\ x_{t-1} \\ x_{t-2} \\ \vdots \\ x_{t-(r-1)} \end{bmatrix} + \begin{bmatrix} \varepsilon_t \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (24.11)$$

$$= [x_t \ x_{t-1} \ \cdots \ x_{t-(r-1)}]^T \quad (24.12)$$

$$= \hat{\mathbf{x}}_{t|t-1} \quad (24.13)$$

We note that $z_{t|t-1} = H\hat{\mathbf{x}}_{t|t-1} + \mu$.²

Then the linear stochastic dynamical system

$$\hat{\mathbf{x}}_{t+1|t} = F\hat{\mathbf{x}}_{t|t-1} + w_t \quad (24.14)$$

$$z_{t|t-1} = H\hat{\mathbf{x}}_{t|t-1} + \mu \quad (24.15)$$

describes the same process as the original ARMA model.

Note

²For a proof of this fact, see Additional Materials.

Equation 24.15 involves a deterministic component, namely μ . The Kalman filter theory developed in the previous lab, however, assumed $\mathbb{E}[\varepsilon_t] = 0$ for the observations $z_{t|t-1}$. This means you should subtract off the mean μ of the error from the time series observations $z_{t|t-1}$ when using them in the predict and update steps.

Likelihood via Kalman Filter

We assumed in Equation 24.10 that the error terms of the model are Gaussian. This means that each conditional distribution in 24.5 is also Gaussian, and is completely characterized by its mean and variance. These two quantities are easily found via the Kalman filter:

$$\text{mean } H\hat{x}_{t|t-1} + \mu \quad (24.16)$$

$$\text{variance } HP_{t|t-1}H^T \quad (24.17)$$

where $\hat{x}_{t|t-1}$ and $P_{t|t-1}$ are found during the Predict step. Given that each conditional distribution is Gaussian, the likelihood can then be found as follows:

$$p(\{z_t\}|\Theta) = \prod_{t=1}^n N(z_t | H\hat{x}_{t|t-1} + \mu, HP_{t|t-1}H^T) \quad (24.18)$$

Problem 2. Write a function `arma_likelihood()` that returns the log-likelihood of an ARMA model, given a time series $\{z_t\}_{t=1}^T$. This function should accept `filename` which contains the observations, and it should accept as parameters each parameter in Θ . In this case, the time series should be the change in temperature of `weather.npy`, which is the first difference of the time series found in `weather.npy`, as was done in the first problem. Return the log-likelihood of the ARMA(p, q) model as a `float`.

Use the `state_space_rep()` function provided to generate F, Q , and H . The function `kalman()` has also been provided to calculate the means and covariances of each observation.

Hint: calling the function `kalman()` on a time series will return an array whose values are $x_{k|k-1}$ and an array whose values are $P_{k|k-1}$ for each $k \leq n$. Remember that the time series should have μ subtracted when using `kalman()`.

When done correctly, your function should match the following output:

```
>>> arma_likelihood(filename='weather.npy', phis=np.array([0.9]),
                     thetas=np.array([0]), mu=17., std=0.4)
-1375.1805469978776
```

Model Identification

Now that we can compute the likelihood of a given ARMA model, we want to find the best choice of parameters given our time series. In this lab, we define the model with the "best" choice of parameters as the model which minimizes the AIC. The benefit of minimizing the AIC is that it rewards goodness of fit while penalizing overfitting. The AIC is expressed by

$$2k \left(1 + \frac{k+1}{n-k} \right) - 2\ell(\Theta) \quad (24.19)$$

where n is the sample size, $k = p + q + 2$ is the number of parameters in the model, and $\ell(\Theta)$ is the maximum likelihood for the model class.

To compute the maximum likelihood for a model class, we need to optimize 24.18 over the space of parameters Θ . We can do so by using an optimization routine such as `scipy.optimize.minimize` on the function `arma_likelihood()` from Problem 2. Use the following code to run this routine.

```
>>> from scipy.optimize import minimize

>>> # assume p, q, and time_series are defined
>>> def f(x): # x contains the phis, thetas, mu, and std
>>>     return -1*arma_likelihood(filename, phis=x[:p], thetas=x[p:p+q],
>>>                             mu=x[-2], std=x[-1])
>>> # create initial point
>>> x0 = np.zeros(p + q + 2)
>>> x0[-2] = time_series.mean()
>>> x0[-1] = time_series.std()
>>> sol = minimize(f, x0, method = "SLSQP")
>>> sol = sol['x']
```

This routine will return a vector `sol` where the first p values are $\{\phi_i\}_{i=1}^p$, the next q values are $\{\theta_i\}_{i=1}^q$, and the last two values are μ and σ , respectively. Note the wrapper $f(x)$ returns the negative log-likelihood. This is because `scipy.optimize.minimize` finds the minimizer of $f(x)$ and we are solving for the maximum likelihood.

To minimize the AIC, we perform model identification. This is choosing the order of our model, p and q , from some admissible set. The order of the model which minimizes the AIC is then the optimal model.

Problem 3. Write a function `model_identification()` that accepts `filename` containing the time series data and parameters p_{max} and q_{max} as integers. Return each parameter in Θ that minimizes the AIC of an ARMA(i, j) model, given that $1 \leq i \leq p_{max}$ and $1 \leq j \leq q_{max}$.

Your code should produce the following output (it may take awhile to run):

```
>>> model_identification(filename='weather.npy', p_max=4, q_max=4)
(array([ 0.7213538]), array([-0.26246426]), 0.359785001944352, ←
    1.5568374351425505)
```

Forecasting with Kalman Filter

We now have identified the optimal ARMA(p, q) model. We can use this model to predict future states. The Kalman filter provides a straightforward way to predict future states by giving the mean and variance of the conditional distribution of future observations. Observations can be found as follows

$$z_{t+k}|z_1, \dots, z_t \sim N(z_{t+k}; H\hat{x}_{t+k|t} + \mu, HP_{t+k|t}H^T) \quad (24.20)$$

To evolve the Kalman filter, recall the predict and update rules of a Kalman filter.

Predict	$\hat{\mathbf{x}}_{k k-1} = F\hat{\mathbf{x}}_{k-1 k-1} + \mathbf{u}$
	$P_{k k-1} = FP_{k-1 k-1}F^T + Q$
Update	$\tilde{\mathbf{y}}_k = \mathbf{z}_k - H\hat{\mathbf{x}}_{k k-1}$
	$S_k = HP_{k k-1}H^T + R$
	$K_k = P_{k k-1}H^T S_k^{-1}$
	$\hat{\mathbf{x}}_{k k} = \hat{\mathbf{x}}_{k k-1} + K_k\tilde{\mathbf{y}}_k$
	$P_{k k} = (I - K_kH)P_{k k-1}$

Achtung!

Recall that the values returned by `kalman()` are conditional on the previous observation. To compute the mean and variance of future observations, the values $x_{n|n}$ and $P_{n|n}$ MUST be computed using the update step. Once computed, only the predict step is needed to find the future means and covariances.

Problem 4. Write a function `arma_forecast()` that accepts `filename` containing a time series, the parameters for an ARMA model, and the number n of observations to forecast. Calculate the mean and covariance of the future n observations using a Kalman filter. Plot the original observations as well as the `mean` for each future observation. Plot a 95% confidence interval (2 standard deviations away from the mean) around the means of future observations. Return the means and standard deviations calculated. Hint: the standard deviation is the square root of the covariance calculated.

The following code should create a plot similar to Figure 24.3.

```
>>> # Get optimal model as found in the previous problem
>>> phis, thetas, mu, std = np.array([0.72135856]), np.array([
>>> [-0.26246788]), 0.35980339870105321, 1.5568331253098422

>>> # Forecast optimal mode
>>> arma_forecast(filename='weather.npy', phis=phis, thetas=thetas,
>>>                 mu=mu, std=std, n=30)
```

How does this plot compare to the naive ARMA model made in Problem 1?

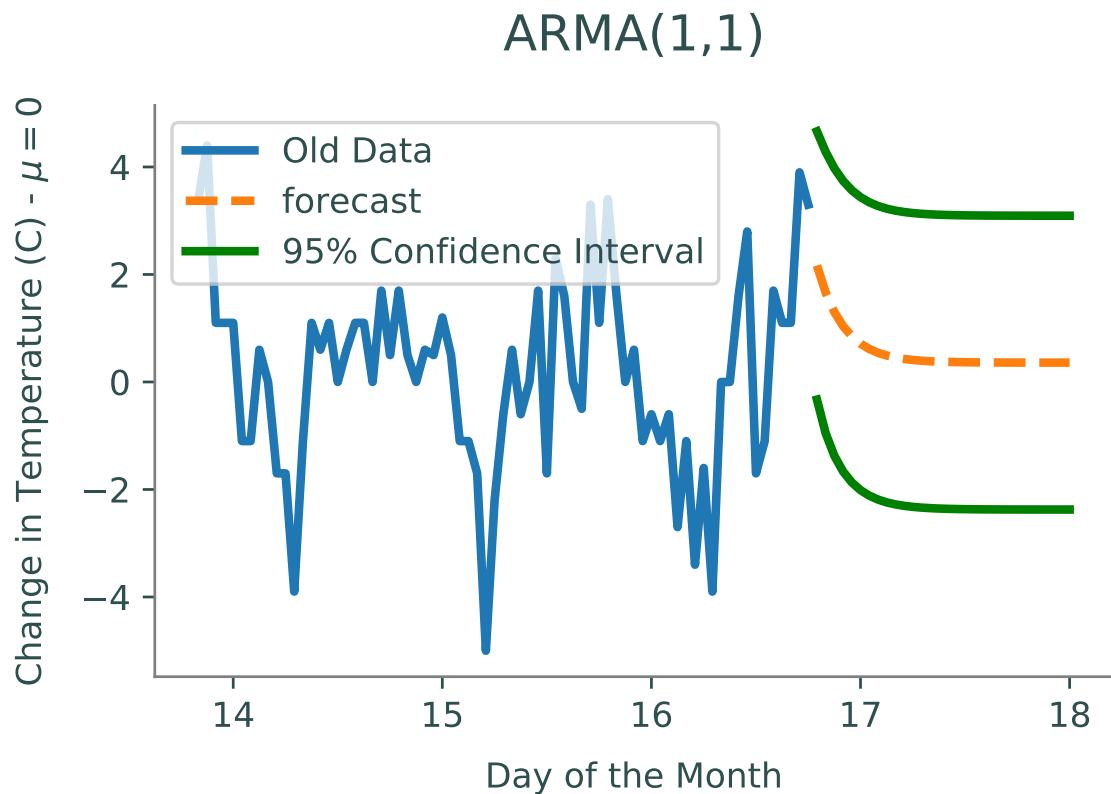


Figure 24.3: ARMA(1,1) forecast on `weather.npy`

Statsmodel ARMA

The module `statsmodels` contains a package that includes an ARMA model class. This is accessed through ARIMA model, which stands for Autoregressive Integrated Moving Average. This class also uses a Kalman Filter to calculate the MLE. When creating an ARIMA object, initialize the variables `endog` (the data) and `order` (the order of the model). The order is of the form (p, d, q) where d is the differences. To create an ARMA model, set $d = 0$. The object can then be fitted based on the MLE using a Kalman Filter.

```
from statsmodels.tsa.arima.model import ARIMA
# Initialize the object with weather data and order (1,1)
model = ARIMA(z,order=(p,0,q),trend='c').fit(method='innovations_mle')

# Access p and q
>>> model.specification.k_ar
p
>>> model.specification.k_ma
q
```

As in other problems, the data passed in should be the time series stationary. The AIC of an ARMA model object is saved as the attribute `aic`. Since the AIC is much faster to compute using `statsmodels`, model identification is much faster. Once a model is chosen, the method `predict` will forecast n observations, where n is the number of known observations. It will return the mean of each future observation.

```
# Predict from the beginning of the model to 30 observations in the future
model.predict(start=0,end=len(data)+30)
```

Problem 5. Write a function `sm_arma()` that accepts `filename` containing a time series, integer values for `p_max` and `q_max`, and the number n of values to predict. Use `statsmodels` to perform model identification as in Problem 3, where the order of $\text{ARMA}(i,j)$ satisfies $1 \leq i \leq p_{\text{max}}$ and $1 \leq j \leq q_{\text{max}}$. Ensure the model is fit using the MLE.

Use the optimal model to predict n future observations of the time series. Plot the original observations along with the mean of each future observations given by `statsmodels`. Return the AIC of the optimal model.

For $p_{\text{max}} = 3$, $q_{\text{max}} = 3$, and $n = 30$, your graph should look similar to Figure 24.4. How does this graph compare to Problem 1? Problem 4?

