

# UCLA Data Science Intensive

Joshua Cook

August 29, 2018



# Contents

<b>1</b>	<b>Introductory Lectures</b>	<b>9</b>
1.1	Vectors and Matrices . . . . .	9
1.2	Statistical Learning . . . . .	17
1.3	The Bias-Variance Tradeoff . . . . .	28
1.4	Probability . . . . .	36
1.5	Probabilistic Model Selection / In-Sample Model Selection . . . .	40
1.6	Using Masks to Create a Vector . . . . .	49
1.7	The Curse of Dimensionality . . . . .	53
1.8	Introduction to the bias-variance tradeoff . . . . .	57
<b>2</b>	<b>Case Study: Seeds</b>	<b>81</b>
2.1	Structure of the Case Study . . . . .	81
2.2	Load and Verify Data . . . . .	83
2.3	Exploratory Data Analysis . . . . .	91
2.4	Building a Predictive Model . . . . .	101
<b>3</b>	<b>Configuring AWS</b>	<b>103</b>
3.1	Amazon Web Services . . . . .	103
3.2	Configure your Local System . . . . .	104
3.3	Create a New Key Pair . . . . .	105
3.4	Configure your AWS Account . . . . .	108
3.5	Launch a New EC2 Instance . . . . .	110
3.6	Git and Github . . . . .	114
3.7	Learning to read the Bash Prompt . . . . .	117
3.8	Test your SSH Connection to Github . . . . .	118
3.9	Docker . . . . .	118



# List of Figures

1.1	Plot the line defining the relationship between $x$ and $y$ . . . . .	19
1.2	Plot new $x$ and $y$ . . . . .	20
1.3	Plot the Advertising Data . . . . .	21
1.4	Plot the line of best fit . . . . .	24
1.5	Plot lines of best fit for higher order functions . . . . .	24
1.6	Plot line of best fit . . . . .	25
1.7	Plot higher order lines of best fit . . . . .	26
1.8	Model Flexibility . . . . .	27
1.9	Plot the bias versus the variance . . . . .	35
1.10	A population of 1000 patients . . . . .	39
1.11	William of Occam . . . . .	44
1.12	Plot Fit Time and Prediction Time . . . . .	56
1.13	Plot Fit Time and Prediction Time with Logarithmic Scale for $n$ . . . . .	57
1.14	Plot true morale over the span of a course . . . . .	59
1.15	Plot Measured Values for Student A versus True Morale . . . . .	62
1.16	Plot the modeled relationship between days and morale . . . . .	63
1.17	Morale over time (difference from true function) . . . . .	67
1.18	Plot Measured Values for Students A and B versus True Morale . . . . .	68
1.19	Morale over time (difference between estimates) . . . . .	69
1.20	Morale over Time . . . . .	71
1.21	Morale over Time . . . . .	73
1.22	Run the Polynomial Plot for Student A . . . . .	75
1.23	Run the Polynomial Plot for Student B . . . . .	75
1.24	Run the Polynomial Plot for Student C . . . . .	75
1.25	Plot 16 <sup>th</sup> -order Regression Models for Students A and B versus True Morale . . . . .	78
1.26	Illustrating the Bias-Variance Tradeoff . . . . .	78
3.1	Connecting with SSH Keys . . . . .	105
3.2	Access EC2 Dashboard . . . . .	108
3.3	Access Key Pairs in the EC2 Dashboard . . . . .	109
3.4	Import a New Public Key . . . . .	110
3.5	Begin the launch process for a new instance . . . . .	110
3.6	Choose the latest stable Ubuntu Server release as AMI . . . . .	111

3.7	Choose <code>t2.micro</code> for Instance Type . . . . .	111
3.8	Choose the latest stable Ubuntu Server release as AMI . . . . .	112
3.9	Review and launch the new instance . . . . .	113
3.10	Add a key pair to the instance . . . . .	113
3.11	Note the IP address of the new instance . . . . .	114
3.12	SSH Connections . . . . .	115

# List of Tables

1.1	Display the matrix $X$ , now a <b>DataFrame</b> . . . . .	12
1.2	Display the first three columns of $X$ . . . . .	13
1.3	Display the first and third columns of $X$ . . . . .	13
1.4	Display the transpose of $X$ . . . . .	14
1.5	Multiply $A$ times $B$ . . . . .	17
1.6	Display results . . . . .	35
1.7	Display the first $n$ rows of the <b>DataFrame</b> . . . . .	42
1.8	Display a sample of the results . . . . .	55
2.1	Show the <code>.head()</code> of <code>seeds_df</code> . . . . .	87





# Chapter 1

## Introductory Lectures

### 1.1 Vectors and Matrices

#### The Vector

A vector is a multi-dimensional element.

e.g.  $(1)$ ,  $(1, 3, 5, 8)$ ,  $(1, 2, \dots, 50)$

As opposed to a scalar

e.g.  $1$ ,  
 $\sqrt{2}$ ,  $e$ ,  $\pi$ ,  $101010$ ,  $47$

#### Scalar Arithmetic

- add
- subtract
- multiply
- divide

#### Representing Data with Matrices

We will use  $n$  to represent the number of distinct data points, or observations, in our data set. We will let  $p$  denote the number of variables that are available for use in making predictions.

In general, we will let  $x_{ij}$  represent the value of the  $j^{th}$  variable for the  $i^{th}$  observation, where  $i \in [1 : n]$  and  $j \in [1 : p]$ .

$i$  will be used to index the samples or observations (from 1 to  $n$ ) and  $j$  will be used to index the variables (from 1 to  $p$ ). We let  $X$  denote a  $n \times p$  matrix whose  $(i, j)^{th}$  element is  $x_{ij}$ .

$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}$$

It is useful to visualize  $X$  as a spreadsheet of numbers with  $n$  rows and  $p$  columns.

## Representing a Matrix in Python

In Python, we can represent a matrix using the `numpy` library and its `array` class.

Listing 1.1: Import `numpy`

```
In [1]: import numpy as np
```

Listing 1.2: Generate a list of random numbers

```
In [2]: import random

n = 5
p = 3
list_of_numbers = []
for _ in range(15):
    list_of_numbers.append(random.randint(1, 100))
```

Listing 1.3: Display the list of random numbers

```
In [3]: list_of_numbers

Out[3]: [37, 68, 61, 10, 10, 5, 67, 69, 97, 30, 67, 60, 60, 96,
        32]
```

Here, we turn our list of numbers into an `np.array` object, then reshape that object into an  $n \times p$  matrix.

Listing 1.4: Reshape the list of random numbers into an  $n \times p$  matrix

```
In [4]: X = np.array(list_of_numbers).reshape((n, p))
```

Listing 1.5: Display the new matrix,  $X$

```
In [5]: X
Out[5]: array([[37, 68, 61],
               [10, 10,  5],
               [67, 69, 97],
               [30, 67, 60],
               [60, 96, 32]])
```

Listing 1.6: Display the shape of the new matrix,  $X$

```
In [6]: X.shape
Out[6]: (5, 3)
```

## The Elements of $X$

$X$  is comprised of two different kinds of vectors:

- Each row of  $X$  is a vector of length  $p$  (contains  $p$  elements).
- Each row of  $X$  is an observation representing a single sample or point from our data set
- Each column of  $X$  is a vector of length  $n$  (contains  $n$  elements).
- Each column of  $X$  is all of the observations for a given variable over our entire data set

All of the rows of  $X$  are the elements of  $X$  considered as row vectors. This is probably how you are most comfortable thinking about  $X$  i.e. some data set.

At the same time, all of the column of  $X$  are the elements of  $X$  considered as column vectors.

Suppose we consider a column of  $X$  as

$$\mathbf{x}_{\cdot j} = \begin{pmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{pmatrix}$$

Note that this is taking into account that we always think of a single vector as a column.

We can then think of  $X$  as

$$X = (\mathbf{x}_{\cdot 1} \quad \mathbf{x}_{\cdot 2} \quad \dots \quad \mathbf{x}_{\cdot p})$$

## Column-Oriented Matrices in Python

In Python, it is easier to work with the columns of a matrix using the Pandas library, and its `DataFrame` class.

Listing 1.7: Import pandas

```
In [7]: import pandas as pd
```

Listing 1.8: Convert  $X$  into a `DataFrame`

```
In [8]: X = pd.DataFrame(X)
```

Listing 1.9: Display the matrix  $X$ , now a `DataFrame`

```
In [9]: X
```

Table 1.1: Display the matrix  $X$ , now a `DataFrame`

	0	1	2
0	37	68	61
1	10	10	5
2	67	69	97
3	30	67	60
4	60	96	32

Listing 1.10: Display the first column of  $X$ 

```
In [10]: X[0]

Out[10]: 0    37
          1    10
          2    67
          3    30
          4    60
          Name: 0, dtype: int64
```

Listing 1.11: Display the first three columns of  $X$ 

```
In [11]: X[[0,1,2]]
```

Table 1.2: Display the first three columns of  $X$ 

	0	1	2
0	37	68	61
1	10	10	5
2	67	69	97
3	30	67	60
4	60	96	32

Listing 1.12: Display the first and third columns of  $X$ 

```
In [12]: X[[0,2]]
```

Table 1.3: Display the first and third columns of  $X$ 

	0	2
0	37	61
1	10	5
2	67	97
3	30	60
4	60	32

## Referencing a Row in a DataFrame

It may still be necessary to reference the row of a matrix represented as a **DataFrame**. This can be done using the `.loc()` method. With `.loc()` rows can be referenced by their row index.

Listing 1.13: Display the second row of  $X$

```
In [13]: X.loc[2]

Out[13]: 0      67
         1      69
         2      97
         Name: 2, dtype: int64
```

Note that even here, the row vector is represented as a column of numbers.

## The Transpose

We obtain the transpose of a matrix by flipping it around its axis.

So that

$$X^T = \begin{pmatrix} x_{11} & x_{21} & \dots & x_{n1} \\ x_{12} & x_{22} & \dots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1p} & x_{2p} & \dots & x_{np} \end{pmatrix}$$

Listing 1.14: Display the transpose of  $X$

```
In [14]: X.T
```

Table 1.4: Display the transpose of  $X$

	0	1	2	3	4
0	37	10	67	30	60
1	68	10	69	67	96
2	61	5	97	60	32

## The Dot Product

It is possible to multiply two vectors of equal length together. This operation is called the **dot product** and yields a scalar value, that is, a single number.

This number is obtained by

1. multiplying each of the corresponding values from the two vectors
2. summing the results

Written in Python, the process would look as follows. Note that we are making use of a convention of writing vectors as a repeated lowercase letter.

Listing 1.15: Prepare a list of the multiples of two vectors

```
In [15]: aa = np.array([1,2,3])
        bb = np.array([-2,-1,0])

        multiples = []

        for a, b in zip(aa, bb):
            multiples.append(a*b)
```

Listing 1.16: Display the list `multiples`

```
In [16]: multiples

Out[16]: [-2, -2, 0]
```

Listing 1.17: Calculate the sum of the list `multiples` to perform the dot product

```
In [17]: dot_product = sum(multiples)
```

Listing 1.18: Display the dot product

```
In [18]: dot_product

Out[18]: -4
```

## The Dot Product as a Row Vector Times a Column Vector

Strictly speaking, the dot product is the multiplication of a row vector times a column vector. So that if we have

$$\mathbf{a} = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \text{ and } \mathbf{b} = \begin{pmatrix} -2 \\ -1 \\ 0 \end{pmatrix}$$

then the dot product we just performed would be written as

$$\mathbf{a}^T \mathbf{b} = (1 \ 2 \ 3) \begin{pmatrix} -2 \\ -1 \\ 0 \end{pmatrix} = 1 \cdot -2 + 2 \cdot -1 + 3 \cdot 0 = -4$$

or

Listing 1.19: Calculate the dot product using `.dot()`

```
In [19]: aa.dot(bb)

Out [19]: -4
```

## Matrix Multiplication

At the heart of the vast majority of work that we will be doing is matrix multiplication. Let us consider two matrices,  $A$  and  $B$ .

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \text{ and } B = \begin{pmatrix} 5 & 6 \\ 7 & 8 \end{pmatrix}$$

Then the product of  $A$  and  $B$  is denoted  $AB$ .

The  $(i, j)^{th}$  element of  $AB$  is computed by computing the dot product of the  $i$ th row of  $A$  with the  $j$ th row of  $B$ .

Listing 1.20: Create two matrices,  $A$  and  $B$ , as `DataFrames`

```
In [20]: A = pd.DataFrame([[1, 2], [3, 4]])
         B = pd.DataFrame([[5, 6], [7, 8]])
```

So the (1,1)th element of  $AB$  should be  $1 \cdot 5 + 2 \cdot 7 = 19$ .



Listing 1.21: Perform the dot product of the first row of  $A$  with the first column of  $B$

```
In [21]: A.loc[0].dot(B[0])
Out[21]: 19
```

We would repeat this process for every element of  $AB$ . As you can imagine this would be very tedious by hand.

Listing 1.22: Multiply  $A$  times  $B$

```
In [22]: A.dot(B)
```

Table 1.5: Multiply  $A$  times  $B$

	0	1
0	19	22
1	43	50

## The importance of Matrices

### Transformations

- Rotation Matrix (project data from two axes to one)
- Compute the difference vector

## 1.2 Statistical Learning

You may be more used to performing calculations where you have a known formula and you wish to use some **input** to calculate an **output**.

For example, you might be calculating the force of an object

$$f = ma$$

Listing 1.23: Define a function for calculating force

```
In [1]: def force(mass, acceleration):  
        return mass * acceleration
```

Listing 1.24: Calculate the force on an object given mass and acceleration

```
In [2]: m = 10  
        a = 9.8  
  
        force(m,a)  
  
Out[2]: 98.0
```

Statistical learning is different. Here, we are given the inputs and outputs:

- $X$  / input / features
- $y$  / output / target

We make an assumption that there is some function  $f$  such that

$$y = f(X) + \varepsilon$$

Note:

1.  $X$  is often a vector
2.  $\varepsilon$  is the "error term" and represents both known and unknown sources of noise

From a known  $X$  and  $y$  we seek the formula or function,  $f$ , that best describes the relationship between them.

This can be trivial.

Listing 1.25: Core imports for numerical computing in Python

```
In [3]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import seaborn as sns

%matplotlib inline
```

Listing 1.26: Define  $x$  and  $y$

```
In [4]: x = [-1,3]
y = [3,5]
```

Here, we have two  $x$  values and two  $y$  values. It is trivial to find a formula describing the relationship between  $x$  and  $y$ .

Listing 1.27: Plot the line defining the relationship between  $x$  and  $y$

```
In [5]: fig, ax = plt.subplots(1,1,figsize=(20,5))

ax.plot(x[0], y[0], 'g^', ms=25)
ax.plot(x[1], y[1], 'g^', ms=25)
ax.plot((x[0], y[0]), (x[1], y[1]))
plt.show()
```

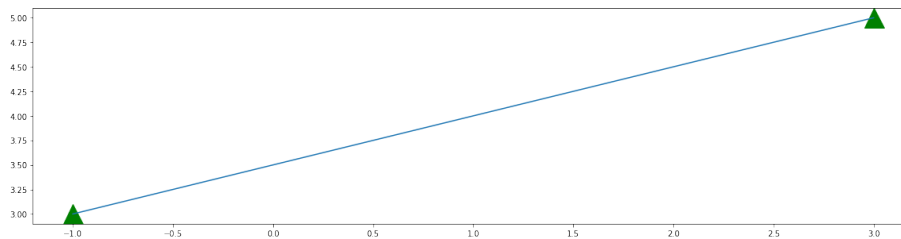


Figure 1.1: Plot the line defining the relationship between  $x$  and  $y$

This is a perfectly determined system. We have two variables and two data points i.e.  $n = p$ .

The situation is more complicated when  $n \neq p$ .

Listing 1.28: Add more values to  $x$  and  $y$ 

```
In [6]: x += [0,1,2]
        y += [3.25, 4, 4.75]

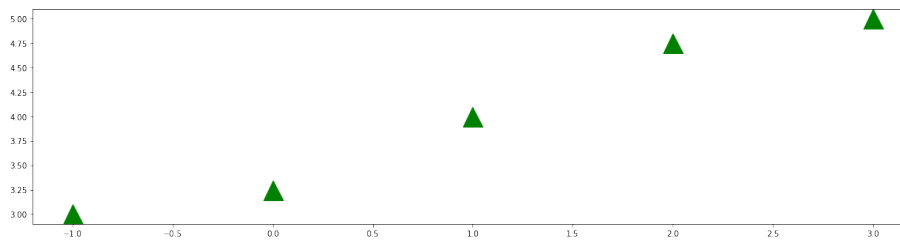
        points = [
            (x[0], y[0]),
            (x[1], y[1]),
            (x[2], y[2]),
            (x[3], y[3]),
            (x[4], y[4])
        ]
```

Listing 1.29: Plot new  $x$  and  $y$ 

```
In [7]: fig, ax = plt.subplots(1,1,figsize=(20,5))

        for point in points:
            ax.plot(*point, 'g^', ms=25)

        plt.show()
```

Figure 1.2: Plot new  $x$  and  $y$ 

The line here does not fit perfectly. The system is **overdetermined** i.e.  $n > p$ .

### The best $f$ aka $\hat{f}$

We seek "the best"  $f$ , called  $\hat{f}$  or "eff hat".

We will quantify what "best" mean in terms of error.

For  $\hat{y} = \hat{f}(X)$ , we consider the error over our data

$$\mathbb{E}[(y - \hat{y})^2] = \left[ f(X) - \hat{f}(X) \right]^2 + \text{Var}(\varepsilon)$$

## Reasons to estimate $f$

1. Prediction - Given some new value for  $X$ , what would we expect  $y$  to be?
2. Inference - Which features ( $X_1, X_2, \dots$ ) are associated with the target  $y$ ?
3. Inference - What is the relationship between the target  $y$  and each feature ( $X_1, X_2, \dots$ )?
4. Inference - Can the relationship between the target  $y$  and each feature be adequately summarized using a linear equation, or is the relationship more complicated?

Listing 1.30: Load Advertising data as a DataFrame

```
In [8]: Advertising = pd.read_csv('data/Advertising.csv',
                                index_col=0)
```

## Inference on Advertising Data

Here we create a series of **linear models** to help us make inferences. Linear models allow for relatively simple and interpretable inference, but may not yield as accurate predictions as some other approaches

Listing 1.31: Plot the Advertising Data

```
In [9]: fig, ax = plt.subplots(1,3,figsize=(20,5))
        advertising_types = ['TV', 'radio', 'newspaper']
        for i, feature in enumerate(advertising_types):
            sns.regplot(Advertising[feature],
                        Advertising['sales'], ax=ax[i])

        plt.show()
```

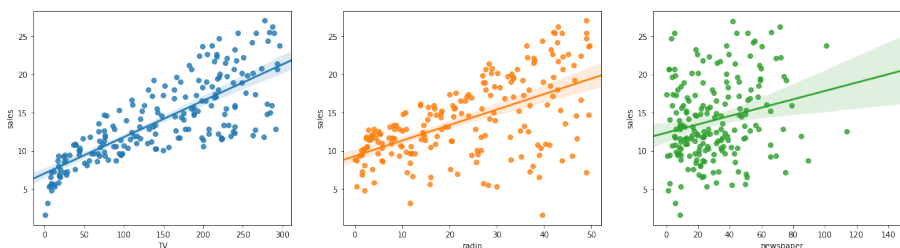


Figure 1.3: Plot the Advertising Data

- Which media contribute to sales?
- Which media generate the biggest boost in sales?
- How much increase in sales is associated with a given increase in TV advertising?

## Parametric versus Non-Parametric

### Parametric Models

1. We make an assumption about the functional form, or shape, of  $f$ .
2. After a model has been selected, we need a procedure that uses the training data to fit or train the model.

The linear model is an important example of a parametric model. We consider

$$\hat{f}(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p$$

- A linear model is specified in terms of  $p + 1$  parameters  $\beta_0, \beta_1, \dots, \beta_p$
- We estimate the parameters by fitting the model to training data.
- Although it is almost never correct, a linear model often serves as a good and interpretable approximation to the unknown true function  $f(X)$ .
- the most common method to fitting linear models is the ordinary least squares (OLS) method

### Non-Parametric Models

- Non-parametric methods do not make explicit assumptions about the functional form of  $f$ .
- Any parametric approach brings with it the possibility that the functional form used to estimate  $f$  is very different from the true  $f$ , in which case the resulting model will not fit the data well.
- In contrast, non-parametric approaches completely avoid this danger, since essentially no assumption about the form of  $f$  is made.

- But non-parametric approaches do suffer from a major disadvantage: since they do not reduce the problem of estimating  $f$  to a small number of parameters, a very large number of observations (far more than is typically needed for a parametric approach) is required in order to obtain an accurate estimate for  $f$  ... **THE CURSE OF DIMENSIONALITY**.

K Nearest Neighbors is an important example of a non-parametric model.

## Tradeoffs

### The Trade-Off Between Overfitting and Underfitting

A potential disadvantage of a parametric approach is that the model we choose will usually not match the true unknown form of  $f$ .

If the chosen model is too far from the true  $f$ , then our estimate will be poor.

We can try to address this problem by choosing flexible models that can fit many different possible functional forms for  $f$ .