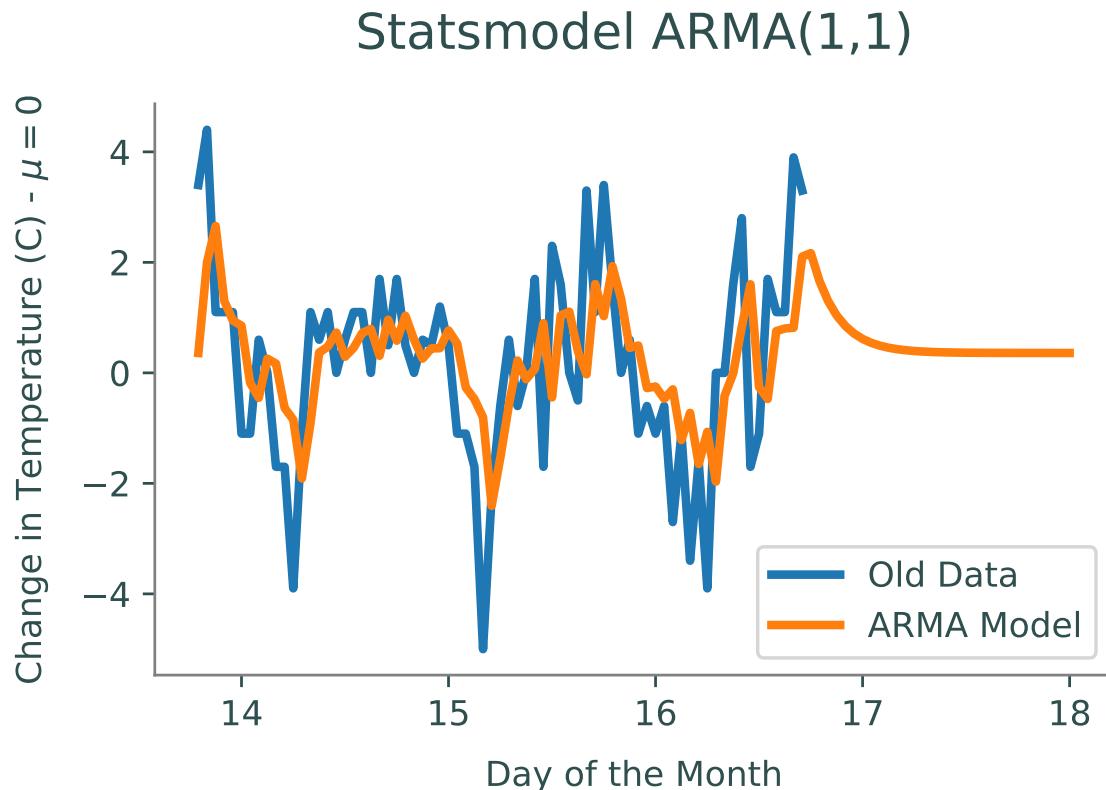


Figure 24.4: Statsmodel ARMA(3,3) forecast on `weather.npy`.

Statsmodel VARMA

Until now we have been dealing with univariate ARMA models. Multivariate ARMA models are used when we have multiple time series that can be useful in predicting one another. For example say we have two time series $z_{t,1}$ and $z_{t,2}$. The multivariate ARMA(1,1) model is as follows:

$$z_{t,1} = c_1 + \phi_{11}z_{t-1,1} + \phi_{12}z_{t-1,2} + \theta_{11}\varepsilon_{t-1,1} + \theta_{12}\varepsilon_{t-1,2} \quad (24.21)$$

$$z_{t,2} = c_1 + \phi_{21}z_{t-1,1} + \phi_{22}z_{t-1,2} + \theta_{21}\varepsilon_{t-1,1} + \theta_{22}\varepsilon_{t-1,2} \quad (24.22)$$

This can be written in matrix form as shown in equation 24.1. The module `statsmodels` contains a package that includes an VARMAX model class which can be used to create a multivariate ARMA model. This stands for Vector Autoregression Moving Average with Exogenous Regressors. An exogenous regressor is a time series that affects the model but is not affected by it. In the example below we have two time series corresponding to the price of copper and aluminum. Since aluminum is a substitute for copper, it is reasonable to assume the price of aluminum may help us predict the price of copper and vice versa. Note that when fitting a VARMAX model setting the parameter `ic` to `aic` selects parameters based on AIC criterion.

```
>>>from statsmodels.tsa.api import VARMAX
>>>import statsmodels.api as sm

>>> # Load in world copper data
>>> data = sm.datasets.copper.load_pandas().data
>>> # Create index compatible with VARMAX model
>>> idx = pd.period_range(start='1951', end='1975', freq = 'Y')
>>> data.index = idx

>>> # Initialize and fit model
>>> mod = VARMAX(data[['ALUMPRICE', 'COPPERPRICE']])
>>> mod = mod.fit(maxiter=1000, disp=False, ic = 'aic')
>>> # Predict until the price of aluminium and copper until 1985
>>> pred = mod.predict('1951','1985')

>>> # Get confidence intervals
>>> forecast_obj = mod.get_forecast('1981')
>>> all_CI = forecast_obj.conf_int(alpha=0.05)
>>> all_CI

>>> # Plot predictions against true price
>>> pred.plot()
>>> plt.plot(data['ALUMPRICE'], 'r--', label = 'ALUMPRICE prediction')
>>> plt.plot(data['COPPERPRICE'], 'r--', label = 'COPPERPRICE prediction')
>>> plt.legend()
>>> plt.title('VARMA Predictions for World Copper Market Dataset')
```

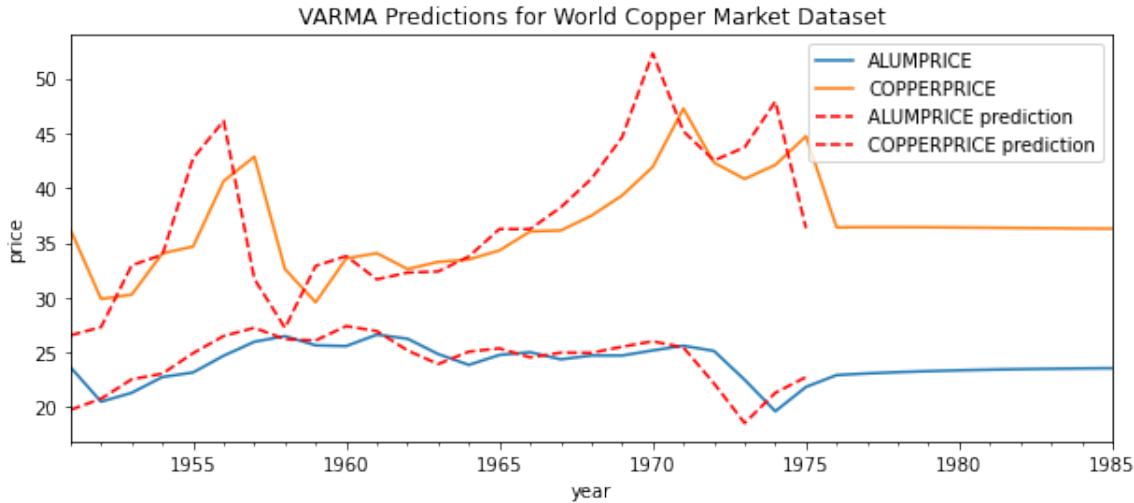


Figure 24.5: Statsmodel VAR(1) forecast.

Problem 6. Write a function `sm_varma()` that accepts start and end dates for forecasting. Use the statsmodels VARMAX class to forecast on macroeconomic data between the start and end dates. Use AIC as the criterion for model selection when fitting the model. Plot the prediction, original data and a 95% confidence interval (2 standard deviations away from the mean) around the future observations. Return the AIC of the chosen model. The plot should be similar to Figure 24.6.

The following code shows how to obtain the data.

```
>>> # Load in data
>>> df = sm.datasets.macrodata.load_pandas().data
>>> # Create DatetimeIndex
>>> dates = df[['year', 'quarter']].astype(int).astype(str)
>>> dates = dates["year"] + "Q" + dates["quarter"]
>>> dates = dates_from_str(dates)
>>> df.index = pd.DatetimeIndex(dates)
>>> # Select columns used in prediction
>>> df = df[['realgdp', 'realcons', 'realinv']]
```

The dataset '`realgdp`' contains the real gross domestic product, '`realcons`' contains real personal consumption expenditures, and '`realinv`' contains real gross private domestic investment. Since personal consumption and domestic investment are components of gross domestic product, it is reasonable to assume these time series will be useful in predicting one another.

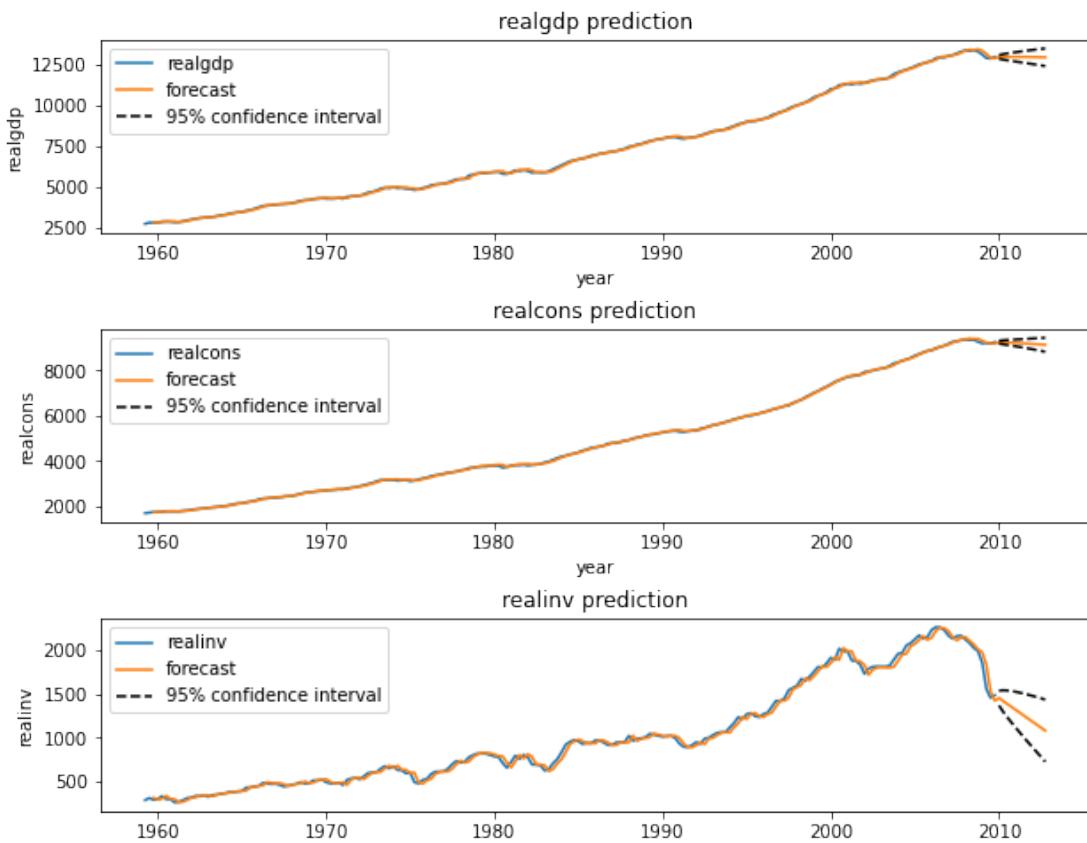


Figure 24.6: Macroeconomic data is forecasted 12 years in the future using statsmodels.

Optional

The `statsmodels` package can help us perform model identification. The method `arma_order_select_ic` will find the optimal order of the ARMA model based on certain criteria. The first parameter `y` is the data. The data must be a NumPy array, not a Pandas DataFrame. The parameter `ic` defines the criteria trying to be minimized. The method will return a dictionary, where the minimal order of each criteria can be accessed.

```
>>> import statsmodels.api as sm
>>> from statsmodel.tsa.stattools import arma_order_select_ic as order_select
>>> import pandas as pd

>>> # Get sunspot data and give DateTimeIndex
>>> sunspot = sm.datasets.sunspots.load_pandas().data
>>> sunspot.index = pd.Index(pd.date_range("1700", end="2009", freq="A-DEC"))
>>> sunspot.drop(columns = ["YEAR"], inplace = True)

>>> # Find best order where p < 5 and q < 5
>>> # Use AICc as basis for minimization
>>> order = order_select(sunspot.values,max_ar=4,max_ma=4,ic=['aic','bic'],←
    fit_kw={'method': 'mle'})
```

```

>>> print(order['aic_min_order'])
(4,2)
>>> print(order['bic_min_order'])
(4,2)

>>> # Fit model
>>> # Note that we need to set the dimensionality to zero in order to have an ←
      ARMA model.
>>> model = ARIMA(sunspot,order = (4,0,2)).fit(method='innovations_mle')

>>> # Predict values from 1950 to 2012.
>>> prediction = model.predict(start='1950',end='2012')

>>> # Plot the prediction along with the sunspot data.
>>> fig, ax = plt.subplots(figsize=(13,7))
>>> plt.plot(prediction)
>>> plt.plot(sunspot['1950':'2009'])
>>> ax.set_title('Sunspot Dataset')
>>> ax.set_xlabel('Year')
>>> ax.set_ylabel('Number of Sunspots')
>>> plt.show()

```

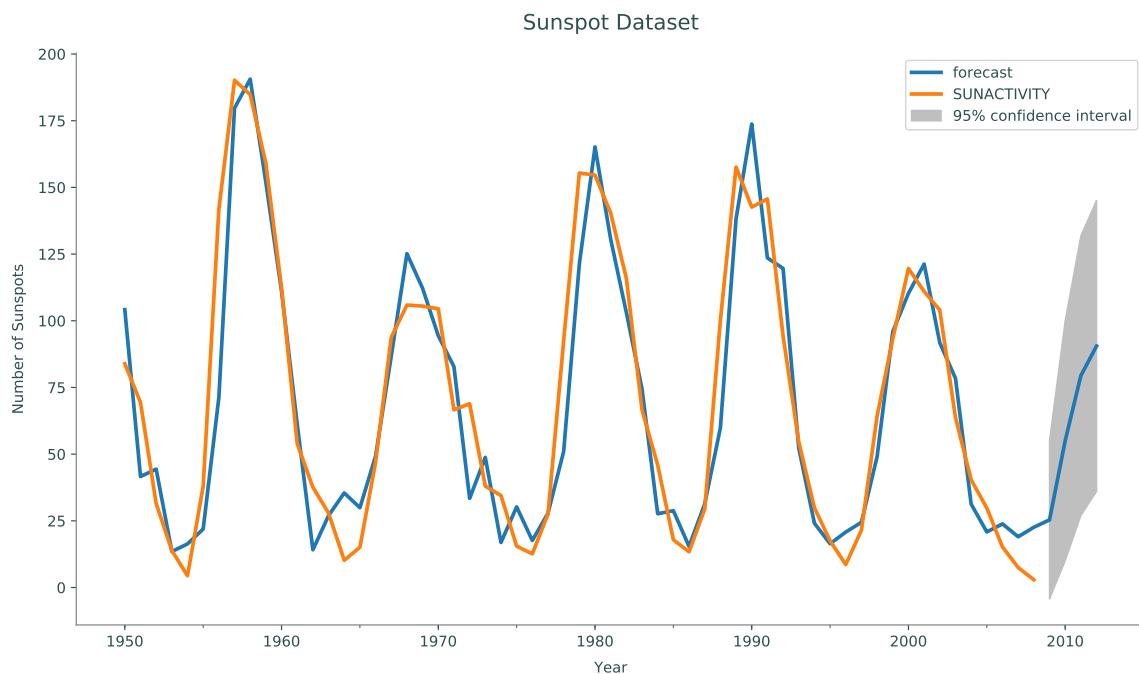


Figure 24.7: Sunspot activity data is forecasted four years in the future using `statsmodels`.