Listing 1.32: Define correlated data with random noise

```
In [10]: xx = np.linspace(-2,2,20)
        yy = xx**3 + 2*np.random.random(20) + 100
```

Listing 1.33: Import the Linear Regression model from Scikit-Learn

```
In [11]: from sklearn.linear_model import LinearRegression
```

Listing 1.34: Instantiate an instance of the Linear Regression model

```
In [12]: lr = LinearRegression()
```

Listing 1.35: Fit the new Linear Regression model

```
In [13]: lr.fit(xx.reshape(-1, 1), yy)

Out[13]: LinearRegression(copy_X=True, fit_intercept=True,
                          n_jobs=1, normalize=False)
```

Listing 1.36: Display the coefficients and intercept of the fit model

```
In [14]: lr.coef_, lr.intercept_

Out[14]: (array([ 2.91264981]), 101.030290844449574)
```

Listing 1.37: Plot the line of best fit

```
In [15]: fig, ax = plt.subplots(1,1,figsize=(20,5))
        ax.scatter(xx, yy)
        ax.plot(xx, xx*lr.coef_ + lr.intercept_)

        plt.show()
```

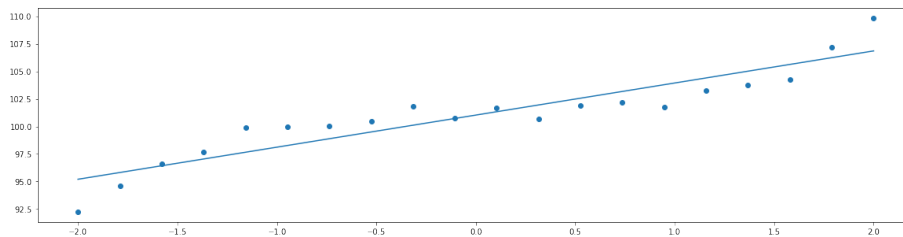


Figure 1.4: Plot the line of best fit

Listing 1.38: Import a helper function

```
In [16]: from lib.lec_3_helper import fit_and_plot_poly_to_degree
```

Listing 1.39: Plot lines of best fit for higher order functions

```
In [17]: fig, ax = plt.subplots(1,1,figsize=(20,5))
        ax.scatter(xx, yy)
        ax.plot(xx, xx*lr.coef_ + lr.intercept_)
        for i in [2,3,7,20]:
            fit_and_plot_poly_to_degree(xx, yy, i, ax)
        plt.legend()

        plt.show()
```

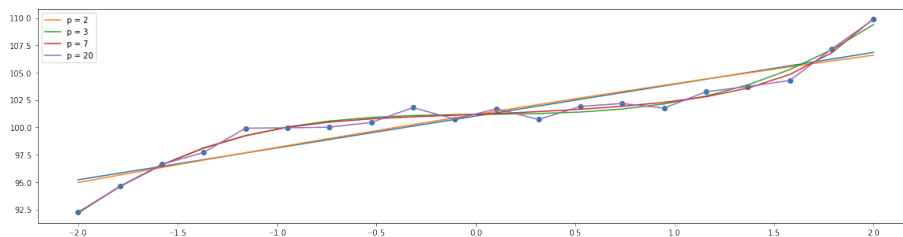


Figure 1.5: Plot lines of best fit for higher order functions



Listing 1.40: Instantiate and Fit Linear Regression on Sales using TV data

```
In [18]: lr = LinearRegression()
         lr.fit(Advertising[['TV']], Advertising.sales)
         lr.coef_, lr.intercept_

Out[18]: (array([ 0.04753664]), 7.0325935491276947)
```

Listing 1.41: Plot line of best fit

```
In [19]: fig, ax = plt.subplots(1,1,figsize=(20,5))

         predictions = lr.predict(Advertising[['TV']])

         ax.scatter(Advertising.TV, Advertising.sales,
                    label='actual')
         ax.scatter(Advertising.TV, predictions,
                    label='prediction')
         plt.legend()

         plt.show()
```

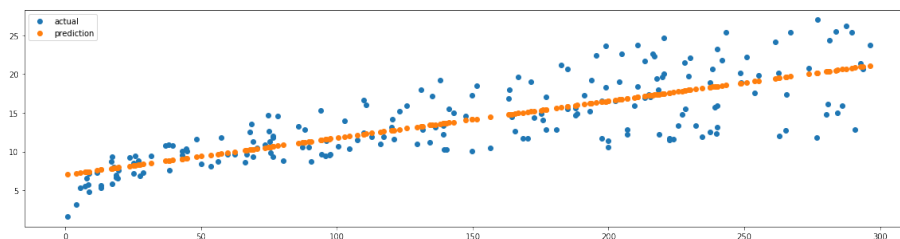


Figure 1.6: Plot line of best fit

Listing 1.42: Plot higher order lines of best fit

```

In [20]: fig, ax = plt.subplots(1,1,figsize=(20,5))

        for i in [1,2,3,20,30,50]:
            fit_and_plot_poly_to_degree(
                Advertising.TV,
                Advertising.sales,
                i, ax, scatter=True)

        ax.scatter(Advertising.TV, Advertising.sales,
                    label='actual', s=100, alpha=0.6)

        plt.legend()

        plt.show()

```

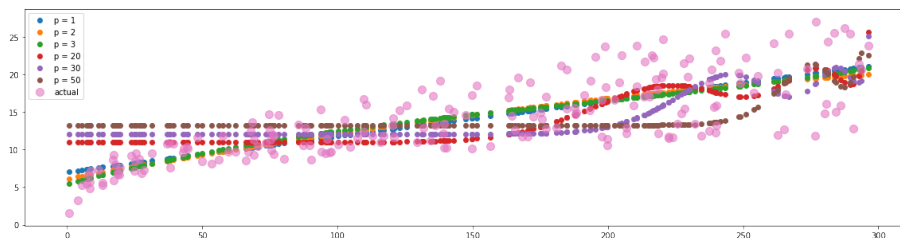


Figure 1.7: Plot higher order lines of best fit

But in general, fitting a more flexible model requires estimating a greater number of parameters.

These more complex models can lead to a phenomenon known as **overfitting** the data, which essentially means they follow the errors, or noise, too closely.

We can generalize this in terms of model flexibility. In general, less flexible models may overfit, whereas more flexible models may overfit.

### The Trade-Off Between Prediction Accuracy and Model Interpretability

We can also think about flexibility and accuracy in terms of our ability to interpret models.

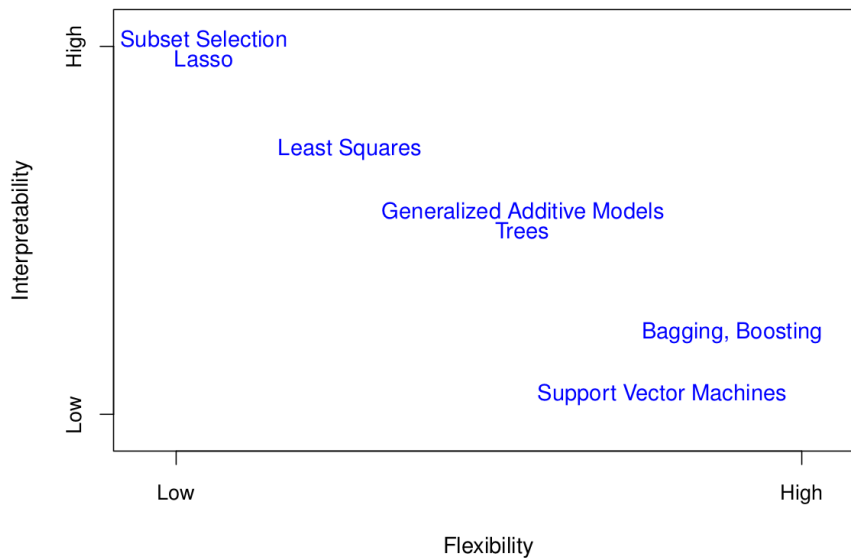


Figure 1.8: Model Flexibility

Models on the right can often be "black box" models that are difficult to interpret. Their increased flexibility may make them stronger at making accurate predictions.

Models on the left will be much more simple, but they may also be less accurate.

## Statistical Learning Domains

1. Supervised Learning
2. Unsupervised Learning
3. Regression
4. Classification
5. Cluster Analysis

Why is it necessary to introduce so many different statistical learning approaches, rather than just a single best method?

**TANSTAAFL**

## 1.3 The Bias-Variance Tradeoff

Listing 1.43: Set up Numerical Python environment

```
In [1]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import seaborn as sns

%matplotlib inline
plt.rc('figure', figsize=(20, 6))
```

Listing 1.44: Load Iris dataset

```
In [2]: from sklearn.datasets import load_iris

feature_names = [
    'sepal_length',
    'sepal_width',
    'petal_length',
    'petal_width'
]

IRIS = load_iris()
iris_df = pd.DataFrame(IRIS.data, columns=feature_names)
labels = IRIS.target_names
```

In his 1996 paper, "The Lack of A Priori Distinctions between Learning Algorithms," David Wolpert asserts that "There is No Free Lunch in Machine Learning", essentially saying that the performance of on machine learning models are equivalent when averaged across all possible problems. In other words, there is no way to know which model is going to perform at best for a particular problem. In practice this means that we must have strong methods for assessing the performance of our models.

Thinking about model performance is complex. We cannot simply choose the model that performs best with the data that we have. The reason for this is that the data we have represents a sample of the actual data that we could possibly collect now or in the future. The "best model" is not necessarily one that performs best on the data that we have. The best model is a model that performs extremely well on the data that we have but it's also capable of generalizing to new data. In statistical learning, we have a framework for thinking about this called The Bias-Variance Tradeoff.

This framework is difficult to understand. Adding to the difficulty is the fact that both "bias" and "variance" are important concepts for working in applied

statistics **and** the meaning of these terms is difficult to reconcile with their meaning when thinking about the Bias-Variance Tradeoff.

Let us try to come to a high level understanding of these terms in this setting before looking at them in application. You might think of **bias** as the extent to which a particular model can learn to represent an underlying physical phenomenon. A low bias means that the model was able to learn the phenomenon well. For example, looking at the Iris data set, if we are building a regression model to predict the petal width then bias would be the degree to which our model was able to learn the relationship between petal length and petal width or sepal width and petal width.

**Variance** on the other hand is the extent to which a particular model would change if fit with different data. We previously looked at sampling our data set. We saw the measured means of value change with each sample. From this we can infer that a model predicting the meaning with only a few points of data has a high variance.

## The Gory Details

If we have split our data into a training set and a testing set, then we can think of choosing the best model in terms of optimizing the expected test error,  $MSE_{test}$

Let's consider sources of possible test error:

$$MSE_{test} = \mathbb{E} [(y - \hat{y})^2] = \text{Var}(\hat{y}) + (\text{Bias}(\hat{y}))^2 + \text{Var}(\epsilon)$$

This is intended to be a conceptual and not an actual calculation to be performed. Let's think about what each of these terms might represent. The variance is error introduced to the model by the specific choice of training data. Of course this isn't something that we choose, at least not without using randomness, but the training data that is used will impact the model. By nature, variance is a squared value and that's always positive. Bias is introduced by choosing a specific model. Note that it is squared here and thus also always positive. The last term is the variance caused by noise in the system. We have no way of controlling this, nor of actually knowing what is truly noise and what is model variance or bias. Again this term is always positive.

The important thing is that all three of these terms are always positive. The impact of this is that one kind of error cannot be offset by another kind of error. A high variance cannot be offset by a low bias. **In order to choose the best model, we are going to need to simultaneously minimize both bias and variance.** The problem is changing an aspect of our model to decrease one will typically increase the other – The Bias-Variance Tradeoff.

## Model Assessment

Consider these eight models for predicting petal width  $x_{pw}$  on the Iris data set, using the other three features, petal length  $x_{pl}$ , sepal length  $x_{sl}$ , and sepal width  $x_{sw}$

$$\begin{aligned}\hat{f}_{.0} &= \hat{x}_{pw} = \beta_{.0} \\ \hat{f}_{.1} &= \hat{x}_{pw} = \beta_{.0} + \beta_{.1}x_{pl} \\ \hat{f}_{.2} &= \hat{x}_{pw} = \beta_{.0} + \beta_{.1}x_{sl} \\ \hat{f}_{.3} &= \hat{x}_{pw} = \beta_{.0} + \beta_{.1}x_{sw} \\ \hat{f}_{.4} &= \hat{x}_{pw} = \beta_{.0} + \beta_{.1}x_{pl} + \beta_{.1}x_{sl} \\ \hat{f}_{.5} &= \hat{x}_{pw} = \beta_{.0} + \beta_{.1}x_{pl} + \beta_{.1}x_{sw} \\ \hat{f}_{.6} &= \hat{x}_{pw} = \beta_{.0} + \beta_{.1}x_{sl} + \beta_{.1}x_{sw} \\ \hat{f}_{.7} &= \hat{x}_{pw} = \beta_{.0} + \beta_{.1}x_{pl} + \beta_{.1}x_{sl} + \beta_{.1}x_{sw}\end{aligned}$$

We will split our dataset into ten randomly chosen subsets of 50 points each. We will assess the performance of each of these models. We will use the mean of the performance to represent bias and the standard deviation of the performance to represent variance.

Listing 1.45: Create 10 sample datasets from full data set

```
In [3]: samples = []
        for _ in range(10):
            samples.append(iris_df.sample(50))
```

Listing 1.46: Show length of each sample in the list `samples`

```
In [4]: [len(sample_set) for sample_set in samples]

Out[4]: [50, 50, 50, 50, 50, 50, 50, 50, 50, 50]
```

Listing 1.47: Import the design matrices library, `dmatrices` from `patsy`

```
In [5]: from patsy import dmatrices
```

Listing 1.48: Import the Linear Regression Model from Sklearn

```
In [6]: from sklearn.linear_model import LinearRegression
```

Listing 1.49: Simple example of computing a list of squared differences

```
In [7]: a = np.array((1,2,3))  
        b = np.array((4,5,6))  
  
        (b - a)**2  
  
Out[7]: array([9, 9, 9])
```

Listing 1.50: Define a Mean Squared Error function

```
In [8]: def MSE(actual, predicted):  
        return sum((actual - predicted)**2)/len(actual)
```

Listing 1.51: Define a function for running our tests

```

In [9]: def test_func(model_description):

    test = dict()

    test['samples'] = [
        dmatrices(model_description, sample)
        for sample in samples
    ]

    test['models'] = [
        LinearRegression(fit_intercept=False)
        for _ in range(10)
    ]

    test['scores'] = []

    models_and_samples = zip(test['models'],
                              test['samples'])

    for model, sample in models_and_samples:
        target = sample[0]
        features = sample[1]
        model.fit(features, target)

        mse = MSE(model.predict(features), target)

        test['scores'].append(mse)

    test['scores'] = np.array(test['scores'])

    results = { 'description' : model_description }

    results['bias'] = test['scores'].mean()
    results['variance'] = test['scores'].std()

    return test, results

```

Listing 1.52: Get the results of running the test on the minimal model

```

In [10]: model_1 = "petal_width ~ 1"
         test_1, results_1 = test_func(model_1)

```



Listing 1.53: Display the results of the MSE scores on the minimal model

```
In [11]: test_1['scores']

Out[11]: array([[ 0.554964],
                 [ 0.4993   ],
                 [ 0.527264],
                 [ 0.575156],
                 [ 0.542644],
                 [ 0.627796],
                 [ 0.502564],
                 [ 0.604    ],
                 [ 0.548164],
                 [ 0.500084]])
```

Listing 1.54: Show how we are collecting results

```
In [12]: test_1['scores'].mean(), test_1['scores'].std()

Out[12]: (0.54819360000000006, 0.041808768121531649)
```

Listing 1.55: Display the results returned by the test

```
In [13]: results_1

Out[13]: {'description': 'petal_width ~ 1',
          'bias': 0.54819360000000006,
          'variance': 0.041808768121531649}
```

Listing 1.56: Define all models

```
In [14]: model_2 = "petal_width ~ 1 + petal_length"
         model_3 = "petal_width ~ 1 + sepal_length"
         model_4 = "petal_width ~ 1 + sepal_width"
         model_5 = "petal_width ~ 1 + petal_length + sepal_width"
         model_6 = "petal_width ~ 1 + petal_length + sepal_length"
         model_7 = "petal_width ~ 1 + sepal_length + sepal_width"
         model_8 = "petal_width ~ 1 + petal_length + sepal_length
         + sepal_width"
```

Listing 1.57: Execute the test on the other seven model and collect the results

```
In [15]: test_2, results_2 = test_func(model_2)
         test_3, results_3 = test_func(model_3)
         test_4, results_4 = test_func(model_4)
         test_5, results_5 = test_func(model_5)
         test_6, results_6 = test_func(model_6)
         test_7, results_7 = test_func(model_7)
         test_8, results_8 = test_func(model_8)
         results = [
             results_1,
             results_2,
             results_3,
             results_4,
             results_5,
             results_6,
             results_7,
             results_8
         ]
```

Listing 1.58: Turn the results into a `DataFrame` for ease of display

```
In [16]: results = pd.DataFrame(results)
```

Listing 1.59: Add colors to results

```
In [17]: results['color'] = ['red', 'orange', 'yellow', 'green',
                             'blue', 'cyan', 'purple', 'black']
```

Listing 1.60: Display results

```
In [18]: results
```

Table 1.6: Display results

	bias	description	variance	color
0	0.548194	petal_width ~ 1	0.041809	red
1	0.043638	petal_width ~ 1 + petal_length	0.007811	orange
2	0.192337	petal_width ~ 1 + sepal_length	0.028752	yellow
3	0.475903	petal_width ~ 1 + sepal_width	0.056484	green
4	0.041390	petal_width ~ 1 + petal_length + sepal_width	0.008234	blue
5	0.041939	petal_width ~ 1 + petal_length + sepal_length	0.007304	cyan
6	0.144526	petal_width ~ 1 + sepal_length + sepal_width	0.022443	purple
7	0.035526	petal_width ~ 1 + petal_length + sepal_length ...	0.005027	black

Listing 1.61: Plot the bias versus the variance

```

In [19]: plt.figure(figsize=(10,5))
         for i in results.index:
             plt.scatter(results.loc[i].bias,
                 results.loc[i].variance, c=results.loc[i].color,
                 label=results.loc[i].description)
         plt.xlabel('Bias')
         plt.ylabel('Variance')
         plt.legend()

         plt.show()

```

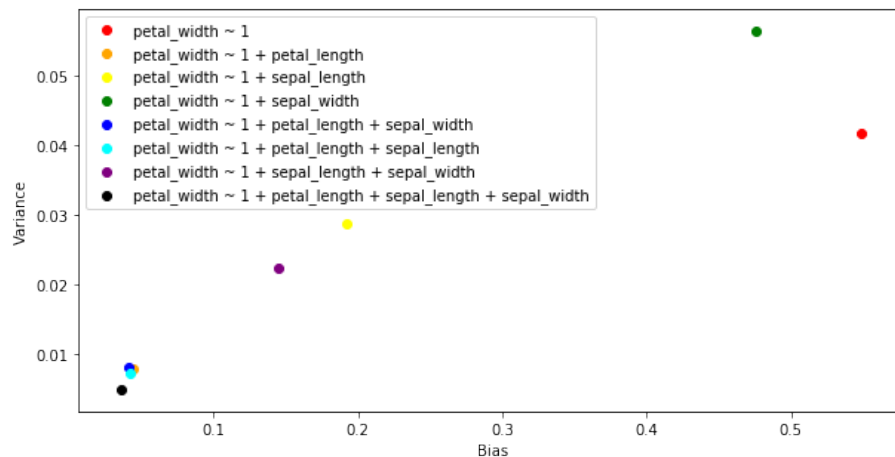


Figure 1.9: Plot the bias versus the variance

## 1.4 Probability

We are all familiar with the phrase “the probability that a coin will land heads is 0.5”. But what does this mean? There are actually at least two different interpretations of probability. One is called the **frequentist** interpretation. In this view, probabilities represent long run frequencies of events. For example, the above statement means that, if we flip the coin many times, we expect it to land heads about half the time.

The other interpretation is called the **Bayesian** interpretation of probability. In this view, probability is used to quantify our uncertainty about something; hence it is fundamentally related to information rather than repeated trials. In the Bayesian view, the above statement means we believe the coin is equally likely to land heads or tails on the next toss.

One big advantage of the Bayesian interpretation is that it can be used to model our uncertainty about events that do not have long term frequencies. For example, we might want to compute the probability that the polar ice cap will melt by 2020 CE. This event will happen zero or one times, but cannot happen repeatedly. Nevertheless, we ought to be able to quantify our uncertainty about this event; based on how probable we think this event is, we will (hopefully!) take appropriate actions. To give some more machine learning oriented examples, we might have received a specific email message, and want to compute the probability it is spam. Or we might have observed a “blip” on our radar screen, and want to compute the probability distribution over the location of the corresponding target (be it a bird, plane, or missile). In all these cases, the idea of repeated trials does not make sense, but the Bayesian interpretation is valid and indeed quite natural.

The basic rules of probability theory are the same, no matter which interpretation is adopted.

### Basic Probability

The expression  $p(A)$  denotes the probability that the event  $A$  is true. For example,  $A$  might be the logical expression “it will rain tomorrow”.

We have:

- $0 \leq p(A) \leq 1$
- $p(A) = 0$  means the event definitely will not happen
- $p(A) = 1$  means the event definitely will happen

- $p(\neg A)$  denotes the probability of the event not  $A$ , that is that  $A$  will not occur
- $p(\neg A) = 1 - p(A)$
- We will often write  $A = 1$  to mean the event  $A$  is true, and  $A = 0$  to mean the event  $A$  is false

## Discrete random variables

A **discrete random variable**  $X$  is a set of possible observed events. For example, we might have that  $X$  is the integer age of the students in our class.

We can intuit that certainly  $X \in [0, 100]$  ( $X$  is *in* the set of integers from 0 to 100). We might take a sample from  $X$  and this will signify the age of one member of the class. In terms of probability, we might think of the event  $P(X = x)$ , the probability that our sample is some number  $x$ . We can also call this simply  $p(x)$ . Assuming that no one in the class has the same integer age, we have an equal chance of sampling every student, and there are  $n$  students in the class, we could say  $p(x) = \frac{1}{n}$ . As with all probability,  $0 \leq p(x) \leq 1$ .

**ADVANCED NOTE**  $p(x)$  is called a **probability mass function** or pmf.

## Probability of Two Events Occurring

Given two events,  $A$  and  $B$ , we define the probability of  $A$  or  $B$  as follows:

$$p(A \vee B) = p(A) + p(B) - p(A \wedge B)$$

$$p(A \vee B) = p(A) + p(B) \text{ if } A \text{ and } B \text{ are mutually exclusive}$$

## Joint Probability

Joint probability refers to two events co-occurring.

$$p(A, B) = p(A \wedge B) = p(A|B)p(B) = p(B|A)p(A)$$

This is sometimes called **the product rule**. Note that in both  $p(A|B)p(B)$  and  $p(B|A)p(A)$ , *both events are occurring*. You should read  $p(A|B)p(B)$  as "the probability of  $A$  given  $B$  times the probability of  $B$ ".

## Conditional Probability

$$p(A|B) = \frac{p(A, B)}{p(B)} \text{ if } p(B) > 0$$

## Bayes Rule

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)} \text{ if } p(B) > 0$$

### An Example: A Cancer Detection Test

Suppose a medical institution has developed a test for assessing whether or not a patient has cancer. The test has been around for a long time (meaning we can use frequentist statistics to measure its success) and we know that it is 98% successful in identifying cancer when a patient has cancer and 99% successful in returning a negative when a patient does not have cancer. We can rewrite each of these as a conditional probability:

$$\begin{aligned} p(\text{positive test}|\text{cancer}) &= 0.99 \\ p(\text{negative test}|\text{no cancer}) &= 0.97 \end{aligned}$$

We also know that cancer in the American population is extremely rare. Approximately 0.4% of people develop cancer. We can thus say  $p(\text{cancer}) = 0.004$ . *NOTE: these numbers are made up for demonstration.*

### Probability that a patient has cancer

What we wish to know is, given a positive test, what is the probability that the patient has cancer? What is  $p(\text{cancer}|\text{positive test})$ ?

### Bayes Rule

We can find this probability by calculating

$$p(\text{cancer}|\text{positive test}) = \frac{p(\text{positive test}|\text{cancer})p(\text{cancer})}{p(\text{positive test})}$$

The calculation is pretty straightforward, except for the calculation of  $p(\text{positive test})$ . This calculation must include all of the ways in which we can obtain a positive test. We have to include the false positives in the calculation. The false positive rate is  $1 - p(\text{negative test}|\text{no cancer}) = 0.03$

$$\begin{aligned} p(\text{positive test}) &= p(\text{positive test}|\text{cancer})p(\text{cancer}) + p(\text{positive test}|\text{no cancer})p(\text{no cancer}) \\ &= 0.99 \cdot 0.004 + 0.03 \cdot 0.999 \\ &= 0.03384 \end{aligned}$$

Then,

$$\begin{aligned} p(\text{cancer}|\text{positive test}) &= \frac{p(\text{positive test}|\text{cancer})p(\text{cancer})}{p(\text{positive test})} \\ &= \frac{0.99 \cdot 0.004}{0.03384} \\ &= 0.11702 \end{aligned}$$

### The Counter-intuition of False Positives

Below we visualize a population of 1000 patients to whom this test has been administered. In 1000 patients, we would expect 996 of them to be cancer-free. But according to the test, with 996 patients, we would expect 30 **false positives**. In the same population, we would expect 4 patients to actually have cancer. Luckily, we would expect the test to correctly identify all four of these patients. This would be a total of 34 positive tests, the vast majority of these being false positives.

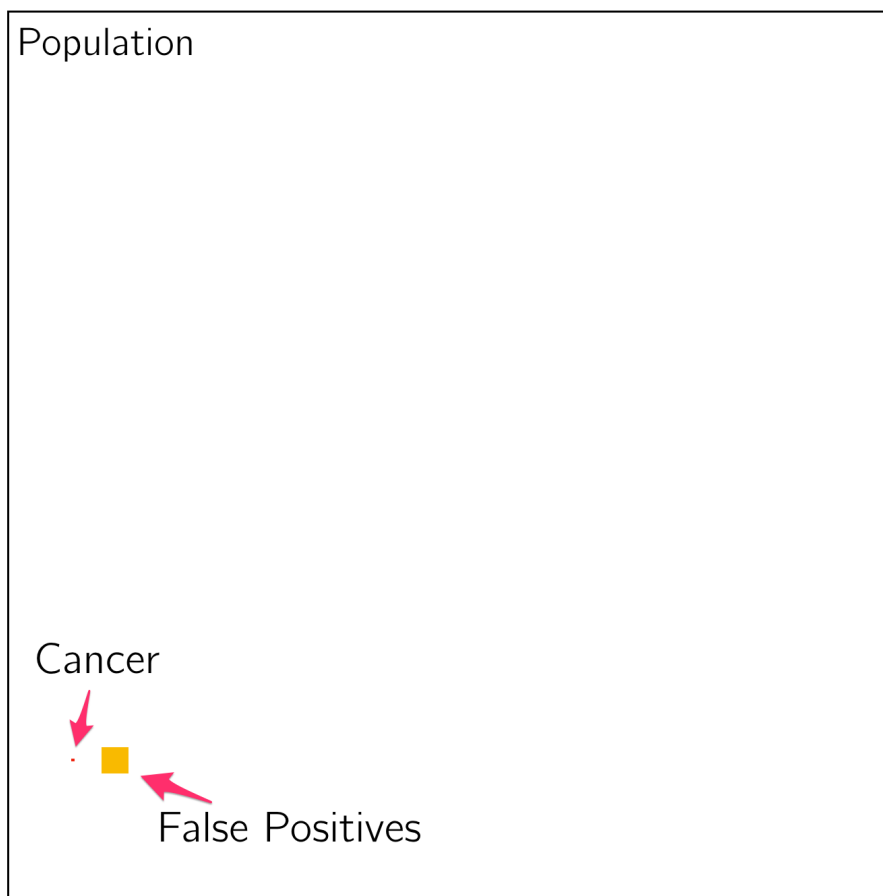


Figure 1.10: A population of 1000 patients

We might also look at these results using a **confusion matrix**

	— True Positive —	True Negative —	— :: — :: — :: —	Predicted Posi-
tive — 4 —	30 —	Predicted Negative —	0 —	968 —

In this particular case, this result may be preferable to lowering the sensitivity of our test to lower the false positive rate, but at the expense of missing true positives. One can imagine a situation in which the opposite were true so that we would want to lower the false positive rate at the expense of missing true positives. For example, we might consider a test to see if a patient is a match for a certain kind of organ donation. In this case, it is preferable that every positive match is a true positive at the expense of possibly missing one.

## 1.5 Probabilistic Model Selection / In-Sample Model Selection

From a probabilistic perspective we seek the model estimate  $\hat{f}$  that is most probable given the data. Let  $H_f$  be the hypothesis that a specific model is the correct model and  $p(D)$  represent the probability of the data, then we seek  $\hat{f}$  that maximizes, that is we seek

$$\hat{f} = \operatorname{argmax}_f (p(H_f|D))$$

We can use Bayes' Rule to invert this probability so that

$$\hat{f} = \operatorname{argmax}_f \left( \frac{p(D|H_f)p(H_f)}{p(D)} \right)$$

If we consider that each model is equally likely and that the probability of the data  $p(D)$  is a constant applied to every calculation, then without loss of generality, we can say

$$\hat{f} = \operatorname{argmax}_f (p(D|H_f))$$

In the Regression setting, it can be shown that maximizing this equation is equivalent to minimizing the residual sum of squares,

$$\text{RSS} = \sum (y_i - \hat{y})^2 = (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}})$$

## Generating Model Hypotheses

For this particular task we will consider the hypothesis space to be all of the linear models possible given our data sets. If we recognize that the set of all of the models possible with `dataset_1` is a subset of all of the models possible with `dataset_2` and likewise with `dataset_3` and `dataset_4`, then we can take the size of this hypothesis space to be all of the possible models made using `dataset_2` and all of the possible models made using `dataset_4`.

This is a vast space.



### 1.5. PROBABILISTIC MODEL SELECTION / IN-SAMPLE MODEL SELECTION 41

To get a sense of the size of this hypothesis space, let us consider a simpler data set, one that has just two features,  $x_1$  and  $x_2$ . We note that there are four possible models using these two features:

$$\begin{aligned}\hat{f}_1 &= \beta_0 \\ \hat{f}_2 &= \beta_0 + \beta_1 x_1 \\ \hat{f}_3 &= \beta_0 + \beta_1 x_2 \\ \hat{f}_4 &= \beta_0 + \beta_1 x_1 + \beta_2 x_2\end{aligned}$$

Similarly for a data set with three features, there will be eight possible models:

$$\begin{aligned}\hat{f}_1 &= \beta_0 \\ \hat{f}_2 &= \beta_0 + \beta_1 x_1 \\ \hat{f}_3 &= \beta_0 + \beta_1 x_2 \\ \hat{f}_4 &= \beta_0 + \beta_1 x_3 \\ \hat{f}_5 &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 \\ \hat{f}_6 &= \beta_0 + \beta_1 x_1 + \beta_2 x_3 \\ \hat{f}_7 &= \beta_0 + \beta_1 x_2 + \beta_2 x_3 \\ \hat{f}_8 &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3\end{aligned}$$

Essentially we are forming what is known as a Power Set of possible feature combinations. The number of elements in a Power Set is given by  $2^p$  where  $p$  is the number of elements in the set.

For  $p = 2$ , this is four models. For  $p = 3$ , this is eight models.

For  $p = 100$ , this is a hypothesis space With a dimension of  $1.27 \times 10^{30}$ . If we trained one model per second, it would take us  $4.02 \times 10^{22}$  years to search the entire hypothesis space. For perspective, physicists estimate that the universe is approximately  $13.8 \times 10^9$  years old.

In other words, we will rarely be able to exhaustively search a hypothesis space.

Listing 1.62: Load Iris Dataset

```
In [1]: iris.data = read.csv("data/iris.csv", row.names='X')
```

Listing 1.63: Display the first  $n$  rows of the `DataFrame`

```
In [2]: head(iris.data)
```

Table 1.7: Display the first  $n$  rows of the `DataFrame`

	sepal_length	sepal_width	petal_length	petal_width	label
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0
5	5.4	3.9	1.7	0.4	0

Listing 1.64: Install the Package `GGally`

```
In [3]: install.packages("GGally")
```

```
Out[3]: Updating HTML index of packages in '.Library'
        Making 'packages.html' ... done
```

Listing 1.65: Load the library `GGally`

```
In [4]: library(GGally)
```

```
Out[4]: Loading required package: ggplot2
```

Listing 1.66: Display `ggpairs` for the Iris Dataset

```
In [5]: ggpairs(iris.data)
```

Listing 1.67: Perform a Logistic Regression on the Species Label

```
In [6]: iris.glm = glm("label ~ 1 + sepal_length + sepal_width +
                      petal_length + petal_width", data = iris.data)
```

Listing 1.68: Display the `summary` of the Logistic Regression

```

In [7]: summary(iris.glm)

Out[7]:
Call:
glm(formula = "label ~ 1 + sepal_length + sepal_width +
  petal_length + petal_width",
     data = iris.data)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-0.59046  -0.15230   0.01338   0.10332   0.55061

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    0.19208    0.20470   0.938 0.349611
sepal_length  -0.10974    0.05776  -1.900 0.059418 .
sepal_width   -0.04424    0.05996  -0.738 0.461832
petal_length   0.22700    0.05699   3.983 0.000107 ***
petal_width    0.60989    0.09447   6.456 1.52e-09 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1
' ' ' 1

(Dispersion parameter for gaussian family taken to be
0.04798457)

Null deviance: 100.0000  on 149  degrees of freedom
Residual deviance:   6.9578  on 145  degrees of freedom
AIC: -22.935

Number of Fisher Scoring iterations: 2

```

## The Log-Likelihood

Without going too far into the math, we can think of the log-likelihood as a **likelihood function** telling us how likely a model is given the data.

This value is not human interpretable but is useful as a comparison.

Listing 1.69: Calculate the Log-Likelihood of the Logistic Regression

```

In [8]: logLik(iris.glm)

Out[8]: 'log Lik.' 17.46751 (df=6)

```

”All models are wrong, but some are useful.” - George Box

We might be concerned with one additional property - the **complexity** of the model.

### William of Occam

**Occam's razor** is the problem-solving principle that, when presented with competing hypothetical answers to a problem, one should select the one that makes the fewest assumptions.



Figure 1.11: William of Occam

We can represent this idea of complexity in terms of both the number of features we use and the amount of data.

### Bayesian Information Criterion

The BIC is formally defined as

$$\text{BIC} = \ln(n)k - 2\ln(\hat{L}).$$

where

-  $\hat{L}$  = the maximized value of the likelihood function of the model  $M$  -  $x$  = the observed data -  $n$  = the number of data points in  $x$ , the number of observations, or equivalently, the sample size; -  $k$  = the number of parameters estimated by the model. For example, in multiple linear regression, the estimated parameters are the intercept, the  $q$  slope parameters, and the constant variance of the errors; thus,  $k = q + 2$ .

It might help us to think of it as

$$\text{BIC} = \text{complexity} - \text{likelihood}$$

## Akaike Information Criterion

The AIC is formally defined as

$$\text{AIC} = 2k - 2\ln(\hat{L}).$$

where

-  $\hat{L}$  = the maximized value of the likelihood function of the model  $M$  -  $x$  = the observed data -  $n$  = the number of data points in  $x$ , the number of observations, or equivalently, the sample size; -  $k$  = the number of parameters estimated by the model. For example, in multiple linear regression, the estimated parameters are the intercept, the  $q$  slope parameters, and the constant variance of the errors; thus,  $k = q + 2$ .

It might help us to think of it as

$$\text{AIC} = \text{complexity} - \text{likelihood}$$

We can think of - likelihood as bias and complexity as variance.

Listing 1.70: Calculate the BIC using R

```
In [9]: BIC(iris.glm)

Out[9]: [1] -4.871215
```

Listing 1.71: Calculate the BIC manually

```
In [10]: n = length(iris.glm$fitted.values)
         p = length(coefficients(iris.glm))

         likelihood = 2 * logLik(iris.glm)
         complexity = log(n)*(p+1)

         bic = complexity - likelihood
         bic

Out[10]: 'log Lik.' -4.871215 (df=6)
```

Listing 1.72: Define a function to calculate the BIC

```
In [11]: BIC_of_model = function (model) {
         n = length(model$fitted.values)
         p = length(coefficients(model))

         likelihood = 2 * logLik(model)
         complexity = log(n)*(p+1)

         bic = complexity - likelihood
         return(bic)
}
```

Listing 1.73: Calculate the BIC with function

```
In [12]: BIC_of_model(iris.glm)

Out[12]: 'log Lik.' -4.871215 (df=6)
```

## Model Selection

Here, we choose the optimal model by removing features one by one.

Listing 1.74: Define full model

```
In [13]: model_1 = "label ~ 1 + sepal_length + sepal_width +
                 petal_length + petal_width"
```

Listing 1.75: Define four models each with one feature removed

```
In [14]: model_2a = gsub('\\+ petal_width','', model_1)
        model_2b = gsub('\\+ petal_length', '', model_1)
        model_2c = gsub('\\+ sepal_width', '', model_1)
        model_2d = gsub('\\+ sepal_length', '', model_1)
```

model\_2a

Listing 1.76: Perform Logistic Regression on these models

```
In [16]: iris.glm.1 = glm(model_1, data=iris.data)
        iris.glm.2a = glm(model_2a, data=iris.data)
        iris.glm.2b = glm(model_2b, data=iris.data)
        iris.glm.2c = glm(model_2c, data=iris.data)
        iris.glm.2d = glm(model_2d, data=iris.data)
```

Listing 1.77: Display the results of the BIC on these models

```
In [17]: print(c('model_1', BIC_of_model(iris.glm.1)))
        print(c('model_2a', BIC_of_model(iris.glm.2a )))
        print(c('model_2b', BIC_of_model(iris.glm.2b )))
        print(c('model_2c', BIC_of_model(iris.glm.2c )))
        print(c('model_2d', BIC_of_model(iris.glm.2d )))

Out[17]: [1] "model_1"          "-4.87121487462612"
        [1] "model_2a"         "28.0137935908893"
        [1] "model_2b"         "5.69337438932066"
        [1] "model_2c"         "-9.31979403027607"
        [1] "model_2d"         "-6.1930960954627"
```

Listing 1.78: Display the results of the BIC on these models using R

```
In [18]: print(c('model_1', BIC(iris.glm.1)))
        print(c('model_2a', BIC(iris.glm.2a )))
        print(c('model_2b', BIC(iris.glm.2b )))
        print(c('model_2c', BIC(iris.glm.2c )))
        print(c('model_2d', BIC(iris.glm.2d )))

Out[18]: [1] "model_1"          "-4.87121487462612"
        [1] "model_2a"         "28.0137935908893"
        [1] "model_2b"         "5.69337438932066"
        [1] "model_2c"         "-9.31979403027607"
        [1] "model_2d"         "-6.1930960954627"
```

Listing 1.79: Define three models with two features removed

```
In [19]: model_3a = gsub('\\+ petal_width','', model_2c)
        model_3b = gsub('\\+ petal_length','', model_2c)
        model_3c = gsub('\\+ sepal_length','', model_2c)
```

Listing 1.80: Perform Logistic Regressions

```
In [20]: iris.glm.3a = glm(model_3a, data=iris.data)
        iris.glm.3b = glm(model_3b, data=iris.data)
        iris.glm.3c = glm(model_3c, data=iris.data)
```

Listing 1.81: Display the results of the BIC using R

```
In [21]: print(c('model_1', BIC(iris.glm.1)))
        print(c('model_2c', BIC(iris.glm.2c)))
        print(c('model_3a', BIC(iris.glm.3a)))
        print(c('model_3b', BIC(iris.glm.3b)))
        print(c('model_3c', BIC(iris.glm.3c)))

Out[21]: [1] "model_1"          "-4.87121487462612"
        [1] "model_2c"         "-9.31979403027607"
        [1] "model_3a"         "25.3174210943167"
        [1] "model_3b"         "15.4504250116728"
        [1] "model_3c"         "-5.0467304546584"
```

Listing 1.82: Display best three feature model

```
In [22]: model_2c

Out[22]: [1] "label ~ 1 + sepal_length + petal_length +
        petal_width"
```

Listing 1.83: Display best two feature model

```
In [23]: model_3c

Out[23]: [1] "label ~ 1 + petal_length + petal_width"
```



Listing 1.84: Display Log-Likelihood's of Models

```

In [24]: print(c('model_1', logLik(iris.glm.1)))
         print(c('model_2c', logLik(iris.glm.2c )))
         print(c('model_3a', logLik(iris.glm.3a )))
         print(c('model_3b', logLik(iris.glm.3b )))
         print(c('model_3c', logLik(iris.glm.3c )))

Out[24]: [1] "model_1"          "17.4675133196018"
         [1] "model_2c"          "17.1864852503787"
         [1] "model_3a"          "-2.63743995896586"
         [1] "model_3b"          "2.29605808235611"
         [1] "model_3c"          "12.5446358155217"

```

Listing 1.85: Display the results of the AIC using R

```

In [25]: print(c('model_1', AIC(iris.glm.1)))
         print(c('model_2a', AIC(iris.glm.2a )))
         print(c('model_2b', AIC(iris.glm.2b )))
         print(c('model_2c', AIC(iris.glm.2c )))
         print(c('model_2d', AIC(iris.glm.2d )))
         print(c('model_3a', AIC(iris.glm.3a )))
         print(c('model_3b', AIC(iris.glm.3b )))
         print(c('model_3c', AIC(iris.glm.3c )))

Out[25]: [1] "model_1"          "-22.9350266392037"
         [1] "model_2a"          "12.960617120408"
         [1] "model_2b"          "-9.35980208116062"
         [1] "model_2c"          "-24.3729705007573"
         [1] "model_2d"          "-21.246272565944"
         [1] "model_3a"          "13.2748799179317"
         [1] "model_3b"          "3.40788383528778"
         [1] "model_3c"          "-17.0892716310434"

```

## 1.6 Using Masks to Create a Vector

### A Vectorized Solution To fizzbuzz

**fizzbuzz** is a canonical "coding interview" problem. You might want to read Joel Grus' humorous attempt to use tensorflow to solve **fizzbuzz**<sup>1</sup>. The challenge is to iterate over the numbers from 1 to 100, printing "fizz" if the number is divisible by 3, "buzz" if the number is divisible by 5, "fizzbuzz" if the number is divisible by 15, and the number itself otherwise. Typically this problem is solved using for-loops and if-else statements and is used as a basic assessment of programming ability. Such a solution might look like this

<sup>1</sup><http://joelgrus.com/2016/05/23/fizz-buzz-in-tensorflow/>

Listing 1.86: a first attempt at `fizzbuzz`

```
In [1]: fizzbuzz = function (n) {  
    for (i in 1:n) {  
        if (i %% 15 == 0) {  
            cat("fizzbuzz", sep="\n")  
        }  
        else if (i %% 5 == 0) {  
            cat("buzz", sep="\n")  
        }  
        else if (i %% 3 == 0) {  
            cat("fizz", sep="\n")  
        }  
        else {  
            cat(i, sep="\n")  
        }  
    }  
}
```

Listing 1.87: display `fizzbuzz` for  $n = 15$ 

```
In [2]: fizzbuzz(15)  
  
Out[2]: 1  
        2  
        fizz  
        4  
        buzz  
        fizz  
        7  
        8  
        fizz  
        buzz  
        11  
        fizz  
        13  
        14  
        fizzbuzz
```

It may be a bit much to come up with a solution to this problem using tensorflow. It is, however, very useful to think about solving this problem using masks and filters. Suppose we begin with a simple solution vector as follows

Listing 1.88: start the `solution` vector

```
In [3]: solution = 1:15  
        solution  
  
Out[3]: [1]  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15
```

The challenge is to replace the values we don't need with the correct strings. Sure we can iterate over this list check the value to see if it's divisible by three or five but using a vectorized solution we can do it all at once.