Problem 7. The dataset `manaus` contains data on the height of the Rio Negro from every month between January 1903 and January 1993. Write a function `manaus()` that accepts the forecasting range as strings `start` and `end`, the maximum parameter for the AR model `p` and the maximum parameter of the MA model `q`. The parameters `start` and `end` should be strings corresponding to a `DateTimeIndex` in the form `Y%M%D`, where `D` is the last day of the month.

The function should determine the optimal order for the ARMA model based on the AIC and the BIC. Then forecast and plot on the range given for both models and compare. Return the order of the AIC model and the order of the BIC model, respectively. For the range '`1983-01-31`' to '`1995-01-31`', your plot should look like Figure 24.8.

(Hint: The data passed into `arma_order_select_ic` must be a NumPy array. Use the attribute `values` of the Pandas DataFrame.)

To get the `manaus` dataset and set it with a `DateTimeIndex`, use the following code:

```
>>> # Get dataset
>>> raw = pydata('manaus')
>>> # Convert to DateTimeIndex
>>> manaus = pd.DataFrame(raw.values, index=pd.date_range('1903-01', '←
    1993-01', freq='M'))
>>> manaus = manaus.drop(0, axis=1)
>>> # Set new column title
>>> manaus.columns = ['Water Level']
```

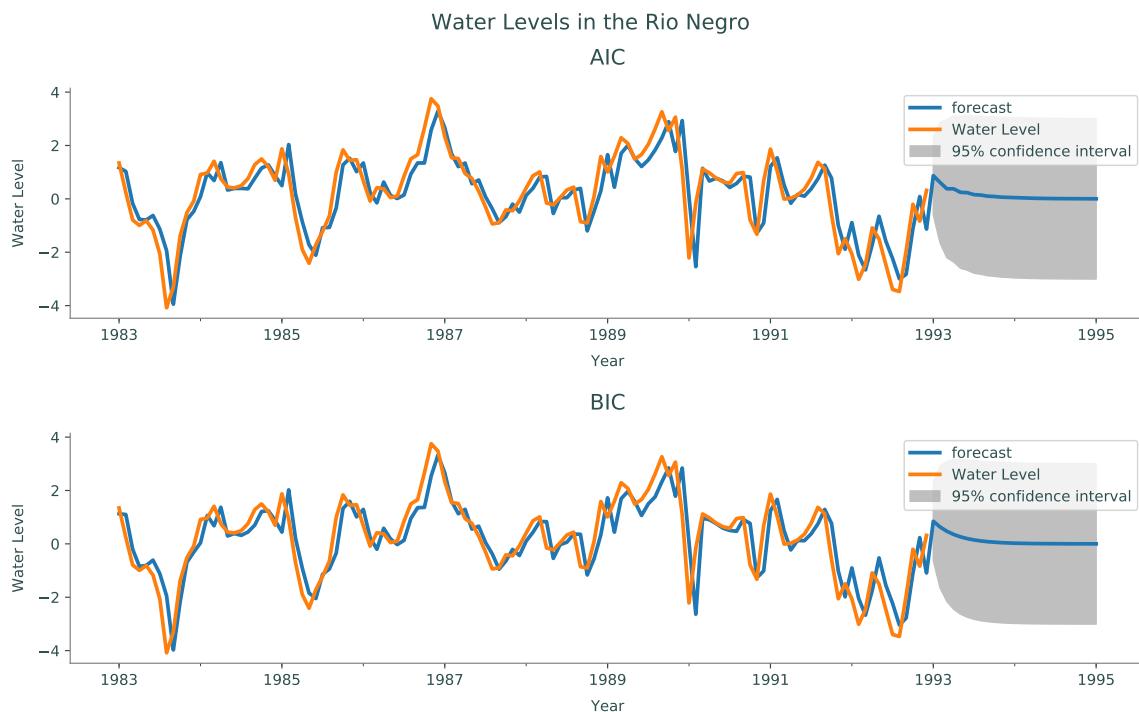


Figure 24.8: AIC and BIC based ARMA models of `manaus` dataset.

Additional Materials

Finding Error Correlation

To find the correlation of the current error with past error, the noise of the data needs to be isolated. Each data point y_t can be decomposed as

$$y_t = T_t + S_t + R_t, \quad (24.23)$$

where T_t is the overall trend of the data, S_t is a seasonal trend, and R_t is noise in the data. The overall trend is what the data tends to do as a whole, while the seasonal trend is what the data does repeatedly. For example, if looking at airfare prices over a decade, the overall trend of the data might be increasing due to inflation. However, we can break this data into individual years. We call each year a season. The seasonal trend of the data might not be strictly increasing, but have increases during busy seasons such as Christmas and summer vacations.

To find T_t , we use an M -fold method. In this case, M is the length of our season. We define the equation

$$T_t = \frac{1}{M} \sum_{-M/2 < i < M/2} y_{i+t}. \quad (24.24)$$

This means for each t , we take the average of the season surrounding y_t .

To find the seasonal trend, first subtract the overall trend from the time series. Define $x_t = y_t - T_t$. The value of the seasonal trend can then be found by averaging each day of the season over every season. For example, if the season was one year, we would find the average value on the first day of the year over all seasons, then the second, and so on. Thus,

$$S_t = \frac{1}{K} \sum_{i \equiv t \pmod{M}} x_i \quad (24.25)$$

where K is the number of seasons.

With the overall and seasonal trend known, the noise of the data is simply $R_t = y_t - T_t - S_t$. To determine the strength of correlations with the current error and the past error, plot y_t vs. R_{t-i} as in Figure 24.1.

Proof of Equation 24.15

$$\sum_{i=1}^p \phi_i(z_{t-i} - \mu) + a_t + \sum_{j=1}^q \theta_j a_{t-j} = \sum_{i=1}^p \phi_i(H\hat{\mathbf{x}}_{t-i}) + a_t + \sum_{j=1}^q \theta_j a_{t-j} \quad (24.26)$$

$$= \sum_{i=1}^r \phi_i(x_{t-i} + \sum_{k=1}^{r-1} \theta_k x_{t-i-k}) + a_t + \sum_{j=1}^{r-1} \theta_j a_{t-j} \quad (24.27)$$

$$= a_t + \sum_{i=1}^r \phi_i(x_{t-i}) + \sum_{j=1}^{r-1} \theta_j \left(\sum_{i=1}^r \phi_i x_{t-j-i} + a_{t-j} \right) \quad (24.28)$$

$$= a_t + \sum_{i=1}^r \phi_i(x_{t-i}) + \sum_{j=1}^{r-1} \theta_j x_{t-k} \quad (24.29)$$

$$= x_t + \sum_{j=1}^{r-1} \theta_j x_{t-k} \theta_k x_{t-k} \quad (24.30)$$

$$= z_t. \quad (24.31)$$

25 Non-negative Matrix Factorization Recommender

Lab Objective: Understand and implement the non-negative matrix factorization for recommendation systems.

Introduction

Collaborative filtering is the process of filtering data for patterns using collaboration techniques. More specifically, it refers to making prediction about a user's interests based on other users' interests. These predictions can be used to recommend items and are why collaborative filtering is one of the common methods of creating a recommendation system.

Recommendation systems look at the similarity between users to predict what item a user is most likely to enjoy. Common recommendation systems include Netflix's Movies you Might Enjoy list, Spotify's Discover Weekly playlist, and Amazon's Products You Might Like.

Non-negative Matrix Factorization

Non-negative matrix factorization is one algorithm used in collaborative filtering. It can be applied to many other cases, including image processing, text mining, clustering, and community detection. The purpose of non-negative matrix factorization is to take a non-negative matrix V and factor it into the product of two non-negative matrices.

For $V \in \mathbb{R}^{m \times n}$, $0 \preceq W$,

$$\begin{array}{ll}\text{minimize} & \|V - WH\| \\ \text{subject to} & 0 \preceq W, 0 \preceq H \\ \text{where} & W \in \mathbb{R}^{m \times k}, H \in \mathbb{R}^{k \times n}\end{array}$$

k is the rank of the decomposition and can either be specified or found using the Root Mean Squared Error (the square root of the MSE), SVD, Non-negative Least Squares, or cross-validation techniques.

For this lab, we will use the Frobenius norm, given by

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}.$$

It is equivalent to the square root of the sum of the diagonal of $A^H A$

Problem 1. Create the `NMFRecommender` class, which will be used to implement the NMF algorithm. Initialize the class with the following parameters: `random_state` defaulting to 15, `tol` defaulting to $1e-3$, `maxiter` defaulting to 200, and `rank` defaulting to 2.

Add a method called `initialize_matrices` that takes in m and n , the dimensions of V . Set the random seed so that initializing the matrices can be replicated.

```
>>> np.random.seed(self.random_state)
```

Then, using `np.random.random`, initialize W and H with randomly generated numbers between 0 and 1, where $W \in \mathbb{R}^{m \times k}$ and $H \in \mathbb{R}^{k \times n}$. Return W and H .

Finally, add a method called `compute_loss()` that takes as parameters `V`, `W`, and `H` and returns the Frobenius norm of $V - WH$.

Multiplicative Update

After initializing W and H , we iteratively update them using the multiplicative update step. There are other methods for optimization and updating, but because of the simplicity and ease of this solution, it is widely used. As with any other iterative algorithm, we perform the step until the `tol` or `maxiter` is met.

$$H_{ij}^{s+1} = H_{ij}^s \frac{((W^s)^T V)_{ij}}{((W^s)^T W^s H^s)_{ij}} \quad (25.1)$$

and

$$W_{ij}^{s+1} = W_{ij}^s \frac{(V(H^{s+1})^T)_{ij}}{(W^s H^{s+1} (H^{s+1})^T)_{ij}} \quad (25.2)$$

Problem 2. Add a method to the `NMF` class called `update_matrices` that takes as inputs matrices V , W , H and returns W_{s+1} and H_{s+1} as described in Equations 25.1 and 25.2.

Problem 3. Finish the NMF class by adding a method `fit` that finds an optimal W and H . It should accept `V` as a numpy array, perform the multiplicative update algorithm until the loss is less than `tol` or `maxiter` is reached, and return W and H .

Finally add a method called `reconstruct` that reconstructs and returns `V` by multiplying `W` and `H`.

Using NMF for Recommendations

Consider the following marketing problem where we have a list of five grocery store customers and their purchases. We want to create personalized food recommendations for their next visit. We start by creating a matrix representing each person and the number of items they purchased in different grocery categories. So from the matrix, we can see that John bought two fruits and one sweet.

$$V = \begin{pmatrix} John & Alice & Mary & Greg & Peter & Jennifer \\ 0 & 1 & 0 & 1 & 2 & 2 \\ 2 & 3 & 1 & 1 & 2 & 2 \\ 1 & 1 & 1 & 0 & 1 & 1 \\ 0 & 2 & 3 & 4 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \begin{array}{l} Vegetables \\ Fruits \\ Sweets \\ Bread \\ Coffee \end{array}$$

After performing NMF on V , we'll get the following W and H .

$$W = \begin{pmatrix} Component1 & Component2 & Component3 \\ 2.1 & 0.03 & 0. \\ 1.17 & 0.19 & 1.76 \\ 0.43 & 0.03 & 0.89 \\ 0.26 & 2.05 & 0.02 \\ 0.45 & 0. & 0. \end{pmatrix} \begin{array}{l} Vegetables \\ Fruits \\ Sweets \\ Bread \\ Coffee \end{array}$$

$$H = \begin{pmatrix} John & Alice & Mary & Greg & Peter & Jennifer \\ 0.00 & 0.45 & 0.00 & 0.43 & 1.0 & 0.9 \\ 0.00 & 0.91 & 1.45 & 1.9 & 0.35 & 0.37 \\ 1.14 & 1.22 & 0.55 & 0.0 & 0.47 & 0.53 \end{pmatrix} \begin{array}{l} Component1 \\ Component2 \\ Component3 \end{array}$$

W represents how much each grocery feature contributes to each component; a higher weight means it's more important to that component. For example, component 1 is heavily determined by vegetables followed by fruit, then coffee, sweets and finally bread. Component 2 is represented almost entirely by bread, while component 3 is based on fruits and sweets, with a small amount of bread. H is similar, except instead of showing how much each grocery category affects the component, it shows a much each person belongs to the component, again with a higher weight indicating that the person belongs more in that component. We can see the John belongs in component 3, while Jennifer mostly belongs in component 1.