The Steps to doing this are as follows:

1. Create a mask for a certain condition we might wish to check
2. Use that mask to restrict the values of the original `solution` we are looking at
3. Replace to values of the restricted vector with the appropriate string

First, we create a mask called `mod15_mask`. Note, that when we display it there is only a single `TRUE` value, in the position where the value is divisible by 15 (and in this case is actually 15).

Listing 1.89: display a mod 15 mask

```
In [4]: solution %% 15 == 0

Out[4]:  [1] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
          FALSE FALSE FALSE FALSE
          [13] FALSE FALSE  TRUE
```

Listing 1.90: create the mod 15 mask

```
In [5]: mod15_mask = solution %% 15 == 0
          mod15_mask

Out[5]:  [1] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
          FALSE FALSE FALSE FALSE
          [13] FALSE FALSE  TRUE
```

Next, we filter the `solution` using the `mod15_mask`.

Listing 1.91: filter `solution` using the mod 15 mask

```
In [6]: solution[mod15_mask]

Out[6]: [1] 15
```

Finally, we assign the filtered values the string `"fizzbuzz"`

Listing 1.92: assign values to the filtered `solution` vector

```
In [7]: solution[mod15_mask] = "fizzbuzz"
```

Let's have a look at the current value of our solution.

Listing 1.93: display `solution` vector

```
In [8]: solution
Out[8]: [1] "1"      "2"      "3"      "4"      "5"
         "6"
         [7] "7"      "8"      "9"      "10"     "11"
         "12"
         [13] "13"     "14"     "fizzbuzz"
```

We can repeat this technique to build an entire solution to the problem.

Listing 1.94: a vectorized `fizzbuzz`

```
In [9]: fizzbuzz = function (n) {  
    solution = 1:n  
    mod3_mask = (solution %% 3 == 0)  
    mod5_mask = (solution %% 5 == 0)  
    mod15_mask = (solution %% 15 == 0)  
  
    solution[mod3_mask] = "fizz"  
    solution[mod5_mask] = "buzz"  
    solution[mod15_mask] = "fizzbuzz"  
  
    cat(solution, sep="\n")  
}  
  
fizzbuzz(15)  
  
Out [9]: 1  
2  
fizz  
4  
buzz  
fizz  
7  
8  
fizz  
buzz  
11  
fizz  
13  
14  
fizzbuzz
```

In terms of the why of doing a vectorized approach, there are tremendous speed gains to be had implementing your algorithms using vectors rather than loops<sup>2</sup>.

## 1.7 The Curse of Dimensionality

Listing 1.95: Define Numerical Python environment

```
In [1]: import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
%matplotlib inline
```

---

<sup>2</sup><http://www.noamross.net/blog/2014/4/16/vectorization-in-r-why.html>

Listing 1.96: Define datasets of varying size using `make_classification`

```
In [2]: from sklearn.datasets import make_classification

        n_vals = np.logspace(1, 4, 20)

        datasets = {
            int(n): make_classification(int(n))
            for n in n_vals
        }
```

Listing 1.97: Load models from Scikit-Learn

```
In [3]: from sklearn.neighbors import KNeighborsClassifier
        from sklearn.tree import DecisionTreeClassifier
        from sklearn.linear_model import LogisticRegression
        from sklearn.svm import SVC
```

Listing 1.98: Define dictionary of models

```
In [4]: models = {
        'knn' : KNeighborsClassifier(),
        'lr' : LogisticRegression(),
        'dt' : DecisionTreeClassifier(),
        'svm' : SVC(),
        }
```

Listing 1.99: import time library

```
In [5]: from time import time
```

Listing 1.100: Define functions to time fit and prediction

```
In [6]: def time_fit(data, model):
        start = time()
        model.fit(*data)
        fit_time = time() - start
        return fit_time

        def time_predict(data, model):
            start = time()
            model.predict(data[0])
            predict_time = time() - start
            return predict_time
```

Listing 1.101: Define function to test all models on a dataset of size  $n$ 

```
In [7]: def fit_predict_model_on_n(model_name, n):

        model = models[model_name]
        dataset = datasets[n]

        fit_time = time_fit(dataset, model)

        predict_time = time_predict(dataset, model)

        return {
            'model_name' : model_name,
            'n' : n,
            'fit_time' : fit_time,
            'predict_time' : predict_time
        }
```

Listing 1.102: Run all models against all datasets

```
In [8]: results = []
        for n in datasets.keys():
            for model in models.keys():
                results.append(fit_predict_model_on_n(model, n))
```

Listing 1.103: Collect results in a DataFrame

```
In [9]: results_df = pd.DataFrame(results)
```

Listing 1.104: Display a sample of the results

```
In [10]: results_df.sample(5)
```

Table 1.8: Display a sample of the results

	fit_time	model_name	n	predict_time
49	0.001512	lr	784	0.000087
61	0.003356	lr	2335	0.000208
41	0.001146	lr	379	0.000084
45	0.001624	lr	545	0.001891
16	0.000275	knn	42	0.000526

Listing 1.105: Plot Fit Time and Prediction Time

```

In [11]: fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(20,5))

        for model_name in results_df.model_name.unique():

            model_mask = (results_df.model_name == model_name)
            model_results_df = results_df[model_mask]

            ax1.set_title('Fit Time')
            ax1.plot(model_results_df.n,
                    model_results_df.fit_time,
                    label=model_name)
            ax1.set_ylim((0,5))
            ax1.legend()

            ax2.set_title('Prediction Time')
            ax2.set_ylim((0, 5))
            ax2.plot(model_results_df.n,
                    model_results_df.predict_time,
                    label=model_name)
            ax2.legend()

        plt.show()

```

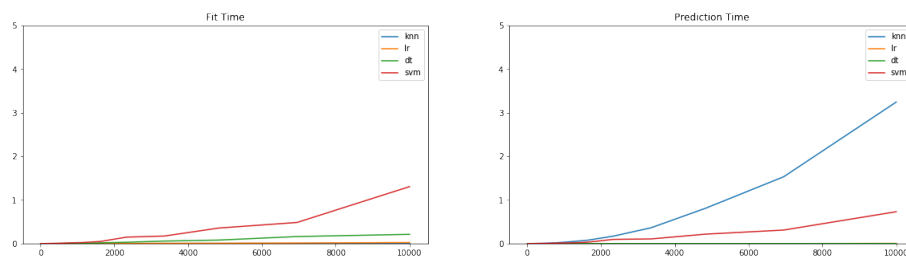


Figure 1.12: Plot Fit Time and Prediction Time



Listing 1.106: Plot Fit Time and Prediction Time with Logarithmic Scale for  $n$ 

```

In [12]: fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(20,5))

        for model_name in results_df.model_name.unique():

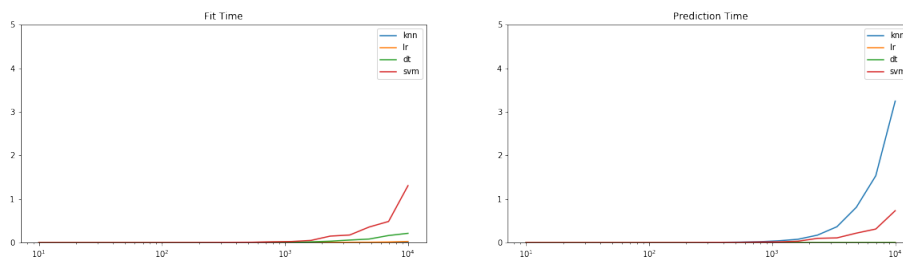
            model_mask = (results_df.model_name == model_name)
            model_results_df = results_df[model_mask]

            ax1.set_title('Fit Time')
            ax1.plot(model_results_df.n,
                    model_results_df.fit_time,
                    label=model_name)
            ax1.set_xscale('log')
            ax1.set_ylim((0, 5))
            ax1.legend()

            ax2.set_title('Prediction Time')
            ax2.plot(model_results_df.n,
                    model_results_df.predict_time,
                    label=model_name)
            ax2.set_xscale('log')
            ax2.set_ylim((0, 5))
            ax2.legend()

        plt.show()

```

Figure 1.13: Plot Fit Time and Prediction Time with Logarithmic Scale for  $n$ 

## 1.8 Introduction to the bias-variance tradeoff

The aim of this lecture is to visually explore the concept of bias and variance in modeling and the tradeoff between the two.

Below we will be fitting models predicting student morale from the day number in the course.

Let's assume that the Bruin gods model student morale as such according to a cubic polynomial interpolated from the following data, one data point for each week:

Listing 1.107: True measurements of morale

```
In [1]: moralepoints = [  
        20, 30, 35, 18, 3, 12,  
        35, 44, 53, 62, 73  
    ]
```

Listing 1.108: Define Numerical Python environment

```
In [2]: import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
from scipy.interpolate import interp1d  
from sklearn.linear_model import LinearRegression  
from sklearn.preprocessing import PolynomialFeatures  
from sklearn.pipeline import Pipeline  
  
np.random.seed(42)  
  
%matplotlib inline
```

Listing 1.109: Define vector of days

```
In [3]: weeks = 10  
days = np.arange(weeks*7+1)
```

## Actual Data According to the Universe

Listing 1.110: Create tuples of weekpoints and measurements of morale

```
In [4]: weekpoints = 7*np.arange(weeks+1)  
  
        list(zip(weekpoints, moralepoints))  
  
Out [4]: [(0, 20),  
          (7, 30),  
          (14, 35),  
          (21, 18),  
          (28, 3),  
          (35, 12),  
          (42, 35),  
          (49, 44),  
          (56, 53),  
          (63, 62),  
          (70, 73)]
```

## Morale function as Interpolated by the Bruin Gods

The `interp1d` function creates an interpolation function for us between the week numbers and morale points. The `kind='cubic'` parameter indicates the smoothing of the interpolation.

Listing 1.111: Build interpolation function using weekpoints and measurements of morale

```
In [5]: morale_func = interp1d(weekpoints,
                               moralepoints,
                               kind='cubic')
```

We can plot out the "true function" of days predicting morale below:

Listing 1.112: Compute true values associated with days

```
In [6]: morale_true = morale_func(days)
```

Listing 1.113: Plot true morale over the span of a course

```
In [7]: plt.figure(figsize=(20, 5))

        plt.plot(days, morale_true,
                 linewidth=5, color='gold',
                 alpha=0.5, label='true function')

        plt.xlabel('days', fontsize=16)
        plt.ylabel('morale', fontsize=16)
        plt.title('Morale over time', fontsize=20)

        plt.legend(loc='upper left')

        plt.show()
```

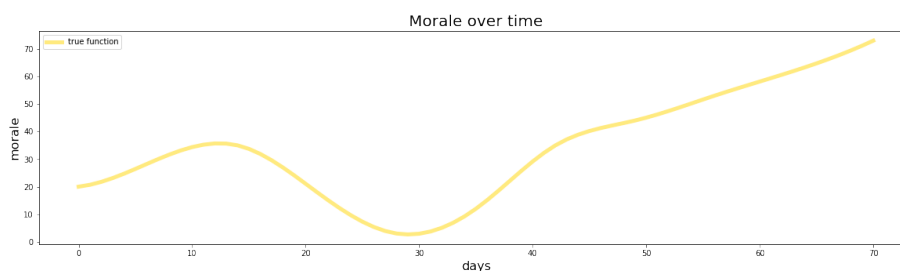


Figure 1.14: Plot true morale over the span of a course

### Important Points

1. We can never know this function.
2. We can only model this function.
3. In building a model, we will introduce error.
4. A model will have three kinds of error:

$$\text{Error} = \text{Error from Bias} + \text{Error from Variance} + \text{Unknown Error}$$

5. All three of these are strictly non-negative.

**The goal when building a model is to minimize the error from bias and the error from variance.**

The perfect model would have zero bias and zero variance. In reality, the bias or the variance (typically) increases while the other decreases.

We can interpret this "true function" in different ways:

1. **As morale without individual variance:** the true function is the baseline morale for each time point that all students vary around to some degree. Students' morale at any given time point is this baseline morale plus or minus some individual deviation. When coding this in python, for each "student" we take the *true* measure of morale from the morale function, and then to create the individual variance we will add random noise to that starting point.
2. **As morale without measurement error:** maybe all student morale at every timepoint is the same, it's just that our way of measuring morale is unreliable. The true function represents the morale at any timepoint for any student without measurement error. If our measurement tool was perfect (a perfect morale survey) we would measure the same morale for every student at each time point. But if the survey form is broken and randomly changes answers by students, for example, then we are adding noise to the measured morale for each observation.
3. **As an average morale across infinite students:** our measurements of morale vary at each time point for each student, but if we had an infinite number of students and averaged all of their morale measurements across all time points, we would end up with the true function of `morale(time)`.

**In each case, we are describing morale as a function of time without error, where each is describing a different source of error:**

1. Error resulting from an imperfect relationship between time and morale.
2. Error resulting from an imperfect ability to measure morale.
3. Error resulting from an insufficient amount of data to quantify the relationship correctly.

We will see that each of these sources of error combined make up the full error in any model we build. Respectively:

1. The **bias**.
2. The **variance**.
3. The **irreducible error**.

You will always have error in your models, it just depends how much and what proportion of each type. We will formalize these components of error in a model further down.

## Generate a sample of students

Say we have four students: student A, B, C, and D.

Each student has had their morale checked 12 different times throughout the course, but not necessarily at the same times.

**Below is a function that will generate the days and morale for each student as a dictionary object.**

*Note: here we are using the morale function to get the true morale for a given day, then adding the individual variance and/or measurement error from a normal distribution.*

Listing 1.114: Define a dictionary of measured morale for four students

```
In [8]: students = {}
        for student in ['A','B','C','D']:

            daysamp = np.random.choice(days, 12)
            morales = morale_func(daysamp)
            morales += np.random.normal(0, 13, 12)

            students[student] = {
                'days' : daysamp,
                'morale': morales
            }
```

### Student A's morale over time

Below we can plot student A's morale at each day. The true function is also plotted in yellow.

Listing 1.115: Plot Measured Values for Student A versus True Morale

```
In [9]: plt.figure(figsize=(20, 5))

        plt.plot(days, morale_true,
                  linewidth=5, color='gold',
                  alpha=0.5, label='true function')

        plt.scatter(students['A']['days'],
                    students['A']['morale'],
                    s=70, c='darkred',
                    label='student A', alpha=0.7)

        plt.xlabel('days', fontsize=16)
        plt.ylabel('morale', fontsize=16)
        plt.title('Morale over time\n', fontsize=20)
        plt.xlim([0, 72])

        plt.legend(loc='upper left')

        plt.show()
```

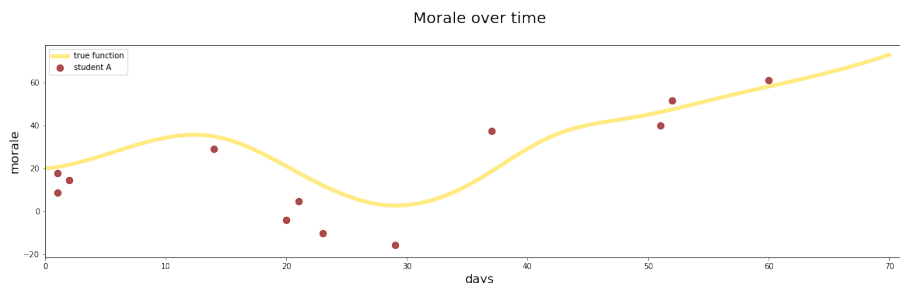


Figure 1.15: Plot Measured Values for Student A versus True Morale

### Build a model for days predicting morale using student A's data

With this student's data, I decide to model the relationship between days and morale with a linear regression.

My model is:

$$\hat{morale} = \beta_0 + \beta_1 days$$

**Construct the model:**

Listing 1.116: Construct a Linear Regression model for Student A

```

In [10]: studA_days = students['A']['days']
        studA_mor = students['A']['morale']

        Amod = LinearRegression()
        Amod.fit(studA_days[:, np.newaxis], studA_mor)

Out[10]: LinearRegression(copy_X=True, fit_intercept=True,
        n_jobs=1, normalize=False)

```

Listing 1.117: Plot the modeled relationship between days and morale

```

In [11]: plt.figure(figsize=(20, 5))

        plt.plot(days, morale_true,
                  linewidth=5, color='gold',
                  alpha=0.5, label='true function')

        plt.plot(days, Amod.predict(days[:, np.newaxis]),
                  lw=7., c='darkred',
                  alpha=0.5, label='model')

        plt.scatter(students['A']['days'],
                    students['A']['morale'],
                    s=70, c='darkred',
                    label='student A', alpha=0.7)

        plt.xlabel('days', fontsize=16)
        plt.ylabel('morale', fontsize=16)
        plt.title('Morale over time\n', fontsize=20)
        plt.xlim([0, 72])

        plt.legend(loc='upper left')

        plt.show()

```

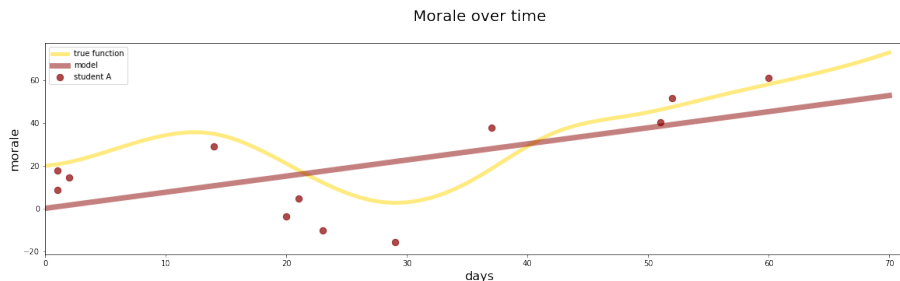


Figure 1.16: Plot the modeled relationship between days and morale

### The total error of a model

As you can see above, our regression line is an imperfect representation of the true function. Furthermore, the regression model fails to perfectly capture the variance of the morale in our sample data.

The model is our blueprint for the estimation of morale. We have chosen to estimate morale from simply the number of days that have elapsed in the course. In doing so, we have made an assumption: morale is a linear function of days.

**When we talk about the bias-variance tradeoff in modeling, and the bias and variance components of error, it is important to think about this in the context of building our model on many samples of the data.** For example, we take our model parameterization `morale ~ days` and build this on the data for student A, then student B, then student C, then student D, and so on. Think of the "students" as random samples of morale and days elapsed from the overall population.

**There are three sources of error in a model:**

$$\text{Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

$$E\left[(\hat{f}(x) - f(x))^2\right] = \left(E[\hat{f}(x)] - f(x)\right)^2 + E\left[(\hat{f}(x) - E[\hat{f}(x)])^2\right] + E\left[(y - f(x))^2\right]$$

Where:

- $f(x)$  is the true function of  $y$  given the predictors.
- $\hat{f}(x)$  is the estimate of  $y$  with the model fit on a random sample of the predictors.
- $E\left[(\hat{f}(x) - f(x))^2\right]$  is the average squared error across multiple models fit on different random samples between the model and the true function.
- $E[\hat{f}(x)]$  is the average of estimates for given predictors across multiple models fit on different random samples.
- $E\left[(y - f(x))^2\right]$  is the average squared error between the true values and the predictions from the true function of the predictors. This is the **irreducible error**.



- $(E[\hat{f}(x)] - f(x))^2$  is the squared error between the average predictions across multiple models fit on different random samples and the prediction of the true function. This is the **bias** (squared).
- $E[(\hat{f}(x) - E[\hat{f}(x)])^2]$  is the average squared distance between individual model predictions and the average prediction of models across multiple random samples. This is the **variance**.

**The irreducible error is "noise" – error in the measurement of our target that cannot be accounted for by our predictors.**

- The true function represents the most perfect relationship between predictors and target, but that does not mean that our variables can perfectly predict the target.
- The irreducible error can be thought of as the measurement error: variation in the target that we cannot represent.

We will go into the bias and variance components individually, in more detail.

## Bias

The *bias*<sup>2</sup> is the source of error in our model that represents how *oversimplified* our model is. In our example, we have built a function to predict the morale of students as a linear function of days elapsed. However, we can see from the true function that the relationship between days and morale is not a line. This error is encapsulated in the bias.

$$Bias^2 = (E[\hat{f}(x)] - f(x))^2$$

**Remember:** we have to think about measuring the bias and variance using different fits of our model across multiple random samples! The bias represents the deviation of the *average predictions across models* from the true function.

**What does having a high vs. low bias mean?**

- If our models are consistently wrong, then the bias will be large.
- Alternatively, if our models are consistently correct then the bias will be small.

- Bias will be small if the errors across our models built on random samples and tested using the same predictors are incorrect in different directions that average out close to 0.

Linear methods like regression tend to have a high bias because we construct a simplification of the true function.

**Below we can plot the error differences between the true function and our model built on student A's data.**

No matter how many lines we fit on different students, those regression lines are never going to average together into the nonlinear true function! Our average estimate across models will always deviate from the true function somewhere across the days.

Listing 1.118: Morale over time (difference from true function)

```
In [12]: plt.figure(figsize=(20, 5))

        plt.plot(days, morale_true,
                  linewidth=5, color='gold',
                  alpha=0.5, label='true function')

        predictions = Amod.predict(days[:, np.newaxis])

        plt.plot(days, predictions,
                  lw=7., c='darkred',
                  alpha=0.5, label='model')

        plt.scatter(students['A']['days'],
                    students['A']['morale'],
                    s=70, c='darkred',
                    label='student A', alpha=0.7)

        for d in students['A']['days']:
            p = Amod.predict(d)
            plt.plot([d, d],
                    [p, morale_func(d)],
                    c='black', lw=3.)

        plt.xlabel('days', fontsize=16)
        plt.ylabel('morale', fontsize=16)
        plt.xlim([0, 72])

        plt.legend(loc='upper left')

        plt.show()
```

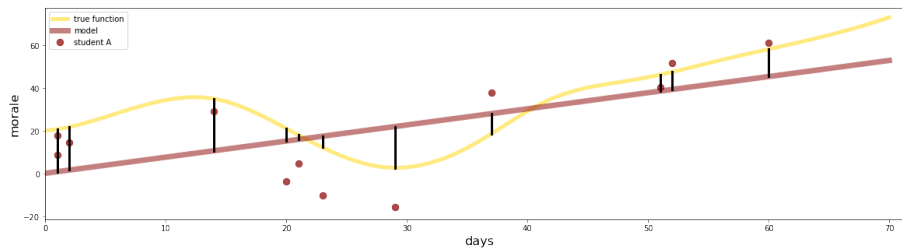


Figure 1.17: Morale over time (difference from true function)

## Variance

The second component of error in the model is the variance of our predictions.

Variance describes the extent to which the individual predictions from models built on different samples (students) deviate from the mean of all the model predictions.

$$\text{Variance} = E[ (\hat{f}(x) - E[\hat{f}(x)])^2 ]$$

### High vs. low variance

- The variance will be large if, for the same observation, models built on different random samples of the data will produce very different predictions.
- Variance is a measure of how *consistent* our model's predictions will be if it were fit on another sample of data.
- Variance is low if the data we train the model on has very little effect on the predictions.

Note that variance is not a measure of how correct or incorrect the predictions are. It is a measure of how variable they are!

### Measuring more students

To better visualise the concept of model variance, let's say we measure a second student, student "B":

Listing 1.119: Plot Measured Values for Students A and B versus True Morale

```

In [13]: plt.figure(figsize=(20, 5))

        plt.plot(days, morale_true,
                  linewidth=5, color='gold',
                  alpha=0.5, label='true function')

        plt.scatter(students['A']['days'],
                    students['A']['morale'],
                    s=70, c='darkred',
                    label='student A', alpha=0.7)

        plt.scatter(students['B']['days'],
                    students['B']['morale'],
                    s=70, c='steelblue',
                    label='student B', alpha=0.7)

        plt.xlabel('days', fontsize=16)
        plt.ylabel('morale', fontsize=16)
        plt.title('Morale over time\n', fontsize=20)
        plt.xlim([0, 72])

        plt.legend(loc='upper left')

        plt.show()

```

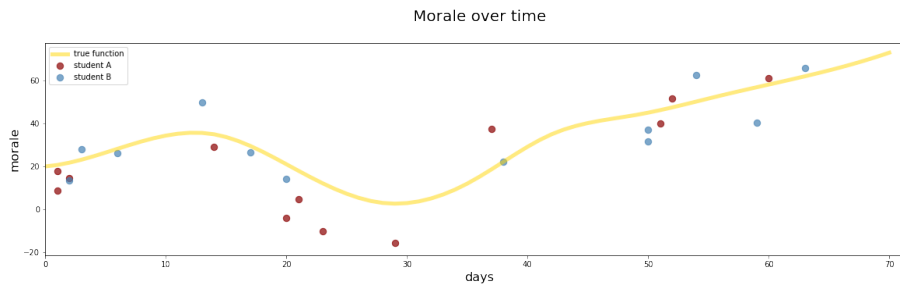


Figure 1.18: Plot Measured Values for Students A and B versus True Morale

Listing 1.120: Fit model for Student B

```

In [14]: studB_days = students['B']['days']
        studB_mor = students['B']['morale']

        Bmod = LinearRegression()
        Bmod.fit(studB_days[:, np.newaxis], studB_mor)

Out[14]: LinearRegression(copy_X=True, fit_intercept=True,
                          n_jobs=1, normalize=False)

```

Listing 1.121: Morale over time (difference between estimates)

```

In [15]: plt.figure(figsize=(20, 5))

plt.plot(days, morale_true,
         linewidth=5, color='gold',
         alpha=0.5, label='true function')

Apred = Amod.predict(days[:, np.newaxis])

plt.plot(days, Apred,
         lw=7., c='darkred',
         alpha=0.5, label='A model')

Bpred = Bmod.predict(days[:, np.newaxis])

plt.plot(days, Bpred,
         lw=7., c='steelblue',
         alpha=0.5, label='B model')

plt.scatter(students['A']['days'],
           students['A']['morale'],
           s=70, c='darkred',
           label='student A', alpha=0.7)

plt.scatter(students['B']['days'],
           students['B']['morale'],
           s=70, c='steelblue',
           label='student B', alpha=0.7)

plt.fill_between(days, Apred, Bpred,
               color='grey', hatch='//',
               edgecolor=None)

plt.xlabel('days', fontsize=16)
plt.ylabel('morale', fontsize=16)
plt.xlim([0, 72])

plt.legend(loc='upper left')

plt.show()

```

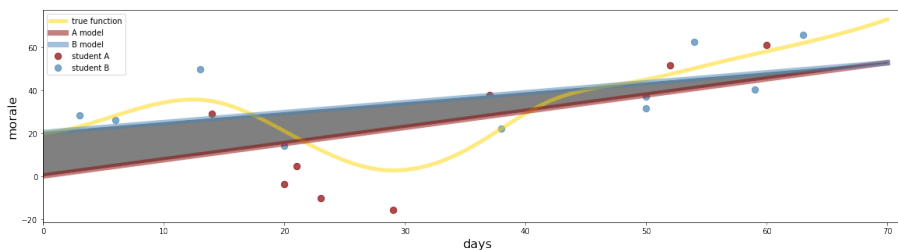


Figure 1.19: Morale over time (difference between estimates)

## Bias and variance using 3 models

Below we will assess the morale of 3 different students over the days in the course at different times.

We can build these simple `morale ~ time` models for each and plot the regression lines.

These models are **high bias and low variance**. This is because there is a considerable amount of difference between the average of the model predictions and the true function, but not a lot of variation in predictions at time points across our models for the 3 students.

Listing 1.122: Morale over Time

```
In [16]: plt.figure(figsize=(20, 5))

plt.plot(days, morale_true,
         linewidth=5, color='gold',
         alpha=0.5, label='true function')

plt.scatter(students['A']['days'],
           students['A']['morale'],
           s=70, c='darkred',
           label='student A', alpha=0.7)

plt.scatter(students['B']['days'],
           students['B']['morale'],
           s=70, c='steelblue',
           label='student B', alpha=0.7)

plt.scatter(students['C']['days'],
           students['C']['morale'],
           s=70, c='darkgreen',
           label='student C', alpha=0.7)

plt.xlabel('days', fontsize=16)
plt.ylabel('morale', fontsize=16)
plt.xlim([0, 72])

plt.legend(loc='upper left')

plt.show()
```

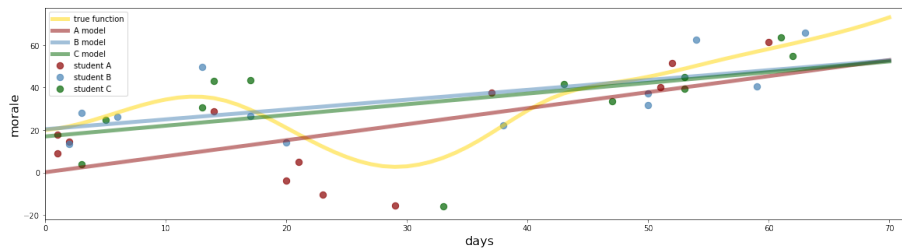


Figure 1.20: Morale over Time

Listing 1.123: Fit model for Student B

```
In [17]: studC_days = students['C']['days']
        studC_mor = students['C']['morale']

        Cmod = LinearRegression()
        Cmod.fit(studC_days[:, np.newaxis], studC_mor)

Out[17]: LinearRegression(copy_X=True, fit_intercept=True,
                          n_jobs=1, normalize=False)
```

Listing 1.124: Morale over Time

```
In [18]: plt.figure(figsize=(20, 5))

plt.plot(days, morale_true,
         linewidth=5, color='gold',
         alpha=0.5, label='true function')

Apred = Amod.predict(days[:, np.newaxis])

plt.plot(days, Apred,
         lw=5., c='darkred',
         alpha=0.5, label='A model')

Bpred = Bmod.predict(days[:, np.newaxis])

plt.plot(days, Bpred,
         lw=5., c='steelblue',
         alpha=0.5, label='B model')

Cpred = Cmod.predict(days[:, np.newaxis])

plt.plot(days, Cpred,
         lw=5., c='darkgreen',
         alpha=0.5, label='C model')

plt.scatter(students['A']['days'],
            students['A']['morale'],
            s=70, c='darkred',
            label='student A', alpha=0.7)

plt.scatter(students['B']['days'],
            students['B']['morale'],
            s=70, c='steelblue',
            label='student B', alpha=0.7)

plt.scatter(students['C']['days'],
            students['C']['morale'],
            s=70, c='darkgreen',
            label='student C', alpha=0.7)

plt.xlabel('days', fontsize=16)
plt.ylabel('morale', fontsize=16)
plt.xlim([0, 72])

plt.legend(loc='upper left')

plt.show()
```



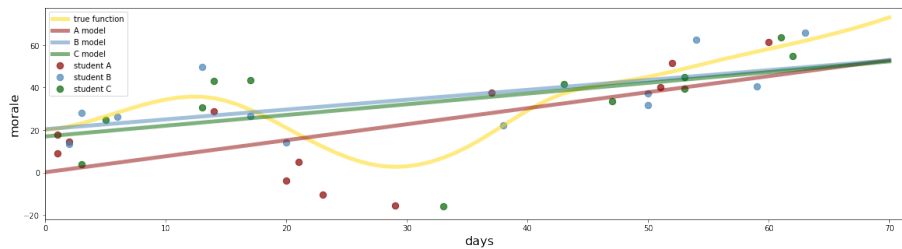


Figure 1.21: Morale over Time

### Increasing complexity to try and capture the true function

Currently our models predict that morale simply increases over time.

Looking at the true function, we can see that there is an increase, then a decrease, and then an increase again.

Just modeling with a linear effect of time alone can't fit a curve; there is only one coefficient being multiplied by time to make our morale prediction. We could, however, add more variables created from time such as  $time^2$ ,  $time^3$ , etc.

$$\hat{morale} = \beta_0 + \beta_1 t + \beta_2 t^2$$

$$\hat{morale} = \beta_0 + \beta_1 t + \beta_2 t^2 + \beta_3 t^3 + \beta_4 t^4$$

$$\hat{morale} = \beta_0 + \beta_1 t + \dots + \beta_{16} t^{16}$$

The plots below show the difference in the fit when you add different numbers of "polynomial" time variables:

Listing 1.125: Define a function to prepare a Polynomial Regression function

```
In [19]: def polynomial_modeler(X, y, degrees):

    poly_feat = PolynomialFeatures(degree=degrees,
                                   include_bias=False)

    linear_regression = LinearRegression()

    pipeline = Pipeline([
        ("polynomial_features", poly_feat),
        ("linear_regression", linear_regression)
    ])

    pipeline.fit(X[:, np.newaxis], y)

    return pipeline
```

Listing 1.126: Define a function to prepare a Plot of Polynomial Regressions for various powers

```
In [20]: def plot_polyfit(X, y, truefunc, degrees=[1,2,4,16],
                        student_color='darkred', name='A'):

    # set the plot size
    plt.subplots(figsize=(20,2))

    # create a plot for each polynomial degree plotted
    for i in range(len(degrees)):
        ax = plt.subplot(1, len(degrees), i + 1)

        poly_model = polynomial_modeler(X, y, degrees[i])

        X_test = np.linspace(1, 70, 200)
        predictions = poly_model.predict(X_test[:,
np.newaxis])

        plt.plot(X_test, predictions,
                 lw=5., c=student_color,
                 label="model", alpha=0.6)

        plt.plot(days, morale_true,
                 linewidth=5, color='gold',
                 alpha=0.5, label='true function')

        plt.scatter(X, y,
                    label="Student observations",
                    c=student_color, s=40)

        plt.xlabel("days")
        plt.ylabel("morale")
        plt.xlim((0, 72))
        plt.ylim((-20, 100))
        plt.legend(loc="best")

        plt.title('Student ' +
                  name +
                  f" (degree {degrees[i]})")
```

Listing 1.127: Run the Polynomial Plot for Student A

```
In [21]: plot_polyfit(students['A']['days'],
                      students['A']['morale'],
                      morale_func,
                      student_color='darkred',
                      name='A')
```

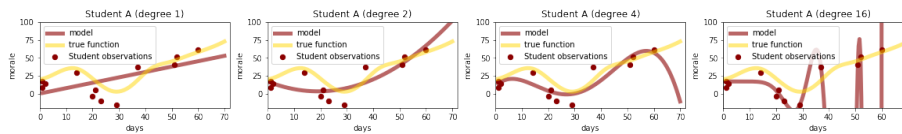


Figure 1.22: Run the Polynomial Plot for Student A

Listing 1.128: Run the Polynomial Plot for Student B

```
In [22]: plot_polyfit(students['B']['days'],
                    students['B']['morale'],
                    morale_func,
                    student_color='steelblue',
                    name='B')
```

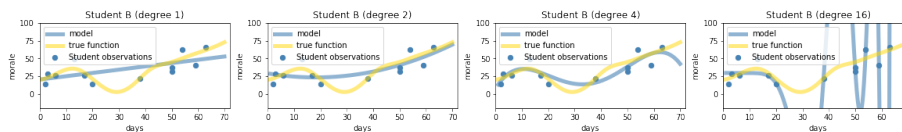


Figure 1.23: Run the Polynomial Plot for Student B

Listing 1.129: Run the Polynomial Plot for Student C

```
In [23]: plot_polyfit(students['C']['days'],
                    students['C']['morale'],
                    morale_func,
                    student_color='darkgreen',
                    name='C')
```

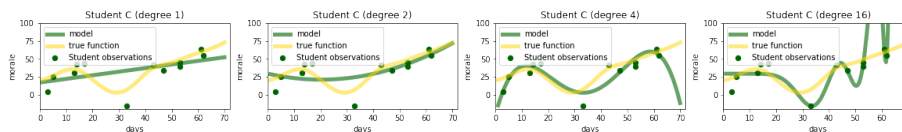


Figure 1.24: Run the Polynomial Plot for Student C

## Higher complexity means higher variance (and lower bias)

The variance of predictions across our models goes up as we increase the model complexity. This is equivalent to saying that the variance of our model is increasing.

Increasing the complexity of the model at the expense of good future predictions is known as "overfitting" the data. High variance and overfitting are

intrinsically related: if your predictions are inconsistent across samples, you are more likely to make the wrong predictions on future data.

Likewise, high bias and underfitting are related. If your model is too basic, it may give very consistent predictions but at the cost of oversimplifying the relationship between the target and predictors.

**Below are student A and student B fit with the 16-polynomial time model and the area showing the difference in predictions at time points between them.** Compare this to the area we saw earlier with the single time term.

Listing 1.130: Plot 16<sup>th</sup>-order Regression Models for Students A and B versus True Morale

```

In [24]: plt.figure(figsize=(20, 5))

plt.plot(days, morale_true,
         linewidth=5, color='gold',
         alpha=0.5, label='true function')

Amod_complex = polynomial_modeler(students['A']['days'],
students['A']['morale'],
                                16)
Bmod_complex = polynomial_modeler(students['B']['days'],
students['B']['morale'],
                                16)

Apred = Amod_complex.predict(days[:, np.newaxis])

plt.plot(days, Apred,
         lw=7., c='darkred',
         alpha=0.5, label='A model')

Bpred = Bmod_complex.predict(days[:, np.newaxis])

plt.plot(days, Bpred,
         lw=7., c='steelblue',
         alpha=0.5, label='B model')

plt.scatter(students['A']['days'],
            students['A']['morale'],
            s=70, c='darkred',
            label='student A', alpha=0.7)

plt.scatter(students['B']['days'],
            students['B']['morale'],
            s=70, c='steelblue',
            label='student B', alpha=0.7)

plt.fill_between(days, Apred, Bpred,
                color='grey', hatch='//',
                edgecolor=None)

plt.xlabel('days', fontsize=16)
plt.ylabel('morale', fontsize=16)
plt.title('Morale over time (16th polynomial)\n',
         fontsize=20)
plt.xlim([0, 72])
plt.ylim([-2500, 2500])

plt.legend(loc='upper left')

plt.show()

```

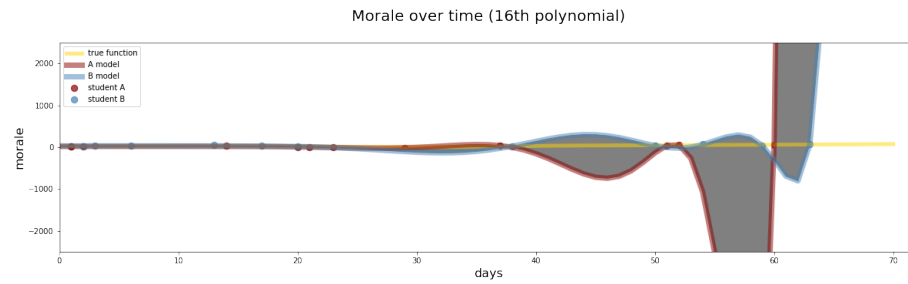


Figure 1.25: Plot 16<sup>th</sup>-order Regression Models for Students A and B versus True Morale

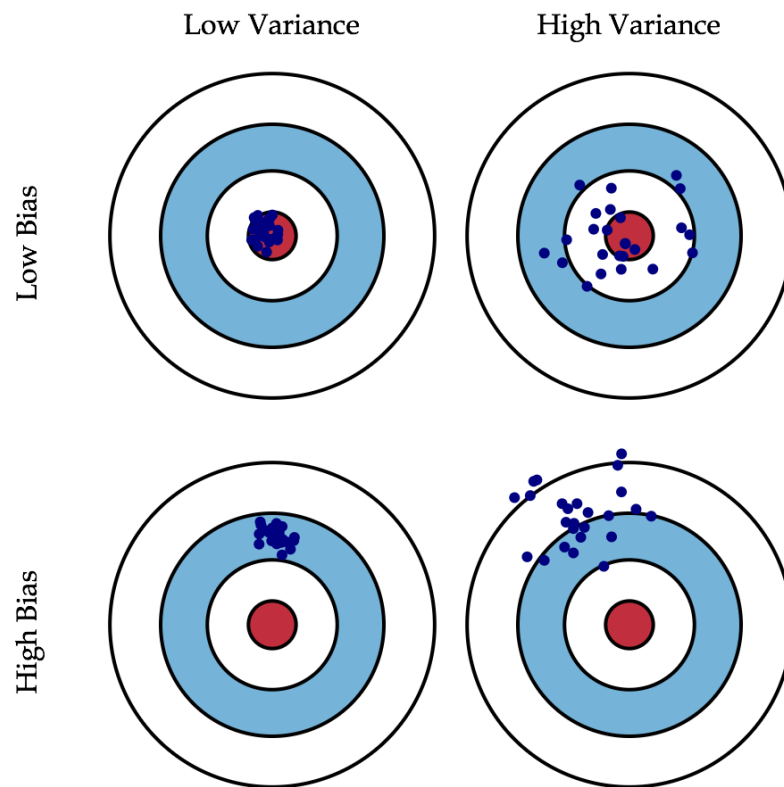


Figure 1.26: Illustrating the Bias-Variance Tradeoff

### An Illustrative Example: Voting Intentions

Imagine we conduct a phone poll to see how voters in a town will vote in the next election and get these results:

— Voting Republican —	Voting Democratic —	Non-Respondent —	Total —
13	16	21	50

Listing 1.131: Create Poll Simulation

```
In [25]: poll_results = {'voting_republican': 13,
                        'voting_democratic': 16,
                        'non-respondent': 21}

        total = sum([val for key, val in poll_results.items()])

        total_voting = sum([
            val for key, val in poll_results.items()
            if key != 'non-respondent'
        ])
```

Listing 1.132: Probability of Voting Republican According to Poll

```
In [26]: probability_of_voting_R =
        (poll_results['voting_republican']/
         float(total_voting))

        probability_of_voting_R

Out[26]: 0.4482758620689655
```

We put out our press release that the Democrats are going to win by over 10 points; but, when the election comes around, it turns out they actually lose by 10 points.





## Chapter 2

# Case Study: Seeds

This will be a study of the Seeds data set available through the UCI Machine Learning Repository<sup>1</sup>, where additional information on the data set can be found<sup>2</sup>. The data set is very similar to Fisher’s Iris data set. It is used here as an alternative to Iris, as a simple and straightforward data set that is useful for introducing the main concepts of statistical learning and data science, especially Supervised and Unsupervised Learning and Classification, Regression, and Clustering.

The UCI ML Repository describes the data as

“The data collected here is measurements of geometrical properties of kernels belonging to three different varieties of wheat. A soft X-ray technique and GRAINS package were used to construct all seven, real-valued attributes.

The **target** (also called labels, response, or dependent variable) fall into one of three classes just as in the Iris data set. The  $n$ , the number of observations (or points or rows), in the data set is very small. Iris contains 150 rows, whereas Seeds contains 210. The  $p$ , the number of salient **features** (also called predictors, input, or independent variables) is a bit larger than in Iris. Iris has four features, whereas Seeds has seven. In both data sets, all of the features are numerical and continuous.