To get our recommendations, we reconstruct V by multiplying W and H .

$$WH = \begin{pmatrix} John & Alice & Mary & Greg & Peter & Jennifer \\ 0.0000 & 0.9723 & 0.0435 & 0.96 & 2.1105 & 1.9011 \\ 2.0064 & 2.8466 & 1.2435 & 0.8641 & 2.0637 & 2.0561 \\ 1.0146 & 1.3066 & 0.533 & 0.2419 & 0.8588 & 0.8698 \\ 0.0228 & 2.0069 & 2.9835 & 4.0068 & 0.9869 & 1.0031 \\ 0.0000 & 0.2025 & 0.0000 & 0.1935 & 0.45 & 0.405 \end{pmatrix} \begin{array}{l} Vegetables \\ Fruits \\ Sweets \\ Bread \\ Coffee \end{array}$$

Most of the zeros from the original V have been filled in. This is the **collaborative filtering** portion of the algorithm. By sorting each column by weight, we can predict which items are more attractive to the customers. For instance, Mary has the highest weight for bread at 2.9835, followed by fruit at 1.2435 and then sweets at .533. So we would recommend bread to Mary.

Another way to interpret WH is to look at a feature and determine who is most likely to buy that item. So if we were having a sale on sweets but only had funds to let three people know, using the reconstructed matrix, we would want to target Alice, John, and Jennifer in that order. This gives us more information than V alone, which says that everyone except Greg bought one sweet.

Problem 4. Use the `NMFRecommender` class to run NMF on V , defined above, with 2 components. Return W , H as matrices, and the number of people who have higher weights in component 2 than in component 1 as a float.

Sklearn NMF

Python has a few packages for recommendation algorithms: Surprise, CaseRecommender and of course SkLearn. They implement various algorithms used in recommendation models. We'll use SkLearn, which is similar to the `NMFRecommender` class, for the last problems.

```
from sklearn.decomposition import NMF

>>> model = NMF(n_components=2, init='random', random_state=0)
>>> W = model.fit_transform(X)
>>> H = model.components_
```

As mentioned earlier, many big companies use recommendation systems to encourage purchasing, ad clicks, or spending more time in their product. One famous example of a Recommendation system is Spotify's Discover Weekly. Every week, Spotify creates a playlist of songs that the user has not listened to on Spotify. This helps users find new music that they enjoy and keeps Spotify at the forefront of music trends.

Problem 5. Read the file `artist_user.csv` as a pandas dataframe. The rows represent users, with the user id in the first column, and the columns represent artists. For each artist j that a user i has listened to, the ij entry contains the number of times user i has listened to artist j .

Identify the rank, or number of components to use. Ideally, we want the smallest rank that minimizes the error. However, this rank may be too computationally expensive, as in this situation. We'll choose the rank by using the following method. First, calculate the frobenius norm of the dataframe and multiply it by .0001. This will be our benchmark value. Next, iterate through $rank = 3, 4, 5, \dots$. For each iteration, run NMF using `n_components=rank` and reconstruct the matrix V . Calculate the root mean square error using `sklearn.metrics.mean_squared_error` of the original dataframe and the reconstructed matrix V . When the RMSE is less than the benchmark value, stop. Return the rank and the reconstructed matrix of this rank.

Problem 6. Write a function `discover_weekly` that takes in a user id and the reconstructed matrix from Problem 5, and returns a list of 30 artists to recommend as strings.

This list of strings should be sorted so that the first artist is the recommendation with the highest weight and the last artist is the least, and it should not contain any artists that the user has already listed to. Use the file `artists.csv` to match the artist ID to their name.

As a check, the Discover Weekly for user 2 should return

[‘Britney Spears’, ‘Avril Lavigne’, ‘Rihanna’, ‘Paramore’, ‘Christina Aguilera’,
‘U2’, ‘The Devil Wears Prada’, ‘Muse’, ‘Hadouken!’, ‘Ke\$ha’, ‘Good Charlotte’,
‘Linkin Park’, ‘Enter Shikari’, ‘Katy Perry’, ‘Miley Cyrus’, ‘Taylor Swift’,
‘Beyoncé’, ‘ Asking Alexandria’, ‘The Veronicas’, ‘Mariah Carey’, ‘Martin L. Gore’,
‘Dance Gavin Dance’, ‘Erasure’, ‘In Flames’, ‘3OH!3’, ‘Blur’, ‘Kelly Clarkson’,
‘Justin Bieber’, ‘Alesana’, ‘Ashley Tisdale’]

Part II

Appendices

A

Getting Started

The labs in this curriculum aim to introduce computational and mathematical concepts, walk through implementations of those concepts in Python, and use industrial-grade code to solve interesting, relevant problems. Lab assignments are usually about 5–10 pages long and include code examples (yellow boxes), important notes (green boxes), warnings about common errors (red boxes), and about 3–7 exercises (blue boxes). Get started by downloading the lab manual(s) for your course from <http://foundations-of-applied-mathematics.github.io/>.

Submitting Assignments

Labs

Every lab has a corresponding specifications file with some code to get you started and to make your submission compatible with automated test drivers. Like the lab manuals, these materials are hosted at <http://foundations-of-applied-mathematics.github.io/>.

Download the `.zip` file for your course, unzip the folder, and move it somewhere where it won't get lost. This folder has some setup scripts and a collection of folders, one per lab, each of which contains the specifications file(s) for that lab. See [Student-Materials/wiki/Lab-Index](#) for the complete list of labs, their specifications and data files, and the manual that each lab belongs to.

Achtung!

Do **not** move or rename the lab folders or the enclosed specifications files; if you do, the test drivers will not be able to find your assignment. Make sure your folder and file names match [Student-Materials/wiki/Lab-Index](#).

To submit a lab, modify the provided specifications file and use the file-sharing program specified by your instructor (discussed in the next section). The instructor will drop feedback files in the lab folder after grading the assignment. For example, the Introduction to Python lab has the specifications file `PythonIntro/python_intro.py`. To complete that assignment, modify `PythonIntro/python_intro.py` and submit it via your instructor's file-sharing system. After grading, the instructor will create a file called `PythonIntro/PythonIntro_feedback.txt` with your score and some feedback.

Homework

Non-lab coding homework should be placed in the `_Homework/` folder and submitted like a lab assignment. Be careful to name your assignment correctly so the instructor (and test driver) can find it. The instructor may drop specifications files and/or feedback files in this folder as well.

Setup

Achtung!

We strongly recommend using a Unix-based operating system (Mac or Linux) for the labs. Unix has a true bash terminal, works well with git and python, and is the preferred platform for computational and data scientists. It is possible to do this curriculum with Windows, but expect some road bumps along the way.

There are two ways to submit code to the instructor: with git (<http://git-scm.com/>), or with a file-syncing service like Google Drive. Your instructor will indicate which system to use.

Setup With Git

Git is a program that manages updates between an online code repository and the copies of the repository, called clones, stored locally on computers. If git is not already installed on your computer, download it at <http://git-scm.com/downloads>. If you have never used git, you might want to read a few of the following resources.

- Official git tutorial: <https://git-scm.com/docs/gittutorial>
- Bitbucket git tutorials: <https://www.atlassian.com/git/tutorials>
- GitHub git cheat sheet: services.github.com/.../github-git-cheat-sheet.pdf
- GitLab git tutorial: <https://docs.gitlab.com/ce/gitlab-basics/start-using-git.html>
- Codecademy git lesson: <https://www.codecademy.com/learn/learn-git>
- Training video series by GitHub: <https://www.youtube.com/playlist?list=PLg7.../>

There are many websites for hosting online git repositories. Your instructor will indicate which web service to use, but we only include instructions here for setup with Bitbucket.

1. Sign up. Create a Bitbucket account at <https://bitbucket.org>. If you use an academic email address (ending in `.edu`, etc.), you will get free unlimited public and private repositories.
2. Make a new repository. On the Bitbucket page, click the `+` button from the menu on the left and, under **CREATE**, select **Repository**. Provide a name for the repository, mark the repository as **private**, and make sure the repository type is **Git**. For **Include a README?**, select **No** (if you accidentally include a README, delete the repository and start over). Under **Advanced settings**, enter a short description for your repository, select **No forks** under forking, and select **Python** as the language. Finally, click the blue **Create repository** button. Take note of the URL of the webpage that is created; it should be something like <https://bitbucket.org/<name>/<repo>>.

3. Give the instructor access to your repository. On your newly created Bitbucket repository page (<https://bitbucket.org/<name>/<repo>> or similar), go to **Settings** in the menu to the left and select **User and group access**, the second option from the top. Enter your instructor's Bitbucket username under **Users** and click **Add**. Select the blue **Write** button so your instructor can read from and write feedback to your repository.
4. Create an SSH key. This step needs to be done only once on each computer that you want to be able to use to access your repository. If you have multiple repositories on the same computer, you do not need to repeat this step for each one. To create an SSH key, in a shell application (Terminal on Linux or Mac, or Git Bash (<https://gitforwindows.org/>) on Windows), enter the following command:

```
$ ssh-keygen
```

Press the Enter or Return key to accept the default file location. It will then prompt to enter a passphrase; this acts as a password to use the SSH key. If you do not want a passphrase, leave it blank and press Enter again. The key will then be created. The file for the key will be placed in in the `/home/<username>/ .ssh` directory on Linux; in `/Users/<username>/ .ssh` on macOS; and in `/c/users/<username>/ .ssh` on Windows.

Now that the key is created, you need to add it to your Bitbucket account. From Bitbucket, choose **Personal settings** and then **SSH keys**. Click **Add key** and enter a label (what it is doesn't matter). Now, using the file explorer, navigate to the SSH key you created, and open the public key file. The file will be called something like `id_rsa.pub`; do NOT use `id_rsa` (without the `.pub` extension). Copy the contents of this file, paste it into the Key field on Bitbucket, and press Save.

For more options and some troubleshooting information, refer to <https://support.atlassian.com/bitbucket-cloud/docs/set-up-an-ssh-key/>.

5. Connect your folder to the new repository. In a shell application (Terminal on Linux or Mac, or Git Bash (<https://gitforwindows.org/>) on Windows), enter the following commands.

```
# Navigate to your folder.
$ cd /path/to/folder # cd means 'change directory'.


# Make sure you are in the right place.
$ pwd # pwd means 'print working directory'.
/path/to/folder
$ ls *.md # ls means 'list files'.
README.md # This means README.md is in the working directory.


# Connect this folder to the online repository.
$ git init
$ git remote add origin git@bitbucket.org:<name>/<repo>.git


# Record your credentials.
$ git config --local user.name "your name"
$ git config --local user.email "your email"


# Add the contents of this folder to git and update the repository.
```

```
$ git add --all
$ git commit -m "initial commit"
$ git push origin master
```

For example, if your Bitbucket username is `greek314`, the repository is called `acmev1`, and the folder is called `Student-Materials/` and is on the desktop, enter the following commands.

```
# Navigate to the folder.
$ cd ~/Desktop/Student-Materials

# Make sure this is the right place.
$ pwd
/Users/Archimedes/Desktop/Student-Materials
$ ls *.md
README.md

# Connect this folder to the online repository.
$ git init
$ git remote add origin git@bitbucket.org:greek314/acmev1.git

# Record credentials.
$ git config --local user.name "archimedes"
$ git config --local user.email "greek314@example.com"

# Add the contents of this folder to git and update the repository.
$ git add --all
$ git commit -m "initial commit"
$ git push origin master
```

At this point you should be able to see the files on your repository page from a web browser. If you enter the repository URL incorrectly in the `git remote add origin` step, you can reset it with the following line:

```
$ git remote set-url origin git@bitbucket.org:<name>/<repo>.git
```

Note

You may get the an error like the following when you run `git push`:

```
remote: Bitbucket Cloud recently stopped supporting account passwords←
      for Git authentication.
...
fatal: Authentication failed for 'https://bitbucket.org/<name>/<repo←
>.git/'
```

If this error occurs, your repository URL is in the wrong format; most likely, you used the `https` version instead of what is shown above. You can use the `git remote set-url origin` command to fix this issue as well.

6. Download data files. Many labs have accompanying data files. To download these files, navigate to your clone and run the `download_data.sh` bash script, which downloads the files and places them in the correct lab folder for you. You can also find individual data files through [Student-Materials/wiki/Lab-Index](#).

```
# Navigate to your folder and run the script.
$ cd /path/to/folder
$ bash download_data.sh
```

7. Install Python package dependencies. The labs require several third-party Python packages that don't come bundled with Anaconda. Run the following command to install the necessary packages.

```
# Navigate to your folder and run the script.
$ cd /path/to/folder
$ bash install_dependencies.sh
```

8. (Optional) Clone your repository. If you want your repository on another computer after completing steps 1–5, use the following commands.

```
# Navigate to where you want to put the folder.
$ cd ~/Desktop/or/something/

# Clone the folder from the online repository.
$ git clone git@bitbucket.org:<name>/<repo>.git <foldername>

# Record your credentials in the new folder.
$ cd <foldername>
$ git config --local user.name "your name"
$ git config --local user.email "your email"

# Download data files to the new folder.
$ bash download_data.sh
```

Setup Without Git

Even if you aren't using git to submit files, you must install it (<http://git-scm.com/downloads>) in order to get the data files for each lab. Share your folder with your instructor according to their directions, and follow steps 6 and 7 of the previous section to download the data files and install package dependencies.