## 2.1 Structure of the Case Study

The Case Study consists of three notebooks, a `README.md` file, and a Python file, `visualization.py`, arranged as follows<sup>3</sup>:

---

<sup>1</sup><https://archive.ics.uci.edu/ml/index.php>

<sup>2</sup><https://archive.ics.uci.edu/ml/datasets/seeds>

<sup>3</sup>The `tree` tool is a useful command line tool that is available via `apt`, `brew`, and `yum`.

Listing 2.1: The Structure of the Case Study

```
$ tree
.
|-- 00-introduction.ipynb
|-- 01-load-and-verify-data.ipynb
|-- 02-exploratory-data-analysis.ipynb
+-- visualization.py
```

The `00-introduction.ipynb` notebook contains information describing the data set and the structure of the project. The `01-load-and-verify-data.ipynb` notebook is the first phase of the Case Study, the **Loading and Verification of Data** and contains all steps required for obtaining, loading, preprocessing, verifying, and saving locally the data prior to any analysis work. The `02-exploratory-data-analysis.ipynb` notebook is the second phase of the Case Study and contains the **Exploratory Data Analysis** (EDA) work. The `visualization.py` file contains two Python functions that generate advanced visualizations.

### Programming Method: Working in Multiple Notebooks

The reader will note that the work in this Case Study has been split across three notebooks. While it is possible to do the work in this Case Study in a single notebook, we will consider it a best practice to use multiple notebooks and to do so to help to structure the work. Each notebook will represent a major phase of the work, while minor phases of work will be designated with subheadings within a specific notebook. This helps with computing resources, especially the execution of JavaScript required to render a Jupyter Notebook, but also with referring to earlier work in later phases of a project.

### Use of Statistical Learning Models for Exploratory Data Analysis

The EDA work done in this Case Study uses several statistical learning models including:

1. Linear Regression (a supervised regression model)
2. Logistic Regression (a supervised classification model)
3. K-Nearest Neighbors Classifier (a supervised classification model)
4. T-SNE (an unsupervised dimensionality reduction model)
5. K-Means (an unsupervised clustering model)

Each of these models is used with a specific EDA-related purpose. Statistical learning models, especially supervised learning models, are often used for building predictive models whose efficacy is evaluated accordingly. It is of note that these models were used here for the purposes of EDA, and as such, have not been subjected to the rigorous model assessment associated with building a predictive model.

## 2.2 Load and Verify Data

### Load the Seeds Data Set from a URL

You will begin the Case Study by loading the data directly from a URL. The data is hosted in a whitespace-separated file on the UCI Machine Learning Repository<sup>4</sup>

### Programming Method: Importing a Python Module as Alias

The first Python code you will run in this Case Study is code importing libraries necessary for the data science worked performed in this case study. Each of the `import` statements you will run is of the following format:

Listing 2.2: Importing a Python Module as Alias

```
import full_library_name as short_library_name
```

This has the effect of importing the entire library to be used under an alias. For example, you will import `numpy` using the import statement

Listing 2.3: Importing `numpy` with an alias

```
import numpy as np
```

This imports the entire `numpy` library, but the library will only be available under the shortened alias `np`. This is a common technique in Python programming and useful when we will be making frequent reference to the library.

---

<sup>4</sup>[https://archive.ics.uci.edu/ml/machine-learning-databases/00236/seeds\\_dataset.txt](https://archive.ics.uci.edu/ml/machine-learning-databases/00236/seeds_dataset.txt)

### Import the Python Numerical Stack

The entire Case Study is performed using the Python Numerical Stack. You first load the libraries that comprise the Python numerical stack, including

- matplotlib - numpy - pandas - seaborn.

Listing 2.4: Import the Python Numerical Stack

```
In [1]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import seaborn as sns
```

### Run the IPython Magic Command `%matplotlib inline`

You make use of the IPython Magic<sup>5</sup> command `matplotlib` to configure the notebook for interactive `matplotlib` work. The `inline` argument specifies that we wish to work interactively with a Jupyter Notebook. Note that as with other IPython Magic commands, we begin this line of code with the `%` symbol. The result is that any images generated by our Python code will be rendered within the notebook while you work.

Listing 2.5: Render images inline with IPython magic

```
In [2]: %matplotlib inline
```

### Define the Dataset URL

Next, you define a variable to hold the URL of the data set as a string. Note that you give the variable an all caps name. While this has no semantic meaning in Python, you do this to signify that this string should be thought of as a constant, that is that it is a value to be stored by the notebook that will never change.

Listing 2.6: Define the Dataset URL

```
In [3]: SEEDS_URL = ('https://archive.ics.uci.edu/' +
                    'ml/machine-learning-databases/00236/seeds_dataset.txt')
```

You display the string stored as `SEEDS_URL` to verify that it matches the expected value.

---

<sup>5</sup><https://ipython.readthedocs.io/en/stable/interactive/magics.html>

Listing 2.7: Display the Dataset URL

```
In [4]: SEEDS_URL

Out [4]:
'https://archive.ics.uci.edu/ml/machine-learning-databases/00236/seeds_dataset.txt'
```

### Display a sample of the Data Set URL

It is worth entering the `SEEDS_URL` into a browser, both to verify that the URL matches the correct URL, but also to get a preliminary sense of what the data should look like. Copy the value stored as `SEEDS_URL` and paste it into the address bar of any browser. If the URL is correct, you will see the data stored at the UCI Machine Learning Repository in your browser window as plain text. The first five rows of the data are as follows:

Listing 2.8: A Sample of the Seeds dataset text file

15.26	14.84	0.871	5.763	3.312	2.221	5.22	1
14.88	14.57	0.8811	5.554	3.333	1.018	4.956	1
14.29	14.09	0.905	5.291	3.337	2.699	4.825	1
13.84	13.94	0.8955	5.324	3.379	2.259	4.805	1
16.14	14.99	0.9034	5.658	3.562	1.355	5.175	1
...							

### Use `pd.read_csv()` to load data as `DataFrame`

Next, use the `pd.read_csv()` method available as part of the Pandas library to read the text file stored at that URL into a `DataFrame`.

Listing 2.9: Use `pd.read_csv()` to load data as `DataFrame`

```
In [5]: seeds_df = pd.read_csv(SEEDS_URL, header=None, sep='\s+')
```

A bit of special handling is required here. First, it is necessary to specify on load that the file has no header row. This is done using the argument `header=None`. Second, it is necessary to specify that this data is whitespace-separated. In other words, whitespace is used to separate the values in a row of data rather than the more conventional commas. This is done using the regular expression `s+`. This argument signifies that one or more whitespace character(s) should be used as separator.

### Manually Define Column Names

Because no header row is specified in the text file, it is necessary to manually specify the names of each column in your new **DataFrame**. These names were obtained from the attribute information section at the UCI Machine Learning Repository. Note that you make all of the names computer-friendly by replacing whitespace with underscores.

The names of the columns of a **DataFrame** are stored as the attribute `df.columns`. You specify these names manually by assigning a list of the correct length to this attribute.

Listing 2.10: Manually Define Column Names

```
In [6]: seeds_df.columns = [  
        "area",  
        "perimeter",  
        "compactness",  
        "length_of_kernel",  
        "width_of_kernel",  
        "asymmetry_coefficient ",  
        "length_of_kernel_groove",  
        "seed_class"  
    ]
```

### Basic Data Verification

As a means of verifying that the data is properly loaded, you display both the type and the shape of the **DataFrame** and the first five rows. This is done using the `type()` function, the `.shape` attribute, and the `.head()` method, respectively. It is worth noting that `.shape` is an attribute of and `.head()` is a method of the **Pandas.DataFrame** class, while `type()` is a Python built-in function.

Listing 2.11: Show the type of `seeds_df`

```
In [7]: type(seeds_df)  
Out[7]: pandas.core.frame.DataFrame
```

Listing 2.12: Show the shape of `seeds_df`

```
In [8]: seeds_df.shape  
Out[8]: (210, 8)
```

As expected, `seeds_df` is a **Pandas DataFrame**. It has a shape of `(210, 8)`. This corresponds to a  $n = 210$  and  $p = 7$ , plus 1 target column.

Listing 2.13: Show the `.head()` of `seeds_df`

```
In [9]: seeds_df.head()
```

Table 2.1: Show the `.head()` of `seeds_df`

	area	perimeter	compactness	length_of_kernel	width_of_kernel	asymmetry_coefficient	length_of_kern
0	15.26	14.84	0.8710	5.763	3.312	2.221	
1	14.88	14.57	0.8811	5.554	3.333	1.018	
2	14.29	14.09	0.9050	5.291	3.337	2.699	
3	13.84	13.94	0.8955	5.324	3.379	2.259	
4	16.14	14.99	0.9034	5.658	3.562	1.355	

After running `seeds_df.head()` make sure to compare the first five rows of the loaded `DataFrame` to the first five rows of the plain text file in your browser.

## Verify Type Casting

During the loading process Pandas will explicitly cast each column into the most appropriate data type<sup>6</sup>. Based on information available at the UCI machine learning repository page for this data set, it is expected that the first seven columns to be floating-point, that is type `float`, and the eighth target column, `seed_class` to be categorical.

### Display `DataFrame.dtypes`

Next, you use the `.dtypes` attribute to display the data types of your `seeds_df` `DataFrame`.

<sup>6</sup><https://rushter.com/blog/pandas-data-type-inference/>

Listing 2.14: Display DataFrame.dtypes

```
In [10]: seeds_df.dtypes

Out[10]: area                float64
         perimeter           float64
         compactness         float64
         length_of_kernel    float64
         width_of_kernel     float64
         asymmetry_coefficient float64
         length_of_kernel_groove float64
         seed_class          int64
         dtype: object
```

Note that the first seven columns have been correctly cast, but `seed_class` has been incorrectly designated as type `int`. This is no doubt due to the fact that the `seed_class` categories are encoded numerically.

### Display Unique Values for the `seed_class` Column

Using the `.unique()` method, you display the unique values of the `seed_class` column.

Listing 2.15: Display Unique Values for the `seed_class` Column

```
In [11]: seeds_df.seed_class.unique()

Out[11]: array([1, 2, 3])
```

Indeed, the seed classes, Kama, Rosa and Canadian, have been encoded as the numbers 1, 2, and 3. As all of the values are integers, Pandas incorrectly intuited that the column should be of type `int`.

### Programming Method: Transform and Reassign a Variable

The programming pattern being used here is one that is fairly common. Here we take some variable, transform that variable, and then reassign the result to the original variable name. In the simplest sense, you might think about doing this with an integer

Listing 2.16: Transform and Reassign `my_int`

```
my_int = 4
my_int = my_int + 5
```



Of course, this pattern is so common that it has a shortcut that works in most popular programming languages

Listing 2.17: Simple Incrementation of `my_int`

```
my_int = 4
my_int += 5
```

In general, where it is not necessary to keep incremental results of a variable, we use a pattern like this one. For example, you might imagine recording the time that a certain process takes using the `time.time()` function.

Listing 2.18: Time a Process

```
import time
start = time.time()
## execute some process
elapsed_time = time.time() - start
```

Note that `elapsed_time` is a number of seconds stored as a `float`, while you may only be interested in the number of seconds as an `int`. A wasteful method would be

Listing 2.19: Cast `elapsed_time` as an integer

```
elapsed_time_int = int(elapsed_time)
```

Using this, you have actually stored two variables `elapsed_time` and `elapsed_time_int`. When transforming and reassigning, you simply overwrite the first variable with the change you wish to make

Listing 2.20: Cast `elapsed_time` as an integer and reassign to `my_int`

```
elapsed_time = int(elapsed_time)
```

Here you have made use of the transformation and reassignment programming pattern. We cast the time `elapsed_time` as an `int` and reassign this to the original `elapsed_time`, thus only ever storing a single value.

**square Note:** It is worth noting that this is not always seen as a best practice, and in many cases it makes more sense to not transform your variables after you have assigned them. The practice is so common in imperative programming (as opposed to functional), that it is worth familiarizing yourself with it, even if you do not intend to handle your variables in this way.

## Fix Typecasting Error

You will need to correct the incorrect typecasting manually. You will do this using the `.astype()` method. This method, applied to a column, returns a copy of that column cast into the requested data type. To complete the casting, you save this typecast copy back as the original column, using the transformation and reassignment programming pattern. This pattern is:

Listing 2.21: Recast the type of a Pandas Series

```
df.my_column = df.my_column.astype(<REQUESTEDTYPE>)
```

### Cast `seed_class` as a Pandas Column of type `category`

Here, you cast `seed_class` as the Pandas data type, `category`<sup>7</sup>.

Listing 2.22: Cast `seed_class` as a Pandas Column of type `category`

```
In [12]: seeds_df.seed_class =
        seeds_df.seed_class.astype('category')
```

Once more, you display the data types of your `seeds_df` DataFrame.

You now see that the target column, `seed_class` is now correctly encoded as categorical.

Listing 2.23: Display `DataFrame.dtypes`

```
In [13]: seeds_df.dtypes

Out[13]: area                float64
         perimeter           float64
         compactness         float64
         length_of_kernel    float64
         width_of_kernel     float64
         asymmetry_coefficient float64
         length_of_kernel_groove float64
         seed_class          category
         dtype: object
```

You now see that the target column, `seed_class` is correctly encoded as a `category`.

<sup>7</sup><https://pandas.pydata.org/pandas-docs/stable/categorical.html>

## Save Data Set as Local Files

So that you do not need to retrieve your data set from remote URL each time you wish work with it, you next use the Pandas `DataFrame` methods, `.to_csv()` and `.to_pickle()`, to store local versions the downloaded file. Each of these takes a string as its argument. This string will be used as the name of the file to be saved locally.

### Programming Method: The Python Pickle Format

While you are no doubt familiar with the CSV format, this may be the first time you are learning about the Pickle format<sup>8</sup>. The Pickle format is a Python specific way to save information. When you write a `DataFrame` to disk as a Pickle you can load it in later, *exactly* as it was when you saved it. Furthermore, just like Pandas has a `pd.read_csv()` method, it also includes the method `pd.read_pickle()` to help you load a Pickled `DataFrame` at some later point.

Listing 2.24: Save the `DataFrame` as a CSV

```
In [14]: seeds_df.to_csv('seeds.csv')
```

Listing 2.25: Save the `DataFrame` as a pickle

```
In [15]: seeds_df.to_pickle('seeds.p')
```

## 2.3 Exploratory Data Analysis

Next, we perform some exploratory data and analysis (EDA) on our data set. This will be the terminal phase of this Case Study. In general, EDA is performed for the purposes of making inferences over a data set. In a project with a goal of developing a predictive model over the data set, EDA is useful for helping to determine which family of predictive models is most appropriate.

### Visualization with a Pair Plot

Because of its simplicity, the first step in EDA is often to prepare a Pair Plot of the data set. This

### Import the Python Numerical Stack

To begin, you import the Python Numerical Stack. Note that this was done in a previous notebook, but it must be done for each notebook.

---

<sup>8</sup><https://docs.python.org/3/library/pickle.html>

```
import matplotlib.pyplot as plt import numpy as np import pandas as pd
import seaborn as sns
%matplotlib inline
```

### Load Saved Data Set

In `01-load-and-verify-data.ipynb`, you retrieved the data set with which you are working from a remote URL and performed the basic preprocessing step of making sure that all columns are stored as the correct data type. After this, you saved the processed `DataFrame` to a local Pickle file called `seeds.p`. Here, you load the saved `DataFrame` for use in this notebook.

```
seeds_df = pd.read_pickle('seeds.p')
```

### Visualization: Prepare a Pair Plot

A simple yet powerful technique for beginning EDA is the Seaborn function `pairplot()`. Seaborn<sup>9</sup> is a Python statistical data visualization library build on top of the popular visualization library `matplotlib`. Essentially, it gives us a simple API for rendering many common statistical plots. `pairplot` takes a `DataFrame` as argument and returns a `PairGrid` plot of the pairwise relationships between each of the variables in the passed `DataFrame`.

Note that the central diagonal axis is treated differently than others. It shows the univariate distribution plot of each feature.

Note that it the plots above the central diagonal axis are symmetric to the plot below the central diagonal axis and it is only necessary to examine one set of plots.

```
sns.pairplot(seeds_df);
```

### Discussion: Interpretation of the PairGrid Visualization

We make note of a few things in the `PairGrid`.

There are several highly linear relationships in this dataset, for example, `area-perimeter`, as shown here:



Note that there are other strongly linear relationships that are not as linear as `area-perimeter`, but are still strongly linear, e.g. `length_of_kernel-width_of_kernel`, as shown here:



---

<sup>9</sup><https://seaborn.pydata.org>

While, there does not appear to be a clear linear separation between the classes there is some indication that we will be able to separate the classes by feature variables. Here, we can see that **area** may be a feature we can use to separate the classes:



Discuss that these linearities can help us to choose a model.

## Modeling for Inference

As you continue the exploration of this data set, you will use a machine learning model for the purposes of EDA, in this case a simple single variable classification model using Logistic Regression. In *Introduction to Statistical Learning with R*, James et al. propose that an inferential investigation of a data set may be interested in the following:

- Which predictors are associated with the response?
- What is the relationship between the response and each predictor?
- Can the relationship between  $Y$  and each predictor be adequately summarized using a linear equation, or is the relationship more complicated?

Here, you use single-variable logistic regression to assess the performance of each feature individually as a predictor of **seed\_class**. This is a simple technique for making inferences about the strength of associations between each feature and the target class.

You will use Scikit-Learn for a simple implementation of the complex mathematical model, Logistic Regression. Logistic Regression is a regression model that creates the best linear separator of classes in a model. In other words, the model you build will find the line or lines that best separates each target class from the others.

### Create features and target Subsets of Data Set

First, you create two subsets of the complete data set: 1. **features**, a **DataFrame** that includes everything but the target column 2. **target**, a **Series** or Pandas Column, that is just the target column

```
features = seeds_df.drop('seed_class', axis=1) target = seeds_df.seed_class
```

### Programming Method: Classes and Objects

For those who are not familiar with object-oriented programming (OOP), this next step may be a bit of a challenge. Using the OOP paradigm, classes of objects are defined in code. To use a class, we must instantiate it as an object. To make this a bit more concrete, you might think of the class, Toyota Prius. This class is a kind of car, the definition of which exists in blueprints stored by the Toyota Corporation. The Prius that you drive, however, is not the class, but an object of the class Toyota Prius. It conforms to the definition as defined

in the blueprints, but it is a specific instance of that class that you are driving around.

You do the same thing here with `lr_model` and `LogisticRegression()`. `LogisticRegression` is a class defined in the `sklearn.linear_model` library. `lr_model` is an object of class `LogisticRegression` that we will use to do our work.

### Import the LogisticRegression Class and instantiate a LogisticRegression Model

We import the `LogisticRegression` class from the `sklearn.linear_model` library and instantiate the `lr_model` object.

```
from sklearn.linear_model import LogisticRegression
lr_model = LogisticRegression()
```

### Programming Method: for-loop Over Features

Using a `for`-loop you iterate through each feature and use it to fit a Logistic Regression model on the data. You then score each of these models using Accuracy (the default scoring method for Scikit-Learn's classification models). Finally, each result is appended to a list called `lm_scores`. After running this cell, `lm_scores` will be a list of dictionaries containing the results of each test.

### Programming Method: Fit and Score

With `lr_model` you are fitting and scoring a model for each step in the `for`-loop. This is a common technique with Scikit-Learn predictive models. You perform the `model.fit()` function on the data. This will execute the particular learning algorithm for your chosen model and store the fit model. You can then use the `model.score()` to score a particular fit model.

For a Supervised Learning model, both fitting and scoring take both features and target as an argument e.g. `model.fit(X, y)` and `model.score(X, y)`. Here, `X` is the features and `y` is the target. Note that any combination of features and target can be used. Care must be taken to make sure that the indices of `X` and `y` match up e.g. that the values of the *ith* row of `X` match the target in the *ith* row of `y`. Additionally, it is required that `X` and `y` have the same number of rows.

During scoring, the interim step of generating predictions is performed. In order to assess the model performance, a vector of `yhat` values is generated internally and compared to the actual values, i.e. the ground truth, stored in the `y` vector.

```
lm_scores = []
for feat in features.columns:
    features_subset = features[[feat]]
    lr_model.fit(features_subset, target)
    lm_scores.append( {'feature': feat, 'score': lr_model.score(features_subset, target)} )
```

**Programming Method: Manually Build a DataFrame**

A list of dictionaries is a useful structure for storing information, but it is not an ideal way for displaying information. This can be seen by displaying `lm_scores`.

```
lm_scores
```

Even this simple set of results becomes unwieldy.

Fortunately, it is trivial to convert a list of dictionaries with the same keys to a pandas **DataFrame**. We know that all of the dictionaries in `lm_scores` have the same keys, but we can confirm this with a simple Python list comprehension.

```
[d.keys() for d in lm_scores]
```

Provided that each dictionary has the same keys, each key will become the name of a column in the created **DataFrame** with the values of each dictionary forming a row.

```
lm_scores_df = pd.DataFrame(lm_scores) lm_scores_df
```

To make interpretation even easier, you perform two transformations on the **DataFrame**. First, you assign the `feature` column to the **Index**, or row names, of the **DataFrame** using the `lm_scores_df.set_index()` function. The `inplace=True` argument specifies that this change is to be written to the **DataFrame** stored in memory. Second, you sort the `score` column in descending order using the `lm_scores_df.sort_values(ascending=False)` function. Again, the `inplace=True` argument specifies that this change is to be written to the **DataFrame** stored in memory.

```
lm_scores_df.set_index('feature', inplace=True) lm_scores_df.sort_values('score',
ascending=False, inplace=True)
```

**Discussion: Interpreting Coefficients of Single Variable Logistic Regression**

From these results, you can see that `area` is the strongest predictor as a single variable, followed by `width_of_kernel`. It is worth noting that these scores are using the features as a single variable predictor. The order of feature strength may be significantly different if using a multi-variable model.

```
lm_scores_df
```

## T-SNE Model for Visualization of High-Dimensional Data

Given a data set with  $p = 7$ , you can imagine that it will be difficult to visualize this data sets using conventional means. Humans are accustomed to looking at data in only two dimensions. We are able to visualize in three dimensions, but three-dimensional visualization is already significantly more difficult than two dimensional. And beyond that, we are out of luck. Luckily for you, there is a fairly simple, at least in implementation, method for visualizing high-dimensional data using a model called T-SNE (often pronounced "tease knee"). The math behind this technique<sup>10</sup> can be fairly involved, but for your purposes it suffices to understand that you are creating a two-dimensional representation of our seven-dimensional features.