Using Git

Git manages the history of a file system through commits, or checkpoints. Use `git status` to see the files that have been changed since the last commit. These changes are then moved to the staging area, a list of files to save during the next commit, with `git add <filename(s)>`. Save the changes in the staging area with `git commit -m "<A brief message describing the changes>"`.

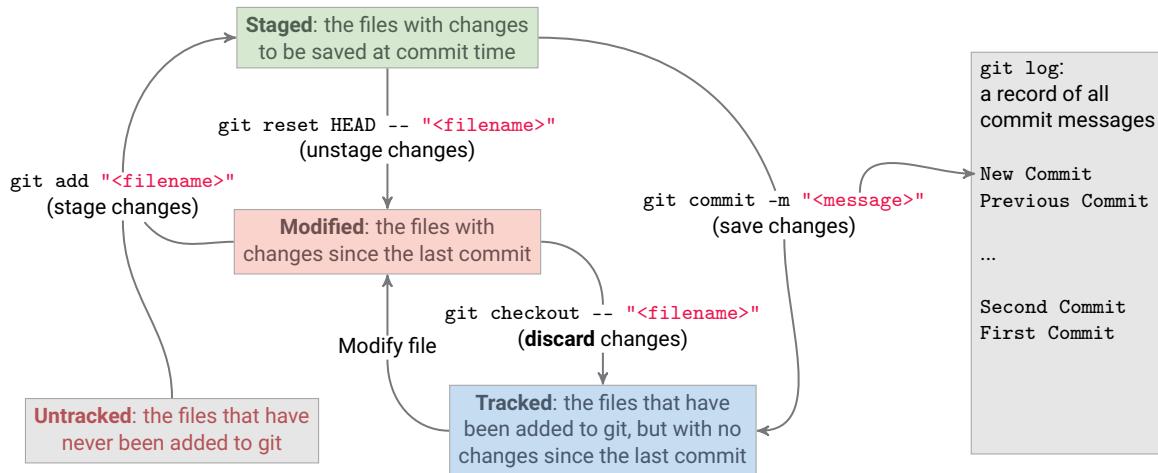


Figure A.1: Git commands to stage, unstage, save, or discard changes. Commit messages are recorded in the log.

All of these commands are done within a clone of the repository, stored somewhere on a computer. This repository must be manually synchronized with the online repository via two other git commands: `git pull origin master`, to pull updates from the web to the computer; and `git push origin master`, to push updates from the computer to the web.

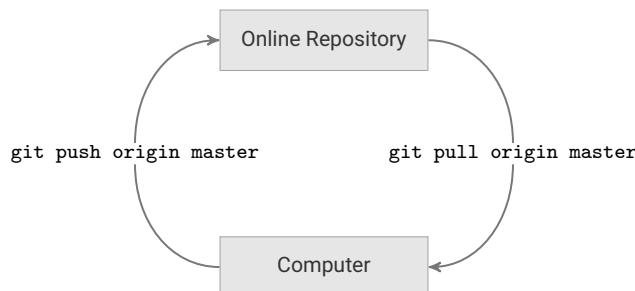


Figure A.2: Exchanging git commits between the repository and a local clone.

Command	Explanation
git status	Display the staging area and untracked changes.
git pull origin master	Pull changes from the online repository.
git push origin master	Push changes to the online repository.
git add <filename(s)>	Add a file or files to the staging area.
git add -u	Add all modified, tracked files to the staging area.
git commit -m "<message>"	Save the changes in the staging area with a given message.
git checkout -- <filename>	Revert changes to an unstaged file since the last commit.
git reset HEAD -- <filename>	Remove a file from the staging area.
git diff <filename>	See the changes to an unstaged file since the last commit.
git diff --cached <filename>	See the changes to a staged file since the last commit.
git config --local <option>	Record your credentials (<code>user.name</code> , <code>user.email</code> , etc.).

Table A.1: Common git commands.

Note

When pulling updates with `git pull origin master`, your terminal may sometimes display the following message.

```
Merge branch 'master' of git@bitbucket.org:<name>/<repo> into master

# Please enter a commit message to explain why this merge is necessary,
# especially if it merges an updated upstream into a topic branch.
#
# Lines starting with '#' will be ignored, and an empty message aborts
# the commit.
~
```

This means that someone else (the instructor) has pushed a commit that you do not yet have, while you have also made one or more commits locally that they do not have. This screen, displayed in vim ([https://en.wikipedia.org/wiki/Vim_\(text_editor\)](https://en.wikipedia.org/wiki/Vim_(text_editor))), is asking you to enter a message (or use the default message) to create a merge commit that will reconcile both changes. To close this screen and create the merge commit, type :wq and press `enter`.

Example Work Sessions

```
$ cd ~/Desktop/Student-Materials/
$ git pull origin master                                # Pull updates.
### Make changes to a file.
$ git add -u                                           # Track changes.
$ git commit -m "Made some changes."                  # Commit changes.
$ git push origin master                            # Push updates.
```

```
# Pull any updates from the online repository (such as TA feedback).
$ cd ~/Desktop/Student-Materials/
$ git pull origin master
From bitbucket.org:username/repo
 * branch            master      -> FETCH_HEAD
Already up-to-date.

### Work on the labs. For example, modify PythonIntro/python_intro.py.

$ git status
On branch master
Your branch is up-to-date with 'origin/master'.
Changes not staged for commit:
  (use "git add <file>..." to update what will be committed)
  (use "git checkout -- <file>..." to discard changes in working directory)

    PythonIntro/python_intro.py

# Track the changes with git.
$ git add PythonIntro/python_intro.py
$ git status
On branch master
Your branch is up-to-date with 'origin/master'.
Changes to be committed:
  (use "git reset HEAD <file>..." to unstage)

    modified:   PythonIntro/python_intro.py

# Commit the changes to the repository with an informative message.
$ git commit -m "Made some changes"
[master fed9b34] Made some changes
  1 file changed, 10 insertion(+) 1 deletion(-)

# Push the changes to the online repository.
$ git push origin master
Counting objects: 3, done.
Delta compression using up to 2 threads.
Compressing objects: 100% (2/2), done.
Writing objects: 100% (3/3), 327 bytes | 0 bytes/s, done.
Total 3 (delta 0), reused 0 (delta 0)
To git@bitbucket.org:username/repo.git
  5742a1b..fed9b34  master -> master

$ git status
On branch master
Your branch is up-to-date with 'origin/master'.
nothing to commit, working directory clean
```

B

Installing and Managing Python

Lab Objective: One of the great advantages of Python is its lack of overhead: it is relatively easy to download, install, start up, and execute. This appendix introduces tools for installing and updating specific packages and gives an overview of possible environments for working efficiently in Python.

Installing Python via Anaconda

A Python distribution is a single download containing everything needed to install and run Python, together with some common packages. For this curriculum, we **strongly** recommend using the Anaconda distribution to install Python. Anaconda includes IPython, a few other tools for developing in Python, and a large selection of packages that are common in applied mathematics, numerical computing, and data science. Anaconda is free and available for Windows, Mac, and Linux.

Follow these steps to install Anaconda.

1. Go to <https://www.anaconda.com/download/>.
2. Download the **Python 3.6** graphical installer specific to your machine.
3. Open the downloaded file and proceed with the default configurations.

For help with installation, see <https://docs.anaconda.com/anaconda/install/>. This page contains links to detailed step-by-step installation instructions for each operating system, as well as information for updating and uninstalling Anaconda.

Achtung!

This curriculum uses Python 3.6, **not** Python 2.7. With the wrong version of Python, some example code within the labs may not execute as intended or result in an error.

Managing Packages

A Python package manager is a tool for installing or updating Python packages, which involves downloading the right source code files, placing those files in the correct location on the machine, and linking the files to the Python interpreter. **Never** try to install a Python package without using a package manager (see <https://xkcd.com/349/>).

Conda

Many packages are not included in the default Anaconda download but can be installed via Anaconda's package manager, `conda`. See <https://docs.anaconda.com/anaconda/packages/pkg-docs> for the complete list of available packages. When you need to update or install a package, **always** try using `conda` first.

Command	Description
<code>conda install <package-name></code>	Install the specified package.
<code>conda update <package-name></code>	Update the specified package.
<code>conda update conda</code>	Update <code>conda</code> itself.
<code>conda update anaconda</code>	Update all packages included in Anaconda.
<code>conda --help</code>	Display the documentation for <code>conda</code> .

For example, the following terminal commands attempt to install and update `matplotlib`.

```
$ conda update conda          # Make sure that conda is up to date.
$ conda install matplotlib    # Attempt to install matplotlib.
$ conda update matplotlib     # Attempt to update matplotlib.
```

See <https://conda.io/docs/user-guide/tasks/manage-pkgs.html> for more examples.

Note

The best way to ensure a package has been installed correctly is to try importing it in IPython.

```
# Start IPython from the command line.
$ ipython
IPython 6.5.0 -- An enhanced Interactive Python. Type '?' for help.

# Try to import matplotlib.
In [1]: from matplotlib import pyplot as plt      # Success!
```

Achtung!

Be careful not to attempt to update a Python package while it is in use. It is safest to update packages while the Python interpreter is not running.

Pip

The most generic Python package manager is called `pip`. While it has a larger package list, `conda` is the cleaner and safer option. Only use `pip` to manage packages that are not available through `conda`.

Command	Description
<code>pip install package-name</code>	Install the specified package.
<code>pip install --upgrade package-name</code>	Update the specified package.
<code>pip freeze</code>	Display the version number on all installed packages.
<code>pip --help</code>	Display the documentation for <code>pip</code> .

See https://pip.pypa.io/en/stable/user_guide/ for more complete documentation.

Workflows

There are several different ways to write and execute programs in Python. Try a variety of workflows to find what works best for you.

Text Editor + Terminal

The most basic way of developing in Python is to write code in a text editor, then run it using either the Python or IPython interpreter in the terminal.

There are many different text editors available for code development. Many text editors are designed specifically for computer programming which contain features such as syntax highlighting and error detection, and are highly customizable. Try installing and using some of the popular text editors listed below.

- Atom: <https://atom.io/>
- Sublime Text: <https://www.sublimetext.com/>
- Notepad++ (Windows): <https://notepad-plus-plus.org/>
- Geany: <https://www.geany.org/>
- Vim: <https://www.vim.org/>
- Emacs: <https://www.gnu.org/software/emacs/>

Once Python code has been written in a text editor and saved to a file, that file can be executed in the terminal or command line.

```
$ ls                               # List the files in the current directory.
hello_world.py
$ cat hello_world.py               # Print the contents of the file to the terminal.
print("hello, world!")
$ python hello_world.py            # Execute the file.
hello, world!

# Alternatively, start IPython and run the file.
$ ipython
```

```
IPython 6.5.0 -- An enhanced Interactive Python. Type '?' for help.  
  
In [1]: %run hello_world.py  
hello, world!
```

IPython is an enhanced version of Python that is more user-friendly and interactive. It has many features that cater to productivity such as tab completion and object introspection.

Note

While Mac and Linux computers come with a built-in bash terminal, Windows computers do not. Windows does come with Powershell, a terminal-like application, but some commands in Powershell are different than their bash analogs, and some bash commands are missing from Powershell altogether. There are two good alternatives to the bash terminal for Windows:

- Windows subsystem for linux: docs.microsoft.com/en-us/windows/wsl/.
- Git bash: <https://gitforwindows.org/>.

Jupyter Notebook

The Jupyter Notebook (previously known as IPython Notebook) is a browser-based interface for Python that comes included as part of the Anaconda Python Distribution. It has an interface similar to the IPython interpreter, except that input is stored in cells and can be modified and re-evaluated as desired. See <https://github.com/jupyter/jupyter/wiki/> for some examples.

To begin using Jupyter Notebook, run the command `jupyter notebook` in the terminal. This will open your file system in a web browser in the Jupyter framework. To create a Jupyter Notebook, click the **New** drop down menu and choose **Python 3** under the **Notebooks** heading. A new tab will open with a new Jupyter Notebook.

Jupyter Notebooks differ from other forms of Python development in that notebook files contain not only the raw Python code, but also formatting information. As such, Jupyter Notebook files cannot be run in any other development environment. They also have the file extension `.ipynb` rather than the standard Python extension `.py`.

Jupyter Notebooks also support Markdown—a simple text formatting language—and L^AT_EX, and can embed images, sound clips, videos, and more. This makes Jupyter Notebook the ideal platform for presenting code.

Integrated Development Environments

An integrated development environment (IDEs) is a program that provides a comprehensive environment with the tools necessary for development, all combined into a single application. Most IDEs have many tightly integrated tools that are easily accessible, but come with more overhead than a plain text editor. Consider trying out each of the following IDEs.