You will be building a T-SNE model over your data set simply for helping to prepare a simple visualization of our data. In other words, you are again building a machine learning model solely to help with the EDA process.

The implementation of takes only a few lines of code using the Python library, Scikit-Learn<sup>11</sup>.

### Import the TSNE Class

Import the model that you need from Scikit-Learn. The TSNE model is in the `sklearn.manifold` library.

```
from sklearn.manifold import TSNE
```

### Instantiate a TSNE Model

Next, you create a `TSNE()` object named `two_dim_model`. Here, you instantiate a TSNE object with 2 components.

```
two_dim_model = TSNE(n_components=2)
```

Note, that if you request the type of the `two_dim_model` object it tells you that it is a TSNE object from the `sklearn.manifold` library.

```
type(two_dim_model)
```

### Fit the Model and Transform the Data

Next, you use the `two_dim_model` object to transform your data. The transformation that you are performing is to project the original data from seven dimensions to two dimensions. Note that in this case you are not interested in

---

<sup>10</sup><http://colah.github.io/posts/2014-10-Visualizing-MNIST/>

<sup>11</sup><http://scikit-learn.org/stable/>



overwriting the original data with this transformation. For this reason, you save the transformed data as a new variable, `features_2d`.

### Programming Method: Fit and Transform

With `two_dim_model` you performed the "fit" and the "transformation" in one step using the method `.fit_transform()`. In this, `two_dim_model` learned the nature of the data, that is, it was "fit" to the data. Then we used the fit model to "transform" the data, resulting in the two-dimensional `np.array`, `featured_2d`.

```
features_2d = two_dim_model.fit_transform(features)
```

If you display the shape of each of the objects, indeed we can see that `features` has  $p = 7$ , while `features_2d` has  $p = 2$ .

```
features.shape, features_2d.shape
```

Scikit-Learn Transformation objects always return `np.array` objects. It is best to continue to work with a Pandas `DataFrame`. For this reason, you turn `features_2d` into a `DataFrame`. You give the columns of the new `DataFrame` the names `Component_1` and `Component_2`. This new `DataFrame` will be used in the next step to prepare a two-dimensional scatter plot by class.

```
features_2d = pd.DataFrame(features_2d) features_2d.columns = ['Component_1', 'Component_2']
```

### Visualization: Scatter Plot by Class

The file `visualization.py` includes a helper function `scatter_plot_by_class`. As the name suggests, this function returns a scatter plot of our data with each point labeled by its class using color.

### Programming Method: Import Functions From File

You have previously imported functions and classes from various libraries available to you using your local Python installation. Here, for the first time, you import a function from a project file, `visualization.py`. This can be done by simply creating a file with the `.py` suffix in the current working directory. In this way, the Python interpreter will be able to easily locate the file and is able to import any functions or classes from that file as if they were part of the installed and configured Python libraries.

**Import** `scatter_plot_by_class`

```
from visualization import scatter_plot_by_class
```

### Programming Method: A Simple Method for Controlling Plots

For most of the plots you will be generating, you will begin the cell with the line

```
' fig, ax = plt.subplots(NROWS, NCOLS, figsize=(XDIM,YDIM)) '
```

This gives a simple and easily repeatable pattern for controlling the appearance of rendered plots within a notebook. Most often, we are interested in controlling two aspects of our plots: 1) the dimension of the plot as rendered and 2) how many and the arrange of subplots to include. Here, you use this pattern in its most trivial sense creating a "set of subplots" with one row and one column. This is a bit of overkill for this simple usage. The benefit is in using a single uniform and easily remembered pattern for working with plots and subplots.

When using this pattern, you will most often work with the `ax` object returned by the function. The `ax` object contains references to the individual subplots in a figure. In the case where a figure is created with 1 row and 1 column of plots i.e. a single plot, the `ax` object is simply a reference to that plot.

```
fig, ax = plt.subplots(1,1,figsize=(20,5)) scatter_plot_by_class(ax, features_2d, target)
```

### Discussion: Two-Dimensional Scatter Plot by Class

Note: TSNE is a randomized-process and your results may vary slightly from the results shown below.

The two-dimensional representation of the data suggests that you may be able to easily classify this data. You can see clear linear separations in the groupings of each of the three seed types. These linear separations indicate that a linear model, such as Logistic Regression, may perform well as a classification model on this data set.



You can also note clear shapes in the data and that class relationships appear to have some locality. This indicates that an instance-based model like K-nearest Neighbors may perform well as a classification model on this data set.



### Visualization: Cluster Model

To continue your exploration of this data set, you next prepare a cluster model. If you are able to create a strong cluster model, this will reinforce the idea that K-Nearest Neighbors is a good classification model for this data set.

To do this you follow a similar pattern as before using the Scikit-Learn library.

You can obtain the cluster model, `KMeans` from the `sklearn.cluster` library. `KMeans`<sup>12</sup> is amongst the most straightforward models in terms of interpretation. Using this model, you seek  $k$  clusters based upon the notion of the most unique *centroids* or mean values over the dataset, thus `KMeans`. As with `TSNE`, the implementation of this can be fairly complex, and you will simply rely on Scikit-Learn to implement this model.

### Import the KMeans Class

```
from sklearn.cluster import KMeans
```

### Instantiate a KMeans Model

After importing the model, you instantiate `cluster_model`, an object of class `KMeans` with `n_clusters` set to 3.

```
cluster_model = KMeans(n_clusters=3)
```

### Fit the Model

With `cluster_model` you will not actually transform the data. Rather you will have the model learn from the data and then use the information learned. `cluster_model` will learn the three most likely cluster centroids over the data set and then use these to create new labels for the data. You will then compare these new labels to the original labels obtained from the UCI Machine Learning Repository stored as `target`.

```
cluster_model.fit(features)
```

### Access Cluster Labels from the Fit Model

Having fit the model you access the labels created during the fit using the model attribute `.labels_`.

```
cluster_labels = cluster_model.labels_
```

Finally, you prepare a plot showing both the original levels and the new labels generated by the cluster model. As expected, the cluster model does a decent job of labeling the data however as expected its struggles some of the boundaries between the different clusters.

---

<sup>12</sup><http://stanford.edu/~cpiech/cs221/handouts/kmeans.html>

### Visualization: Scatter Plots by Class and Cluster

Note that here you use the figure preparation code `fig, ax = plt.subplots(2,1,figsize=(20,10))`. This time, you specify that the figure should have two rows and one column of subplots. As a result, the `ax` object is a list containing two references, one to each of the two subplots in the figure. Note that you refer to these subplots using `ax[0]` and `ax[1]`.

```
fig, ax = plt.subplots(2,1,figsize=(20,10))
scatter_plot_by_class(ax[0], features_2d, target) scatter_plot_by_class(ax[1],
features_2d, cluster_labels, label=False)
ax[0].set_title('Actual Labels') ax[1].set_title('Cluster Model Labels');
```

### Discussion: Scatter Plots by Class and Cluster

A visual inspection of the labels generated by the cluster model as compared to the original labels shows that the cluster model performs very strongly. The cluster model makes some error particularly at separation boundaries, but in general is able to identify the proper clustering of points.

It is of note that the cluster labels themselves are simply numerical references to the stored centroids. In other words, a label of 0 signifies that the point should be associated with the first stored mean. The cluster labels should not be taken to mean an association with a particular species of wheat. The cluster model is an unsupervised model and has no notion of the actual labels. Recall that during the fitting process

```
' cluster_model.fit(features) '
the model was not passed the labels at all.
```

It is also of note that while `features_2d` is used for plotting the labels returned by the cluster model, the model itself was fit using all seven features.

### Plotting the Decision Boundary

During the **Modeling for Inference** phase you noted that the two strongest predictive features when used individually were `area` and `width_of_kernel`. You use this knowledge to prepare the final visualization of this Case Study, a two-dimension scatter plot with a decision boundary. This plot takes a predictive model and uses the model to prepare a colored areas representing partitions of the feature space that should be considered indicative of a given class.

#### Import scatter\_plot\_with\_decision\_boundary

As before you import `scatter_plot_with_decision_boundary` from the local Python file `visualization.py`.

```
from visualization import scatter_plot_with_decision_boundary
```

**Import the KNeighborsClassifier Class and instantiate a KNeighborsClassifier Model**

You will be preparing two decision boundary plots, one for a Logistic Regression model and the other for a K-Nearest Neighbors model. It is necessary to import `KNeighborsClassifier` from the `sklearn.neighbors` module.

```
from sklearn.neighbors import KNeighborsClassifier
knn_model = KNeighborsClassifier(n_neighbors=4)
```

**Visualization: Scatter Plots with Decision Boundary**

```
fig, ax = plt.subplots(2,1,figsize=(20,15))
scatter_plot_with_decision_boundary(ax[0], features, target, 'area', 'width_of_kernel',
lr_model) scatter_plot_with_decision_boundary(ax[1], features, target, 'area', 'width_of_kernel',
knn_model)
ax[0].set_title("Decision Boundary Logistic Regression Area - Width of Kernel") ax[1].set_title("Decision Boundary K Nearest Neighbors Area - Width of Kernel");
```

**Discussion: Scatter Plots with Decision Boundary**

The scatter plots with decision boundary show that both models, each using only `area` and `width_of_kernel` to a reasonably good job of capturing the class relationship in the data. You noted some linear separation in the data class relationships and indeed the model is able to encode this reasonably well. It is of note that this linear separation breaks down in a few notable locations. Either of these linear separations, for example, would have returned fewer errors:



## 2.4 Building a Predictive Model

```
from sklearn.metrics import accuracy_score
```

```
from sklearn.neighbors import KNeighborsClassifier
```

```
from sklearn.model_selection import train_test_split
```

```
knc = KNeighborsClassifier()
```

```
X_train, X_test, y_train, y_test = train_test_split(features, target, random_state=10)
```

```

scores = [] for i in range(2,20): knc = KNeighborsClassifier(i) knc.fit(X_train,
y_train) train_score = knc.score(X_train, y_train) test_score = knc.score(X_test,
y_test) scores.append( 'k' : i, 'train_score' : train_score, 'test_score' : test_score )
scores_df = pd.DataFrame(scores) scores_df.set_index('k', inplace=True) scores_df

```

```

scores_df.plot()

```

```

scores = [] for i in range(1,10): knc = KNeighborsClassifier(i*2) knc.fit(X_train,
y_train) train_score = knc.score(X_train, y_train) test_score = knc.score(X_test,
y_test) scores.append( 'k' : i*2, 'train_score' : train_score, 'test_score' : test_score
) scores_df = pd.DataFrame(scores) scores_df.set_index('k', inplace=True) scores_df

```

```

scores_df.plot()

```

```

best_features = list(lm_scores_df.index)

```

```

lr = LogisticRegression() knc = KNeighborsClassifier(2)
X_train_df = pd.DataFrame(X_train) X_train_df.columns = features.columns
X_test_df = pd.DataFrame(X_test) X_test_df.columns = features.columns
final_performance_scores = [] for i in range(len(best_features)): knc = KNeigh-
borsClassifier(2) feats = best_features[i+1] lr.fit(X_train_df[feats], y_train) knc.fit(X_train_df[feats],
y_train)
final_performance_scores.append( 'features' : ' '.join(feats), 'feature_added' :
best_features[i], 'lr_train_score' : lr.score(X_train_df[feats], y_train), 'lr_test_score'
: lr.score(X_test_df[feats], y_test), 'knc_train_score' : knc.score(X_train_df[feats],
y_train), 'knc_test_score' : knc.score(X_test_df[feats], y_test) )
final_performance_scores_df = pd.DataFrame(final_performance_scores) final_performance_scores_df.se
inplace=True) final_performance_scores_df

```

```

final_performance_scores_df.plot(rot=45, figsize=(20,7))

```

# Chapter 3

## Configuring AWS

In this first chapter, you will configure your local system and create an AWS instance in order to do data science work. To do this work, you will use the containerization technology Docker in order to streamline the work of configuring and provisioning your data science platform. You will use Docker Compose to manage the platform and its two containers and rely on Amazon Web Services (AWS) to manage your hardware.

If you have your own Jupyter Notebook setup previously installed, feel free to use it, in which case, it is not necessary to work through this first chapter.

### 3.1 Amazon Web Services

If you have not already done so, set up an AWS account<sup>1</sup>.

You will be using AWS to manage the hardware upon which your data science platform will run. We will leave the details of what exactly "hardware" means to AWS. This is to say that AWS may be allocating resources as a virtual machine, but for your purposes, the experience will be as if you are using a physical system across the room from you.

The most popular service offered by Amazon Web Services is the Elastic Compute Cloud (EC2), "a web service that provides secure, resizable compute capacity in the cloud"<sup>2</sup>. For our purposes, compute capacity means a cloud-based computer you will use to run your platform.

---

<sup>1</sup>As of 2018/08/25, detailed instructions for doing this can be obtained here: <https://aws.amazon.com/premiumsupport/knowledge-center/create-and-activate-aws-account/>.

<sup>2</sup><https://aws.amazon.com/ec2/>

Readers new to AWS will be able to work through this text using the AWS Free Tier<sup>3</sup>. For the first 12 months following sign up, new users receive 750 Hours per month of EC2 time. This amounts to 31.25 days of availability and, provided that readers keep only one server running at a time, ensures that readers can work through this text at no cost.

## 3.2 Configure your Local System

In this first chapter, with the exception of some AWS system configuration in the browser, all of the work that you will be doing will take place at the command line. Work at the command line is done via a special type of application called a shell. A shell is a user interface that provides access to the operating system of a computer. We prefer the shell to other modes of working with computers because of its simplicity.

The most popular shell is the Bourne Again Shell or Bash. If you are using a Mac OS X system or a Linux system, you will already have Bash available to you in an application called Terminal. If you are using a Windows system, we recommend the use of Git-Bash<sup>4</sup>. We have no preference as to the settings used when configuring Git-Bash. The default settings are fine when installing the program. If you are using a Chromebook, we recommend the use of Termius, available from the Chrome Web Store.

### SSH Keys

All of the work that you will be doing will take place remotely. As such, there is very little configuration to be done for the local system. The one thing that you will need to do is configure a set of SSH Keys to enable secure connection to the remote system you bring online.

---

<sup>3</sup><https://aws.amazon.com/free/>

<sup>4</sup><https://gitforwindows.org>



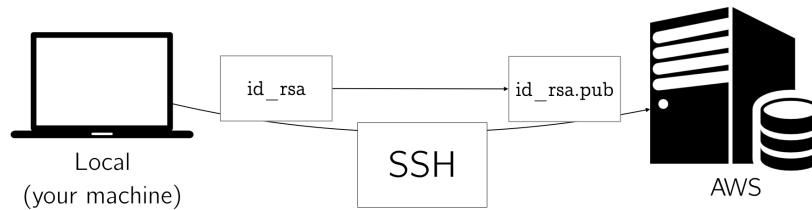


Figure 3.1: Connecting with SSH Keys

An SSH Key is a password-less method of authenticating to a remote system using public-key cryptography. Authentication is done using a key pair consisting of a public key, which can be shared publicly, and a private key, which is known only to the user (See Figure 3.1). One might think of the public key as the lock on your front door, accessible to anyone, and the private key as the key in your pocket so that only you are able to open your door and gain access to your home.

You will generate this key pair on our local system and then provide the public key to AWS so that it can be added to any system you wish to launch. You will keep the private key on our local system and use it whenever you wish to gain access.

### 3.3 Create a New Key Pair

You will use the Bash tool `ssh-keygen` to create a new key pair. To begin open a new terminal session<sup>5</sup>, where you will examine whether or not you already have a key pair. Launching a new Bash session will put you in your home directory.

The canonical location for storing SSH Keys is in a folder called `/.ssh` in our home directory<sup>6</sup>. Note that this directory begins with a `.` which makes it a

<sup>5</sup>On a Mac or Linux system, simply open the Terminal application. On Windows, open Git-Bash. On Chromebook, open Termius.

<sup>6</sup>On every system that we will use, the `sim` symbol is used as an alias for your home directory. The location of the actual home directory will vary by system. Ubuntu users' home directory will be `/home/username`. Mac OS X users' home directory will be `/Users/username`. Windows/Git-Bash users' home directory will be `/c/Users/username`.

hidden directory. In Listing 3.1, you use `cd` to navigate to our home directory and `ls -la` to display all of the contents of our home directory in a list.

home: On every system that we will use, the *sim* symbol is used as an alias for your home directory. The location of the actual home directory will vary by system. Ubuntu users' home directory will be `/home/username`. Mac OS X users' home directory will be `/Users/username`. Windows/Git-Bash users' home directory will be `/c/Users/username`.

*square* **Note:** Occasionally in code listings, I will truncate the output. If you see `...` in the listing, this should be taken to mean that there is additional output generated that is not important for the discussion.

Listing 3.1: List the contents of our home directory

```
$ cd ~
$ ls -la
...
drwxr-xr-x  21 joshuacook  staff    714 Jul  31   2017 .pylint.d
drwx-----  11 joshuacook  staff    374 Jan  28  18:49 .ssh
drwxr-xr-x   6 joshuacook  staff    204 Feb   2  22:01 .vim
-rw-----   1 joshuacook  staff  20788 Feb  10  08:54 .viminfo
-rw-r--r--@   1 joshuacook  staff   1263 Jul  26   2017 .vimrc
drwx-----@   5 joshuacook  staff    170 Aug  26  09:12 Applications
drwx-----+  19 joshuacook  staff    646 Feb  11  09:32 Desktop
drwx-----+   6 joshuacook  staff    204 Feb   4  12:18 Documents
...
```

My local system is running Mac OS X and has the `.ssh` folder already. As the directory already exists, in Listing 3.2, I list the contents of my `.ssh` directory.

Listing 3.2: Display contents of the `.ssh` directory

```
$ ls -la .ssh
total 96
drwx-----  11 joshuacook  staff    374 Jan  28  18:49 .
drwxr-xr-x+  74 joshuacook  staff   2516 Feb  10  08:54 ..
-rw-----   1 joshuacook  staff   1679 Jan  11  16:21 id_rsa
-rw-r--r--   1 joshuacook  staff    418 Jan  11  16:21 id_rsa.pub
```

As can be seen, I already have an SSH Keypair named `id_rsa` and `id_rsa.pub`. If this is true for you, as well, you should skip the next step and not create new SSH Keys (See Listing 3.3).

If when listing the home directory, you do not see a folder called `.ssh` or when displaying the contents of `.ssh` you do not see an SSH Keypair named `id_rsa`

and `id_rsa.pub`, a new SSH Keypair will need to be created. In Listing 3.3, you create a new SSH Keypair using the `ssh-keygen` command line utility.

During the creation of the SSH Keypair, you will be prompted three times. The first asks where you should save the SSH Keypair, defaulting to the `.ssh/id_rsa` in our home directory. In Listing 3.3, you see that this is being done at `/Users/joshuacook/.ssh/id_rsa` on my local system where my username is `joshuacook`. The second and third prompts will ask for a passphrase to be added to the key. For our purposes, leaving this passphrase empty will be fine. In other words, the default options are preferable and you may simply hit `<ENTER>` three times.

Listing 3.3: Create a new SSH Keypair

```
$ ssh-keygen
Generating public/private rsa key pair.
Enter file in which to save the key
    (/Users/joshuacook/.ssh/id_rsa):
Created directory '/Users/joshuacook/.ssh'.
Enter passphrase (empty for no passphrase):
Enter same passphrase again:
Your identification has been saved in
    /Users/joshuacook/.ssh/id_rsa.
Your public key has been saved in
    /Users/joshuacook/.ssh/id_rsa.pub.
The key fingerprint is:
SHA256:MkCnhaAzjcjRHCUC/pdUsxrk7be6+gNXTrEchJyF0qs joshuacook@LOCAL
The key's randomart image is:
+---[RSA 2048]---+
| .==oo..o*o    |
| o+o=o+o o=oo  |
| *...o  +.o. +  |
| o   ...o=   =  |
|   .o+S.+      |
|   . = ....    |
|   . o .       |
|   E  ..       |
|   .o+o        |
+---[SHA256]-----+
```

You can verify the SSH Keypair you just created by displaying the Public Key in our shell (Listing 3.4). Here, you use the `cat` command, which concatenates the contents of `id_rsa.pub` to the shell output.

Listing 3.4: Display Public SSH Key

```
$ cat ~/.ssh/id_rsa.pub
ssh-rsa
AAAAB3NzaC1yc2EAAAADAQABAAQDdnHPEiq1a40sDDY+g9luWQs8pCjBmR
64MmsrQ9MaIaE5shIcFB1Kg3pGwJpypiZjoSh9pS55S9LckNsBfn8Ff42ALLj
R8y+WlJKVvK/ODvDXGVCcC0t/uTvVxVx0bRruYxLW167J89UnxnJuRZDLeY9fD
OfIzSR5eglhCWVqi0zB+OsLqR1W04Xz1oStID78UiY5msW+EFg25Hg1wepYMC
JG/Zr43By0YPGseUrbCqFBS1K1QnzfWRfEKHZbtEe6HbWwz1UDL2NrdFXxZAI
XYYoCVt14WXd/WjDwSjbMmtf3BqenVKZcP2DQ9/W+geIGGjv0TfUdsCHennYI
EUfEEP joshuacook@LOCAL
```

### 3.4 Configure your AWS Account

That is the sum of the local configuration you will need to do in order to get started. The next thing you will need to do is configure our AWS Account. To do this, you will need to configure a Key Pair corresponding to the SSH Keypair on your local system. The AWS Key Pair is slightly misnamed as it is not in fact a pair, but rather is simply the public portion of the SSH Keypair you have on our local system.