- JupyterLab: <http://jupyterlab.readthedocs.io/en/stable/>
- PyCharm: <https://www.jetbrains.com/pycharm/>

- Spyder: <http://code.google.com/p/spyderlib/>
- Eclipse with PyDev: <http://www.eclipse.org/>, <https://www.pydev.org/>

See <https://realpython.com/python-ides-code-editors-guide/> for a good overview of these (and other) workflow tools.

C

NumPy Visual Guide

Lab Objective: NumPy operations can be difficult to visualize, but the concepts are straightforward. This appendix provides visual demonstrations of how NumPy arrays are used with slicing syntax, stacking, broadcasting, and axis-specific operations. Though these visualizations are for 1- or 2-dimensional arrays, the concepts can be extended to n -dimensional arrays.

Data Access

The entries of a 2-D array are the rows of the matrix (as 1-D arrays). To access a single entry, enter the row index, a comma, and the column index. Remember that indexing begins with 0.

$$A[0] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix} \quad A[2,1] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix}$$

Slicing

A lone colon extracts an entire row or column from a 2-D array. The syntax $[a:b]$ can be read as “the a th entry up to (but not including) the b th entry.” Similarly, $[a:]$ means “the a th entry to the end” and $[:b]$ means “everything up to (but not including) the b th entry.”

$$A[1] = A[1,:] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix} \quad A[:,2] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix}$$

$$A[1:,:2] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix} \quad A[1:-1,1:-1] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix}$$

Stacking

`np.hstack()` stacks sequence of arrays horizontally and `np.vstack()` stacks a sequence of arrays vertically.

$$A = \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix}$$

$$B = \begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}$$

$$\text{np.hstack}((A, B, A)) = \begin{bmatrix} \times & \times & \times & * & * & * & \times & \times & \times \\ \times & \times & \times & * & * & * & \times & \times & \times \\ \times & \times & \times & * & * & * & \times & \times & \times \end{bmatrix}$$

$$\text{np.vstack}((A, B, A)) = \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ * & * & * \\ * & * & * \\ * & * & * \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix}$$

Because 1-D arrays are flat, `np.hstack()` concatenates 1-D arrays and `np.vstack()` stacks them vertically. To make several 1-D arrays into the columns of a 2-D array, use `np.column_stack()`.

$$x = [\times \quad \times \quad \times \quad \times]$$

$$y = [* \quad * \quad * \quad *]$$

$$\text{np.hstack}((x, y, x)) = [\times \quad \times \quad \times \quad \times \quad * \quad * \quad * \quad * \quad \times \quad \times \quad \times \quad \times]$$

$$\text{np.vstack}((x, y, x)) = \begin{bmatrix} \times & \times & \times & \times \\ * & * & * & * \\ \times & \times & \times & \times \end{bmatrix}$$

$$\text{np.column_stack}((x, y, x)) = \begin{bmatrix} \times & * & \times \\ \times & * & \times \\ \times & * & \times \\ \times & * & \times \end{bmatrix}$$

The functions `np.concatenate()` and `np.stack()` are more general versions of `np.hstack()` and `np.vstack()`, and `np.row_stack()` is an alias for `np.vstack()`.

Broadcasting

NumPy automatically aligns arrays for component-wise operations whenever possible. See <http://docs.scipy.org/doc/numpy/user/basics.broadcasting.html> for more in-depth examples and broadcasting rules.

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix} \quad x = [10 \quad 20 \quad 30]$$

$$A + x = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \\ + \\ 10 & 20 & 30 \end{bmatrix} = \begin{bmatrix} 11 & 22 & 33 \\ 11 & 22 & 33 \\ 11 & 22 & 33 \end{bmatrix}$$

$$A + x.reshape((1, -1)) = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix} + \begin{bmatrix} 10 \\ 20 \\ 30 \end{bmatrix} = \begin{bmatrix} 11 & 12 & 13 \\ 21 & 22 & 23 \\ 31 & 32 & 33 \end{bmatrix}$$

Operations along an Axis

Most array methods have an `axis` argument that allows an operation to be done along a given axis. To compute the sum of each column, use `axis=0`; to compute the sum of each row, use `axis=1`.

$$A = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

$$A.sum(axis=0) = \left[\begin{array}{c|c|c|c} 1 & 2 & 3 & 4 \\ \hline 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{array} \right] = [4 \quad 8 \quad 12 \quad 16]$$

$$A.sum(axis=1) = \left[\begin{array}{cccc} 1 & 2 & 3 & 4 \\ \hline 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ \hline 1 & 2 & 3 & 4 \end{array} \right] = [10 \quad 10 \quad 10 \quad 10]$$

D

Introduction to Scikit-Learn

Lab Objective: Scikit-learn is one of the fundamental tools in Python for machine learning. In this appendix we highlight and give examples of some popular scikit-learn tools for classification and regression, training and testing, data normalization, and constructing complex models.

Note

This guide corresponds to scikit-learn version 0.20, which has a few significant differences from previous releases. See http://scikit-learn.org/stable/whats_new.html for current release notes. Install scikit-learn (the `sklearn` module) with `conda install scikit-learn`.

Base Classes and API

Many machine learning problems center on constructing a function $f : X \rightarrow Y$, called a **model** or **estimator**, that accurately represents properties of given data. The domain X is usually \mathbb{R}^D , and the range Y is typically either \mathbb{R} (regression) or a subset of \mathbb{Z} (classification). The model is trained on N samples $(\mathbf{x}_i)_{i=1}^N \subset X$ that usually (but not always) have N accompanying labels $(y_i)_{i=1}^N \subset Y$.

Scikit-learn [PVG⁺11, BLB⁺13] takes a highly object-oriented approach to machine learning models. Every major scikit-learn class inherits from `sklearn.base.BaseEstimator` and conforms to the following conventions:

1. The constructor `__init__()` receives hyperparameters for the classifier, which are parameters for the model f that are **not dependent on data**. Each hyperparameter must have a default value (i.e., every argument of `__init__()` is a keyword argument), and each argument must be saved as an instance variable of the **same name** as the parameter.
2. The `fit()` method constructs the model f . It receives an $N \times D$ matrix X and, optionally, a vector \mathbf{y} with N entries. Each row \mathbf{x}_i of X is one sample with corresponding label y_i . By convention, `fit()` always returns `self`.

Along with the `BaseEstimator` class, there are several other “mix in” base classes in `sklearn.base` that define specific kinds of models. The three listed below are the most common.¹

¹See <http://scikit-learn.org/stable/modules/classes.html#base-classes> for the complete list.

- **ClassifierMixin**: for classifiers, estimators that take on discrete values.
- **RegressorMixin**: for regressors, estimators that take on continuous values.
- **TransformerMixin**: for preprocessing data before estimation.

Classifiers and Regressors

The **ClassifierMixin** and **RegressorMixin** both require a `predict()` method that acts as the actual model f . That is, `predict()` receives an $N \times D$ matrix X and returns N predicted labels $(y_i)_{i=1}^N$, where y_i is the label corresponding to the i th row of X . Both of these base class have a predefined `score()` method that uses `predict()` to test the accuracy of the model. It accepts $N \times D$ test data and a vector of N corresponding labels, then reports either the classification accuracy (for classifiers) or the R^2 value of the regression (for regressors).

For example, a **KNeighborsClassifier** from `sklearn.neighbors` inherits from **BaseEstimator** and **ClassifierMixin**. This classifier uses a simple strategy: to classify a new piece of data \mathbf{z} , find the k training samples that are “nearest” to \mathbf{z} , then take the most common label corresponding to those nearest neighbors to be the label for \mathbf{z} . Its constructor accepts hyperparameters such as `n_neighbors`, for determining the number of neighbors k to search for, `algorithm`, which specifies the strategy to find the neighbors, and `n_jobs`, the number of parallel jobs to run during the neighbors search. Again, these hyperparameters are independent of any data, which is why they are set in the constructor (before fitting the model). Calling `fit()` organizes the data X into a data structure for efficient nearest neighbor searches (determined by `algorithm`). Calling `predict()` executes the search, determines the most common label of the neighbors, and returns that label.

```
>>> from sklearn.datasets import load_breast_cancer
>>> from sklearn.neighbors import KNeighborsClassifier
>>> from sklearn.model_selection import train_test_split

# Load the breast cancer dataset and split it into training and testing groups.
>>> cancer = load_breast_cancer()
>>> X_train, X_test, y_train, y_test = train_test_split(cancer.data,
...                                                     cancer.target)
>>> print(X_train.shape, y_train.shape)
(426, 30) # There are 426 training points, each with 30 features.

# Train a KNeighborsClassifier object on the training data.
# fit() returns the object, so we can instantiate and train in a single line.
>>> knn = KNeighborsClassifier(n_neighbors=2).fit(X_train, y_train)
# The hyperparameter 'n_neighbors' is saved as an attribute of the same name.
>>> knn.n_neighbors
2

# Test the classifier on the testing data.
>>> knn.predict(X_test[:6])
array([0, 1, 0, 1, 1, 0]) # Predicted labels for the first 6 test points.
>>> knn.score(X_test, y_test)
0.8951048951048951 # predict() chooses 89.51% of the labels right.
```

The `KNeighborsClassifier` object could easily be replaced with a different classifier, such as a `GaussianNB` object from `sklearn.naive_bayes`. Since `GaussianNB` also inherits from `BaseEstimator` and `ClassifierMixin`, it has `fit()`, `predict()`, and `score()` methods that take in the same kinds of inputs as the corresponding methods for the `KNeighborsClassifier`. The only difference, from an external perspective, is the hyperparameters that the constructor accepts.

```
>>> from sklearn.naive_bayes import GaussianNB

>>> gnb = GaussianNB().fit(X_train, y_train)
>>> gnb.predict(X_test[:6])
array([1, 1, 0, 1, 1, 0])
>>> gnb.score(X_test, y_test)
0.9440559440559441
```

Roughly speaking, the `GaussianNB` classifier assumes all features in the data are independent and normally distributed, then uses Bayes' rule to compute the likelihood of a new point belonging to a label for each of the possible labels. To do this, the `fit()` method computes the mean and variance of each feature, grouped by label. These quantities are saved as the attributes `theta_` (the means) and `sigma_` (the variances), then used in `predict()`. Parameters like these that **are dependent on data** are only defined in `fit()`, not the constructor, and they are always named with a trailing underscore. These “non-hyper” parameters are often simply called model parameters.

```
>>> gnb.classes_           # The collection of distinct training labels.
array([0, 1])
>>> gnb.theta_[:,0]       # The means of the first feature, grouped by label.
array([17.55785276, 12.0354981 ])
# The samples with label 0 have a mean of 17.56 in the first feature.
```

The `fit()` method should do all of the “heavy lifting” by calculating the model parameters. The `predict()` method should then use these parameters to choose a label for test data.

	Hyperparameters	Model Parameters
Data dependence	No	Yes
Initialization location	<code>__init__()</code>	<code>fit()</code>
Naming convention	Same as argument name	Ends with an underscore
Examples	<code>n_neighbors</code> , <code>algorithm</code> , <code>n_jobs</code>	<code>classes_</code> , <code>theta_</code> , <code>sigma_</code>

Table D.1: Naming and initialization conventions for scikit-learn model parameters.

Building Custom Estimators

The consistent conventions in the various scikit-learn classes makes it easy to use a wide variety of estimators with near-identical syntax. These conventions also make it possible to write custom estimators that behave like native scikit-learn objects. This usually only involves writing `fit()` and `predict()` methods and inheriting from the appropriate base classes. As a simple (though poorly performing) example, consider an estimator that either always predicts the same user-provided label, or that always predicts the most common label in the training data. Which strategy to use is independent of the data, so we encode that behavior with hyperparameters; the most common label must be calculated from the data, so that is a model parameter.

```

>>> import numpy as np
>>> from collections import Counter
>>> from sklearn.base import BaseEstimator, ClassifierMixin

>>> class PopularClassifier(BaseEstimator, ClassifierMixin):
...     """Classifier that always guesses the most common training label."""
...     def __init__(self, strategy="most_frequent", constant=None):
...         self.strategy = strategy      # Store the hyperparameters, using
...         self.constant = constant      # the same names as the arguments.
...
...     def fit(self, X, y):
...         """Find and store the most common label."""
...         self.popular_label_ = Counter(y).most_common(1)[0][0]
...         return self                  # fit() always returns 'self'.
...
...     def predict(self, X):
...         """Always guess the most popular training label."""
...         M = X.shape[0]
...         if self.strategy == "most_frequent":
...             return np.full(M, self.popular_label_)
...         elif self.strategy == "constant":
...             return np.full(M, self.constant)
...         else:
...             raise ValueError("invalid value for 'strategy' param")
...
# Train a PopularClassifier on the breast cancer training data.
>>> pc = PopularClassifier().fit(X_train, y_train)
>>> pc.popular_label_
1
# Score the model on the testing data.
>>> pc.score(X_test, y_test)
0.6573426573426573                         # 65.73% of the testing data is labeled 1.

# Change the strategy to always guess 0 by changing the hyperparameters.
>>> pc.strategy = "constant"
>>> pc.constant = 0
>>> pc.score(X_test, y_test)
0.34265734265734266                      # 34.27% of the testing data is labeled 0.

```

This is a terrible classifier, but it is actually implemented as `sklearn.dummy.DummyClassifier` because any legitimate machine learning algorithm should be able to beat it, so it is useful as a baseline comparison.

Note that `score()` was inherited from `ClassifierMixin` (it isn't defined explicitly), so it returns a classification rate. In the next example, a slight simplification of the equally unintelligent `sklearn.dummy.DummyRegressor`, the `score()` method is inherited from `RegressorMixin`, so it returns an R^2 value.

```

>>> from sklearn.base import RegressorMixin

>>> class ConstRegressor(BaseEstimator, RegressorMixin):
...     """Regressor that always predicts a mean or median of training data."""
...     def __init__(self, strategy="mean", constant=None):
...         self.strategy = strategy      # Store the hyperparameters, using
...         self.constant = constant      # the same names as the arguments.
...
...     def fit(self, X, y):
...         self.mean_, self.median_ = np.mean(y), np.median(y)
...         return self               # fit() always returns 'self'.
...
...     def predict(self, X):
...         """Always predict the middle of the training data."""
...         M = X.shape[0]
...         if self.strategy == "mean":
...             return np.full(M, self.mean_)
...         elif self.strategy == "median":
...             return np.full(M, self.median_)
...         elif self.strategy == "constant":
...             return np.full(M, self.constant)
...         else:
...             raise ValueError("invalid value for 'strategy' param")
...
# Train on the breast cancer data (treating it as a regression problem).
>>> cr = ConstRegressor(strategy="mean").fit(X_train, y_train)
>>> print("mean:", cr.mean_, " median:", cr.median_)
mean: 0.6173708920187794  median: 1.0

# Get the R^2 score of the regression on the testing data.
>>> cr.score(X_train, y_train)
0                         # Unsurprisingly, no correlation.

```

Achtung!

Both `PopularClassifier` and `ConstRegressor` wait until `predict()` to validate the `strategy` hyperparameter. The check could easily be done in the constructor, but that goes against scikit-learn conventions: in order to cooperate with automated validation tools, the constructor of any class inheriting from `BaseEstimator` must store the arguments of `__init__()` as attributes—with the same names as the arguments—and do nothing else.

Note

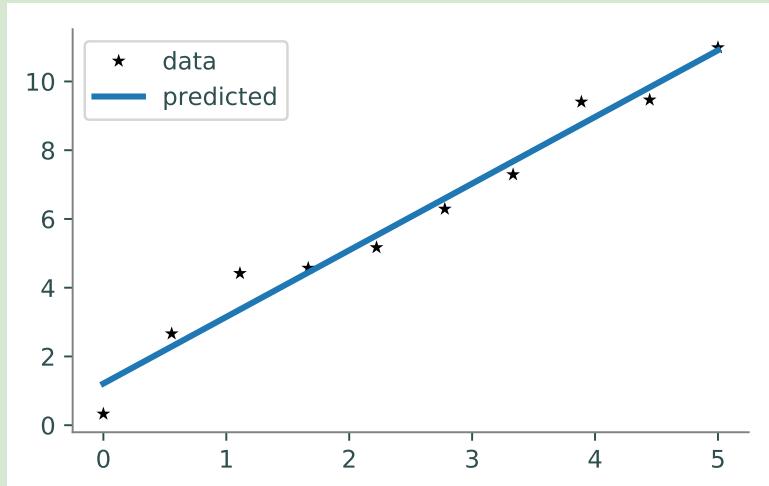
The first input to `fit()` and `predict()` are **always** two-dimensional $N \times D$ NumPy arrays, where N is the number of observations and D is the number of features. To fit or predict on one-dimensional data ($D = 1$), reshape the input array into a “column vector” before feeding it into the estimator. One-dimensional problems are somewhat rare in machine learning, but the following example shows how to do a simple one-dimensional linear regression.

```
>>> from matplotlib import pyplot as plt
>>> from sklearn.linear_model import LinearRegression

# Generate data for a 1-dimensional regression problem.
>>> X = np.linspace(0, 5, 10)
>>> Y = 2*X + 1 + np.random.normal(size=10)

# Reshape the training data into a column vector.
>>> lr = LinearRegression().fit(X.reshape((-1,1)), Y)

# Define another set of points to do predictions on.
>>> x = np.linspace(0, 5, 20)
>>> y = lr.predict(x.reshape((-1,1)))    # Reshape before predicting.
>>> plt.plot(X, Y, 'k*', label="data")
>>> plt.plot(x, y, label="predicted")
>>> plt.legend(loc="upper left")
>>> plt.show()
```



Transformers

A scikit-learn transformer processes data to make it better suited for estimation. This may involve shifting and scaling data, dropping columns, replacing missing values, and so on.

Classes that inherit from the `TransformerMixin` base class have a `fit()` method that accepts an $N \times D$ matrix X (like an estimator) and an optional set of labels. The labels are not needed—in fact the `fit()` method should do nothing with them—but the parameter for the labels remains as a keyword argument to be consistent with the `fit(X, y)` syntax of estimators. Instead of a `predict()` method, the `transform()` method accepts data, modifies it (usually via a copy), and returns the result. The new data may or may not have the same number of columns as the original data.

One common transformation is shifting and scaling the features (columns) so that they each have a mean of 0 and a standard deviation of 1. The following example implements a basic version of this transformer.

```
>>> from sklearn.base import TransformerMixin

>>> class NormalizingTransformer(BaseEstimator, TransformerMixin):
...     def fit(self, X, y=None):
...         """Calculate the mean and standard deviation of each column."""
...         self.mu_ = np.mean(X, axis=0)
...         self.sig_ = np.std(X, axis=0)
...         return self
...
...     def transform(self, X):
...         """Center each column at zero and normalize it."""
...         return (X - self.mu_) / self.sig_
...
# Fit the transformer and transform the cancer data (both train and test).
>>> nt = NormalizingTransformer()
>>> Z_train = nt.fit_transform(X_train) # Or nt.fit(X_train).transform(X_train)
>>> Z_test = nt.transform(X_test)      # Transform test data (without fitting)

>>> np.mean(Z_train, axis=0)[:3]       # The columns of Z_train have mean 0...
array([-8.08951237e-16, -1.72006384e-17,  1.78678147e-15])
>>> np.std(Z_train, axis=0)[:3]        # ...and have unit variance.
array([1., 1., 1.])
>>> np.mean(Z_test, axis=0)[:3]        # The columns of Z_test each have mean
array([-0.02355067,  0.11665332, -0.03996177])                      # close to 0...
>>> np.std(Z_test, axis=0)[:3]         # ...and have close to unit deviation.
array([0.9263711, 1.18461151, 0.91548103])

# Check to see if the classification improved.
>>> knn.fit(X_train, y_train).score(X_test, y_test)                  # Old score.
0.8951048951048951
>>> knn.fit(Z_train, y_train).score(Z_test, y_test)                  # New score.
0.958041958041958
```

This particular transformer is implemented as `sklearn.preprocessing.StandardScaler`. A close cousin is `sklearn.preprocessing.RobustScaler`, which ignores outliers when choosing the scaling and shifting factors.

Like estimators, transformers may have both hyperparameters (provided to the constructor) and model parameters (determined by `fit()`). Thus a transformer looks and acts like an estimator, with the exception of the `predict()` and `transform()` methods.

Achtung!

The `transform()` method should only rely on model parameters derived from the training data in `fit()`, **not** on the data that is worked on in `transform()`. For example, if the `NormalizingTransformer` is fit with the input \hat{X} , then `transform()` should shift and scale any input X by the mean and standard deviation of \hat{X} , not by the mean and standard deviation of X . Otherwise, the transformation is different for each input X .

Scikit-learn Module	Classifier Name	Notable Hyperparameters
<code>discriminant_analysis</code>	<code>LinearDiscriminantAnalysis</code>	<code>solver</code> , <code>shrinkage</code> , <code>n_components</code>
<code>discriminant_analysis</code>	<code>QuadraticDiscriminantAnalysis</code>	<code>reg_param</code>
<code>ensemble</code>	<code>AdaBoostClassifier</code>	<code>n_estimators</code> , <code>learning_rate</code>
<code>ensemble</code>	<code>RandomForestClassifier</code>	<code>n_estimators</code> , <code>max_depth</code>
<code>linear_model</code>	<code>LogisticRegression</code>	<code>penalty</code> , <code>C</code>
<code>linear_model</code>	<code>SGDClassifier</code>	<code>loss</code> , <code>penalty</code> , <code>alpha</code>
<code>naive_bayes</code>	<code>GaussianNB</code>	<code>priors</code>
<code>naive_bayes</code>	<code>MultinomialNB</code>	<code>alpha</code>
<code>neighbors</code>	<code>KNeighborsClassifier</code>	<code>n_neighbors</code> , <code>weights</code>
<code>neighbors</code>	<code>RadiusNeighborsClassifier</code>	<code>radius</code> , <code>weights</code>
<code>neural_network</code>	<code>MLPClassifier</code>	<code>hidden_layer_size</code> , <code>activation</code>
<code>svm</code>	<code>SVC</code>	<code>C</code> , <code>kernel</code>
<code>tree</code>	<code>DecisionTreeClassifier</code>	<code>max_depth</code>
Scikit-learn Module	Regressor Name	Notable Hyperparameters
<code>ensemble</code>	<code>AdaBoostRegressor</code>	<code>n_estimators</code> , <code>learning_rate</code>
<code>ensemble</code>	<code>ExtraTreesRegressor</code>	<code>n_estimators</code> , <code>max_depth</code>
<code>ensemble</code>	<code>GradientBoostingRegressor</code>	<code>n_estimators</code> , <code>max_depth</code>
<code>ensemble</code>	<code>RandomForestRegressor</code>	<code>n_estimators</code> , <code>max_depth</code>
<code>isotonic</code>	<code>IsotonicRegression</code>	<code>y_min</code> , <code>y_max</code>
<code>kernel_ridge</code>	<code>KernelRidge</code>	<code>alpha</code> , <code>kernel</code>
<code>linear_model</code>	<code>LinearRegression</code>	<code>fit_intercept</code>
<code>neural_network</code>	<code>MLPRegressor</code>	<code>hidden_layer_size</code> , <code>activation</code>
<code>svm</code>	<code>SVR</code>	<code>C</code> , <code>kernel</code>
<code>tree</code>	<code>DecisionTreeRegressor</code>	<code>max_depth</code>
Module	Transformer Name	Notable Hyperparameters
<code>decomposition</code>	<code>PCA</code>	<code>n_components</code>
<code>preprocessing</code>	<code>Imputer</code>	<code>missing_values</code> , <code>strategy</code>
<code>preprocessing</code>	<code>MinMaxScaler</code>	<code>feature_range</code>
<code>preprocessing</code>	<code>OneHotEncoder</code>	<code>categorical_features</code>
<code>preprocessing</code>	<code>QuantileTransformer</code>	<code>n_quantiles</code> , <code>output_distribution</code>
<code>preprocessing</code>	<code>RobustScaler</code>	<code>with_centering</code> , <code>with_scaling</code>
<code>preprocessing</code>	<code>StandardScaler</code>	<code>with_mean</code> , <code>with_std</code>

Table D.2: Common scikit-learn classifiers, regressors, and transformers. For full documentation on these classes, see <http://scikit-learn.org/stable/modules/classes.html>.

Validation Tools

Knowing how to determine whether or not an estimator performs well is an essential part of machine learning. This often turns out to be a surprisingly sophisticated issue that largely depends on the type of problem being solved and the kind of data that is available for training. Scikit-learn has validation tools for many situations; for brevity, we restrict our attention to the simple (but important) case of binary classification, where the range of the desired model is $Y = \{0, 1\}$.

Evaluation Metrics

The `score()` method of a scikit-learn estimator representing the model $f : X \rightarrow \{0, 1\}$ returns the accuracy of the model, which is the percent of labels that are predicted correctly. However, accuracy isn't always the best measure of success. Consider the confusion matrix for a classifier, the matrix where the (i, j) th entry is the number of observations with actual label i but that are classified as label j . In binary classification, calling the class with label 0 the negatives and the class with label 1 the positives, this becomes the following.

$$\begin{array}{cc} & \text{Predicted: 0} & \text{Predicted: 1} \\ \text{Actual: 0} & \begin{bmatrix} \text{True Negatives (TN)} & \text{False Positives (FP)} \\ \text{False Negatives (FN)} & \text{True Positives (TP)} \end{bmatrix} \\ \text{Actual: 1} & & \end{array}$$

With this terminology, we define the following metrics.

- Accuracy: $\frac{TN + TP}{TN + FN + FP + TP}$, the percent of labels predicted correctly.
- Precision: $\frac{TP}{TP + FP}$, the percent of predicted positives that are actually correct.
- Recall: $\frac{TP}{TP + FN}$, the percent of actual positives that are predicted correctly.

Precision is useful in situations where false positives are dangerous or costly, while recall is important when avoiding false negatives takes priority. For example, an email spam filter should avoid filtering out an email that isn't actually spam, so precision is a valuable metric for the filter. On the other hand, recall is more important in disease detection: it is better to test positive and not have the disease than to test negative when the disease is actually present. Focusing on a single metric often leads to skewed results (for example, always predicting the same label), so the following metric is also common.

$$\bullet F_\beta \text{ Score: } (1 + \beta^2) \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}} = \frac{(1 + \beta^2)TP}{(1 + \beta^2)TP + FP + \beta^2FN}.$$

Choosing $\beta < 1$ weighs precision more than recall, while $\beta > 1$ prioritizes recall over precision. The choice of $\beta = 1$ yields the common F_1 score, which weighs precision and recall equally. This is an important alternative to accuracy when, for example, the training set is heavily unbalanced with respect to the class labels.

Scikit-learn implements these metrics in `sklearn.metrics`, as well as functions for evaluating regression, non-binary classification, and clustering models. The general syntax for such functions is `some_score(actual_labels, predicted_labels)`. For the complete list and further discussion, see http://scikit-learn.org/stable/modules/model_evaluation.html.

```

>>> from sklearn.metrics import (confusion_matrix, classification_report,
...                               accuracy_score, precision_score,
...                               recall_score, f1_score)

# Fit the estimator to training data and predict the test labels.
>>> knn.fit(X_train, y_train)
>>> knn_predicted = knn.predict(X_test)

# Compute the confusion matrix by comparing actual labels to predicted labels.
>>> CM = confusion_matrix(y_test, knn_predicted)
>>> CM
array([[44,  5],
       [10,  84]])

# Get accuracy (the "usual" score), precision, recall, and f1 scores.
>>> accuracy_score(y_test, knn_predicted)    # (CM[0,0] + CM[1,1]) / CM.sum()
0.8951048951048951
>>> precision_score(y_test, knn_predicted)   # CM[1,1] / CM[:,1].sum()
0.9438202247191011
>>> recall_score(y_test, knn_predicted)       # CM[1,1] / CM[1,:].sum()
0.8936170212765957
>>> f1_score(y_test, knn_predicted)
0.9180327868852459

# Get all of these scores at once with classification_report().
>>> print(classification_report(y_test, knn_predicted))
      precision    recall  f1-score   support

          0       0.81      0.90      0.85       49
          1       0.94      0.89      0.92       94

   micro avg       0.90      0.90      0.90      143
   macro avg       0.88      0.90      0.89      143
weighted avg       0.90      0.90      0.90      143

```

Cross Validation

The `sklearn.model_selection` module has utilities to streamline and improve model evaluation.

- `train_test_split()` randomly splits data into training and testing sets (we already used this).
- `cross_val_score()` randomly splits the data and trains and scores the model a set number of times. Each trial uses different training data and results in a different model. The function returns the score of each trial.
- `cross_validate()` does the same thing as `cross_val_score()`, but it also reports the time it took to fit, the time it took to score, and the scores for the test set as well as the training set.

Doing multiple evaluations with different testing and training sets is extremely important. If the scores on a cross validation test vary wildly, the model is likely overfitting to the training data.

```
>>> from sklearn.model_selection import cross_val_score, cross_validate

# Make (but do not train) a classifier to test.
>>> knn = KNeighborsClassifier(n_neighbors=3)

# Test the classifier on the training data 4 times.
>>> cross_val_score(knn, X_train, y_train, cv=4)
array([0.88811189, 0.92957746, 0.96478873, 0.92253521])

# Get more details on the train/test procedure.
>>> cross_validate(knn, X_train, y_train, cv=4,
...                  return_train_score=False)
{'fit_time': array([0.00064683, 0.00042295, 0.00040913, 0.00040436]),
 'score_time': array([0.00115728, 0.00109601, 0.00105286, 0.00102782]),
 'test_score': array([0.88811189, 0.92957746, 0.96478873, 0.92253521])}

# Do the scoring with an alternative metric.
>>> cross_val_score(knn, X_train, y_train, scoring="f1", cv=4)
array([0.93048128, 0.95652174, 0.96629213, 0.93103448])
```

Note

Any estimator, even a user-defined class, can be evaluated with the scikit-learn tools presented in this section as long as that class conforms to the scikit-learn API discussed previously (i.e., inheriting from the correct base classes, having `fit()` and `predict()` methods, managing hyperparameters and parameters correctly, and so on). Any time you define a custom estimator, following the scikit-learn API gives you instant access to tools such as `cross_val_score()`.

Grid Search

Recall that the hyperparameters of a machine learning model are user-provided parameters that do not depend on the training data. Finding the optimal hyperparameters for a given model is a challenging and active area of research.² However, brute-force searching over a small hyperparameter space is simple in scikit-learn: a `sklearn.model_selection.GridSearchCV` object is initialized with an estimator, a dictionary of hyperparameters, and cross validation parameters (such as `cv` and `scoring`). When its `fit()` method is called, it does a cross validation test on the given estimator with every possible hyperparameter combination.

For example, a k -neighbors classifier has a few important hyperparameters that can have a significant impact on the speed and accuracy of the model: `n_neighbors`, the number of nearest neighbors allowed to vote; and `weights`, which specifies a strategy for weighting the distances between points. The following code tests various combinations of these hyperparameters.

²Intelligent hyperparameter selection is sometimes called metalearning. See, for example, [SGCP⁺18].

```
>>> from sklearn.model_selection import GridSearchCV

>>> knn = KNeighborsClassifier()
# Specify the hyperparameters to vary and the possible values they should take.
>>> param_grid = {"n_neighbors": [2, 3, 4, 5, 6],
...                 "weights": ["uniform", "distance"]}
>>> knn_gs = GridSearchCV(knn, param_grid, cv=4, scoring="f1", verbose=1)
>>> knn_gs.fit(X_train, y_train)
Fitting 4 folds for each of 5 candidates, totalling 20 fits
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent worker.
[Parallel(n_jobs=1)]: Done 20 out of 20 | elapsed: 0.1s finished

# After fitting, the gridsearch object has data about the results.
>>> print(knn_gs.best_params_, knn_gs.best_score_)
{'n_neighbors': 5, 'weights': 'uniform'} 0.9532526583188765
```

The cost of a grid search rapidly increases as the hyperparameter space grows. However, the outcomes of each trial are completely independent of each other, so the problem of training each classifier is embarrassingly parallel. To parallelize the grid search over n cores, set the `n_jobs` parameter to n , or set it to -1 to divide the labor between as many cores as are available.

In some circumstances, the parameter grid can be also organized in a way that eliminates redundancy. Consider an SVC classifier from `sklearn.svm`, an estimator that works by lifting the data into a high-dimensional space, then constructing a hyperplane to separate the classes. The SVC has a hyperparameter, `kernel`, that determines how the lifting into higher dimensions is done, and for each choice of kernel there are additional corresponding hyperparameters. To search the total hyperparameter space without redundancies, enter the parameter grid as a list of dictionaries, each of which defines a different section of the hyperparameter space. In the following code, doing so reduces the number of trials from $3 \times 2 \times 3 \times 4 = 72$ to only $1 + (1 \times 1 \times 3) + (1 \times 4) = 11$.

```
>>> from sklearn.svm import SVC

>>> svc = SVC(C=0.01, max_iter=100)
>>> param_grid = [
...     {"kernel": ["linear"]},
...     {"kernel": ["poly"], "degree": [2,3], "coef0": [0,1,5]},
...     {"kernel": ["rbf"], "gamma": [.01, .1, 1, 100]}
]
>>> svc_gs = GridSearchCV(svc, param_grid,
...                         cv=4, scoring="f1",
...                         verbose=1, n_jobs=-1).fit(X_train, y_train)
Fitting 4 folds for each of 11 candidates, totalling 44 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done 44 out of 44 | elapsed: 2.4s finished

>>> print(svc_gs.best_params_, svc_gs.best_score_)
{'gamma': 0.01, 'kernel': 'rbf'} 0.8909310239174055
```

See https://scikit-learn.org/stable/modules/grid_search.html for more details about `GridSearchCV` and its relatives.

Pipelines

Most machine learning problems require at least a little data preprocessing before estimation in order to get good results. A scikit-learn pipeline (`sklearn.pipeline.Pipeline`) chains together one or more transformers and one estimator into a single object, complete with `fit()` and `predict()` methods. For example, it is often a good idea to shift and scale data before feeding it into a classifier. The `StandardScaler` transformer can be combined with a classifier with a pipeline. Calling `fit()` on the resulting object calls `fit_transform()` on each successive transformer, then `fit()` on the estimator at the end. Likewise, calling `predict()` on the `Pipeline` object calls `transform()` on each transformer, then `predict()` on the estimator.

```
>>> from sklearn.preprocessing import StandardScaler
>>> from sklearn.pipeline import Pipeline

# Chain together a scaler transformer and a KNN estimator.
>>> pipe = Pipeline([("scaler", StandardScaler()),           # "scaler" is a label.
                    ("knn", KNeighborsClassifier())]) # "knn" is a label.
>>> pipe.fit(X_train, y_train)
>>> pipe.score(X_test, y_test)
0.972027972027972                                         # Already an improvement!
```

Since `Pipeline` objects behaves like estimators (following the `fit()` and `predict()` conventions), they can be used with tools like `cross_val_score()` and `GridSearchCV`. To specify which hyperparameters belong to which steps of the pipeline, precede each hyperparameter name with `<stepname>__`. For example, `knn__n_neighbors` corresponds to the `n_neighbors` hyperparameter of the part of the pipeline that is labeled `knn`.

```
# Specify the possible hyperparameters for each step.
>>> pipe_param_grid = {"scaler__with_mean": [True, False],
...                      "scaler__with_std": [True, False],
...                      "knn__n_neighbors": [2,3,4,5,6],
...                      "knn__weights": ["uniform", "distance"]}

# Pass the Pipeline object to the GridSearchCV and fit it to the data.
>>> pipe = Pipeline([("scaler", StandardScaler()),
                    ("knn", KNeighborsClassifier())])
>>> pipe_gs = GridSearchCV(pipe, pipe_param_grid,
...                         cv=4, n_jobs=-1, verbose=1).fit(X_train, y_train)
Fitting 4 folds for each of 40 candidates, totalling 160 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done 160 out of 160 | elapsed:    0.3s finished

>>> print(pipe_gs.best_params_, pipe_gs.best_score_, sep='\n')
{'knn__n_neighbors': 6, 'knn__weights': 'distance',
 'scaler__with_mean': True, 'scaler__with_std': True}
0.971830985915493
```

Pipelines can also be used to compare different transformations or estimators. For example, a pipeline could end in either a `KNeighborsClassifier()` or an `SVC()`, even though they have different hyperparameters. Like before, use a list of dictionaries to specify the hyperparameter space.

Additional Material

Exercises

Problem 1. Writing custom scikit-learn transformers is a convenient way to organize the data cleaning process. Consider the data in `titanic.csv`, which contains information about passengers on the maiden voyage of the RMS Titanic in 1912. Write a custom transformer class to clean this data, implementing the `transform()` method as follows:

1. Extract a copy of data frame with just the `"Pclass"`, `"Sex"`, and `"Age"` columns.
2. Replace `NaN` values in the `"Age"` column (of the copied data frame) with the mean age. The mean age of the training data should be calculated in `fit()` and used in `transform()` (compare this step to using `sklearn.preprocessing.Imputer`).
3. Convert the `"Pclass"` column datatype to pandas categoricals (`pd.CategoricalIndex`).
4. Use `pd.get_dummies()` to convert the categorical columns to multiple binary columns (compare this step to using `sklearn.preprocessing.OneHotEncoder`).
5. Cast the result as a NumPy array and return it.

Ensure that your transformer matches scikit-learn conventions (it inherits from the correct base classes, `fit()` returns `self`, etc.).

Problem 2. Read the data from `titanic.csv` with `pd.read_csv()`. The `"Survived"` column indicates which passengers survived, so the entries of the column are the labels that we would like to predict. Drop any rows in the raw data that have `NaN` values in the `"Survived"` column, then separate the column from the rest of the data. Split the data and labels into training and testing sets. Use the training data to fit a transformer from Problem 1, then use that transformer to clean the training set, then the testing set. Finally, train a `LogisticRegressionClassifier` and a `RandomForestClassifier` on the cleaned training data, and score them using the cleaned test set.

Problem 3. Use `classification_report()` to score your classifiers from Problem 2. Next, do a grid search for each classifier (using only the cleaned training data), varying at least two hyperparameters for each kind of model. Use `classification_report()` to score the resulting best estimators with the cleaned test data. Try changing the hyperparameter spaces or scoring metrics so that each grid search yields a better estimator.

Problem 4. Make a pipeline with at least two transformers to further process the Titanic dataset. Do a gridsearch on the pipeline and report the hyperparameters of the best estimator.

Bibliography

- [ADH⁺01] David Ascher, Paul F Dubois, Konrad Hinsen, Jim Hugunin, Travis Oliphant, et al. Numerical python, 2001.
- [BLB⁺13] Lars Buitinck, Gilles Louppe, Mathieu Blondel, Fabian Pedregosa, Andreas Mueller, Olivier Grisel, Vlad Niculae, Peter Prettenhofer, Alexandre Gramfort, Jaques Grobler, Robert Layton, Jake VanderPlas, Arnaud Joly, Brian Holt, and Gaël Varoquaux. API design for machine learning software: experiences from the scikit-learn project. In ECML PKDD Workshop: Languages for Data Mining and Machine Learning, pages 108–122, 2013.
- [Hun07] J. D. Hunter. Matplotlib: A 2d graphics environment. Computing In Science & Engineering, 9(3):90–95, 2007.
- [Oli06] Travis E Oliphant. A guide to NumPy, volume 1. Trelgol Publishing USA, 2006.
- [Oli07] Travis E Oliphant. Python for scientific computing. Computing in Science & Engineering, 9(3), 2007.
- [PVG⁺11] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12:2825–2830, 2011.
- [SGCP⁺18] Brandon Schoenfeld, Christophe Giraud-Carrier, Mason Poggemann, Jarom Christensen, and Kevin Seppi. Preprocessor selection for machine learning pipelines. arXiv preprint arXiv:1810.09942, 2018.
- [VD10] Guido VanRossum and Fred L Drake. The python language reference. Python software foundation Amsterdam, Netherlands, 2010.