To begin, log in to your AWS control panel and navigate to the EC2 Dashboard (Figure 3.2). First, access “Services” (Figure 3.2, #1) then access “EC2” (Figure 3.2, #2). The Services link can be accessed from any page in the AWS website.

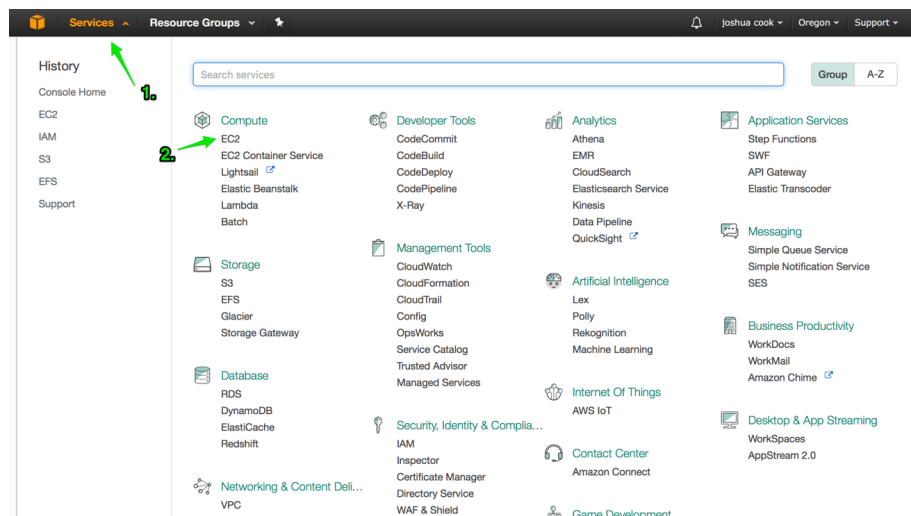


Figure 3.2: Access EC2 Dashboard

## Add the Public Key to AWS

Once at the EC2 control panel, access the Key Pairs pane using either link (Figure 3.3).

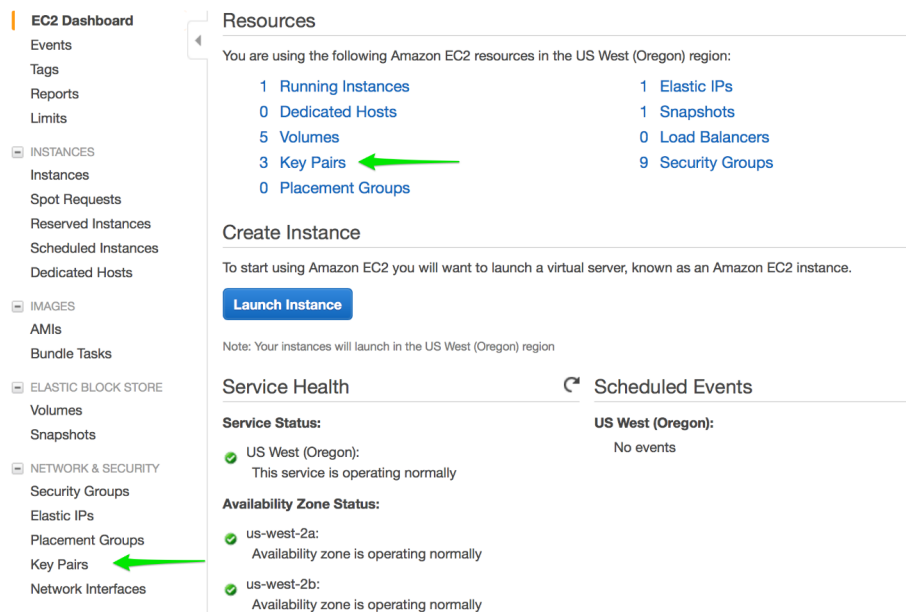


Figure 3.3: Access Key Pairs in the EC2 Dashboard

From the Key Pairs pane, choose “Import Key Pair.” This will activate a modal that you can use to create a new key pair associated with a region on your AWS account. Make sure to give the key pair a computer-friendly name, like `from-MacBook-2018`. Paste the contents of your public key (`id_rsa.pub`) into the public key contents. Prior to clicking Import, your key should appear as in Figure 3.4. Click Import to create the new key.

Import Key Pair

X

Click Browse and navigate to your public key. You may change the name of your key if necessary. Alternatively, you can copy and paste the contents of your public key into the dialog.

Load public key from file

Choose File no file selected

Key pair name

my\_pc\_2017

Public key contents

```
ssh-rsa
AAAAB3NzaC1yc2EAAAADAQABAAQDdnHPEiq1a4OsDDY+g9luWQS8pCjBmR64MmsrQ9
MalaE5shlcFB1Kg3pGwJypipZjoSh9pS55S9LckNsBfn8Ff42ALLjR8y+WlJKVh/0DvDXgGVcCc0t
/uTvXVx0bRruYxLW167J89UnxnJuRZDLey9fDOflzSR5eglhCWVqiOzB+OsLqR1W04Xz1oStlD7
8UiY5msW+EFg25Hg1wepYMCJG/Zr43ByOYPGseUrbCqFBS1KIQnzfWRfEKHZbtEe6HbWwz1
UDL2NrdFXxZAlXYoCVtl4WXd/WjDwSjbMmtf3BqenVKZcP2DQ9/W+geIGGjvOTfUdsCHennYI
EUfEEP ubuntu@ip-172-31-43-19
```

Cancel

Import

Figure 3.4: Import a New Public Key

You have created a key pair between AWS and your local system. When you create a new instance, you will instruct AWS to provision the instance with this public key and thus you will be able to access the cloud-based system from your local system using your private key.

## 3.5 Launch a New EC2 Instance

To create a new instance, start from the EC2 Dashboard and click the Launch Instance button (Figure 3.5).

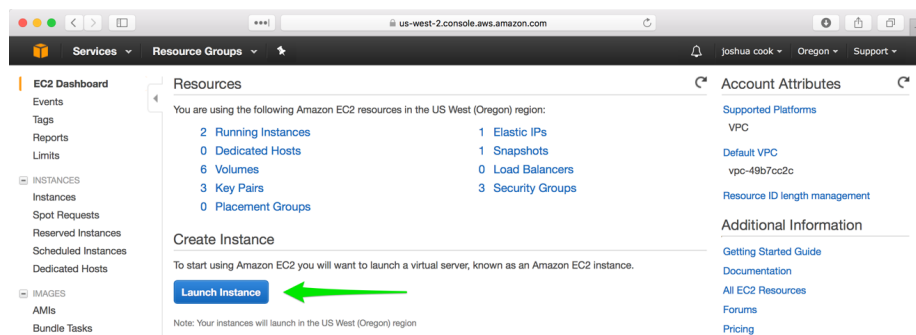


Figure 3.5: Begin the launch process for a new instance

## Step 1: Choose an Amazon Machine Image (AMI)

The launching of a new instance is a multi-step process that walks the user through all configurations necessary. The **first tab** is “Choose AMI.” An AMI is an Amazon Machine Image<sup>7</sup>, and contains the software you will need to run your sandbox machine. I recommend choosing the latest stable Ubuntu Server release that is free-tier eligible. At the time of writing, this was ami-efd0428f, Ubuntu Server 16.04 LTS (HVM), SSD Volume Type (Figure 3.6).

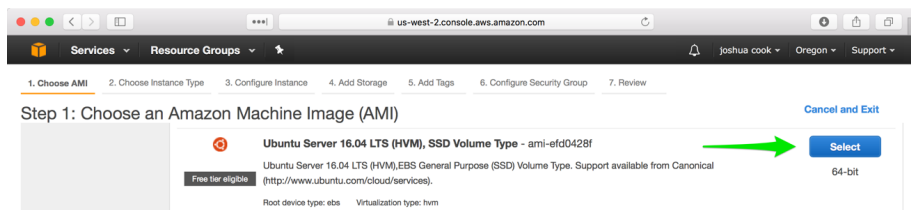


Figure 3.6: Choose the latest stable Ubuntu Server release as AMI

## Step 2: Choose Instance Type

The **second tab** is “Choose Instance Type.” In practice, I have found that the free tier, **t2.micro** (Figure 3.7), is sufficient for many applications. Furthermore, the instance type may always be changed later should the need present itself.

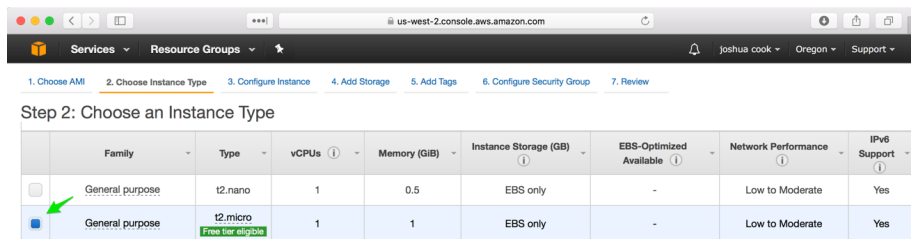


Figure 3.7: Choose t2.micro for Instance Type

## Step 3: Configure Instance Details

The **third tab**, “Configure Instance,” can be safely ignored.

## Step 4: Add Storage

The **fourth tab** is “Add Storage.” This option is also specific to intended usage. It should be noted that Jupyter Docker images can take up more than

<sup>7</sup><http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/AMIs.html>

5GB of disk space in the local image cache. For this reason, it is recommended to raise the value from the default 8GB to 30GB. Furthermore, as noted on this tab:

Free tier eligible customers can get up to 30 GB of EBS General Purpose (SSD) or Magnetic storage.

## Step 5: Add Tags

The fifth tab, “Add Tags,” can be safely ignored.

## Step 6: Configure Security Group

The sixth tab, “Configure Security Group,” is critical for the proper functioning of your systems. By default this tab will be set up to “Create a **new** security group”. This will not work for us! Ultimately, we will be accessing our system via a web browser which we require at a minimum that port 80 is open. We recommend simply using the default group which will open our system on all ports. If greater security is required for your specific application a more restrictive security group may be defined and used.

Select the “default” security group (Figure 3.8).



Figure 3.8: Choose the latest stable Ubuntu Server release as AMI

**square Note:** You may receive a Warning stating, “Rules with source of 0.0.0.0/0 allow all IP addresses to access your instance. We recommend setting security group rules to allow access from known IP addresses only.” This is expected and is okay.

## Step 7: Review Instance Launch

Finally, click “Review and Launch.” Here, you see the specific configuration of the EC2 instance you will be creating. Verify that you are creating a **t2.micro** (Figure 3.9, #2) running the latest free tier-eligible version of Ubuntu Server (Figure 3.9, #1) and that it is available to all traffic (Figure 3.9, #3), and then click the Launch button (Figure 3.9, #4).



**Step 7: Review Instance Launch**  
Please review your instance launch details. You can go back to edit changes for each section. Click **Launch** to assign a key pair to your instance and complete the launch process.

**AMI Details** [Edit AMI](#)  
 1. **Ubuntu Server 18.04 LTS (HVM), SSD Volume Type - ami-1b791862**  
 Ubuntu Server 18.04 LTS (HVM), EBS General Purpose (SSD) Volume Type. Support available from Canonical (<http://www.ubuntu.com/cloud/services>).  
 Root Device Type: ebs    Virtualization type: hvm

**Instance Type** [Edit instance type](#)  
 2. 

Instance Type	ECUs	vCPUs	Memory (GiB)	Instance Storage (GiB)	EBS-Optimized Available	Network Performance
t2.micro	Variable	1	1	EBS only	-	Low to Moderate

**Security Groups** [Edit security groups](#)  

Security Group ID	Name	Description
sg-94cc68f0	default	default VPC security group

All selected security groups inbound rules

Type	Protocol	Port Range	Source	Description
All traffic	All	All	sg-94cc68f0 (default)	

3. **Instance Details** [Edit instance details](#)  
[Cancel](#) [Previous](#) [Launch](#) 4.

Figure 3.9: Review and launch the new instance

## Add an SSH Key

In a final confirmation step, you will see a modal titled “Select an existing key pair or create a new key pair.” Select the key pair you previously created. Check the box acknowledging access to that key pair and launch the instance (Figure 3.10).

**Select an existing key pair or create a new key pair** ✕

A key pair consists of a **public key** that AWS stores, and a **private key file** that you store. Together, they allow you to connect to your instance securely. For Windows AMIs, the private key file is required to obtain the password used to log into your instance. For Linux AMIs, the private key file allows you to securely SSH into your instance.

Note: The selected key pair will be added to the set of keys authorized for this instance. Learn more about [removing existing key pairs from a public AMI](#).

Choose an existing key pair

**Select a key pair**

→ from-MacBook-2017

→ ☒ I acknowledge that I have access to the selected private key file (from-MacBook-2017.pem), and that without this file, I won't be able to log into my instance.

[Cancel](#) [Launch Instances](#) ←

Figure 3.10: Add a key pair to the instance

*square* **Note:** If this step is not done correctly, that is, if the correct key pair is not added to the launching instance, the instance will need to be terminated and a new instance will need to be launched. There is now way to add a key pair to a running instance.

You should see a notification that the instance is now running. Click the View Instances tab in the lower right corner to be taken to the EC2 Dashboard Instances pane, where you should see your new instance running.

## Examining the Newly Launched Instance

Make note of the IP address of the new instance (Figure 3.11).

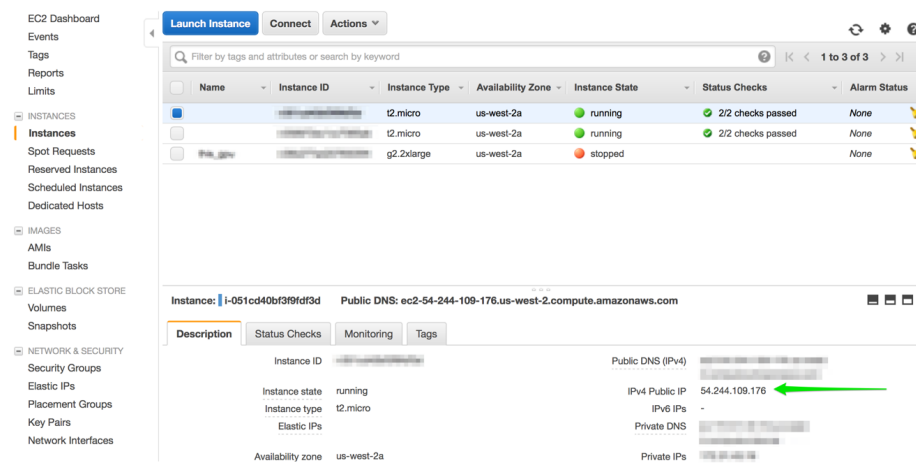


Figure 3.11: Note the IP address of the new instance

## 3.6 Git and Github

As you work through this text, you will be developing a series of data science projects. Tracking software development work is typically done using version control software. One of the most popular version control tools is **git**. Additionally, it can be useful to use a version control hosting service as a remote backup for work being tracked using **git**. The remote service we will use is Github.com. In my experience, learners who are new to version control often confuse **git** and Github, so it bares repeating – we will use **git** to track changes we make to our code and Github as a remote backup for these changes.

## Configuring Github

We will assume that you have a Github account. Once this has been done, you will need to configure an SSH connection between AWS and Github. This next part may potentially create a confusion. We are actually going to need a new SSH Keypair, this one associated with our AWS instance. This is because it is our AWS instance that will be connecting to Github, not our local machine (Figure See 3.12).

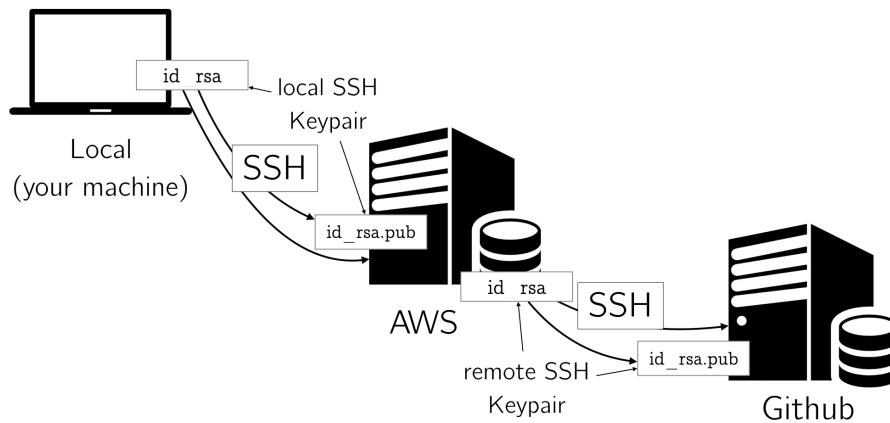


Figure 3.12: SSH Connections

## Create a New Key Pair

In Listing 3.6, you create a new key pair on your remote AWS instance. In Listing 3.5, you connect to your new AWS instance. To do this we will use the IP address we made note of in Figure 3.11. We use SSH to connect to our remote AWS instance. Note that we use the username, **ubuntu**, the default username for the Ubuntu 16 AMI provided by AWS.

Listing 3.5: SSH into New Instance

```
$ ssh ubuntu@54.244.109.176
Welcome to Ubuntu 16.04.2 LTS (GNU/Linux 4.4.0-64-generic x86_64)
```

*square* **Note:** The first time you access your EC2 instance, you should see the following message: The authenticity of host '54.244.109.176 (54.244.109.176)' can't be established ... Are you sure you want to continue connecting (yes/no)? This is expected. You should hit <ENTER> to accept or type yes and hit <ENTER>.

In Listing 3.6, you create a new SSH Keypair on our remote AWS instance. Again, during the creation of the SSH Keypair, you will be prompted three times. The first asks where you should save the SSH Keypair, defaulting to the `.ssh/id_rsa` in our home directory. In Listing 3.6, you see that this is being done at `/home/ubuntu/.ssh/id_rsa`<sup>8</sup>. The second and third prompts will ask for a passphrase to be added to the key. For our purposes, leaving this passphrase empty will be fine. In other words, the default options are preferable and you may simply hit `<ENTER>` three times.

Listing 3.6: Create a new SSH Keypair

```
$ ssh-keygen
Generating public/private rsa key pair.
Enter file in which to save the key (/home/ubuntu/.ssh/id_rsa):
Created directory '/home/ubuntu/.ssh'.
Enter passphrase (empty for no passphrase):
Enter same passphrase again:
Your identification has been saved in /home/ubuntu/.ssh/id_rsa.
Your public key has been saved in /home/ubuntu/.ssh/id_rsa.pub.
The key fingerprint is:
SHA256:ZSpFpgSRgRqlQom8yVBG2dZo1tgkPQdrmUGgMXGDtRY
ubuntu@ip-172-31-43-19
The key's randomart image is:
+---[RSA 2048]-----+
|o=XBE/*..o|
|==+=0**0.|
|o++o *o. o|
|o+ . . . +|
| . S |
| . |
| |
| |
+-----[SHA256]-----+
```

As before, you can verify the SSH Keypair you just created by displaying the Public Key in your shell (Listing 3.7). Again, you use the `cat` command, which concatenates the contents of `id_rsa.pub` to the shell output.

<sup>8</sup>This should be the same for everyone now, as you should be working on an AWS `t2.micro` running an Ubuntu system where the user's name is `ubuntu`

Listing 3.7: Display Public SSH Key

```
$ cat ~/.ssh/id_rsa.pub
ssh-rsa
AAAAB3NzaC1yc2EAAAADAQABAAQDAQ896GUMgCMAIW79gwF3ojRjcUYCKUKc8b+q
iQ1ah2jtr7sOK4WRGjkt0y3lCCHO+iUK/GrzY1Y4VxCKoKJDH3G9N5UzyGhlxa/2Ah
kKxzHht1knyh/mkVGqYUhuHpXfxUQAstCFrIdp3GOMDPiko2qeJcBF7JSv1lLMbIuM
XuVU/Mzq6BU+tEogScYytmLckyEe1j8RJ+e5nBURwmkgj3UANiDzmU/1VwL1ltEpmC
Dl0el4yEXAw8yBwM3GwjahfiBThvBHpsc43HxWrkM8Yi/kdDnvsDZYxU4zhXZPsPab
UY/LfxEod9c6Sui5W8GtAfdi6krnqbzxrKt81Mradh ubuntu@ip-172-31-43-19
```

## Add the Public Key to Github

Previously, you added your local SSH public key to your AWS account. Now, you will add your AWS SSH public key to your Github account<sup>9</sup>. First, access the **Settings** for your account by clicking the profile photo in the upper-right corner of any page on Github. Next, in the user settings sidebar, select **SSH and GPG keys**. On the SSH and GPG Keys page, click **New SSH key**. On the next page, give your key a descriptive title e.g. "AWS Feb 2018" and then paste your AWS public key in to the "Key" field. Finally, click **Add SSH key** and confirm your Github password, if prompted.

## 3.7 Learning to read the Bash Prompt

During your work you will no doubt notice that an idle SSH connection may become disconnected and/or unresponsive. Should this happen, simply close the terminal session, launch a new one, and reconnect to the remote instance.

The most important thing is that you are aware of which system your current shell session is connected to. Shell prompts are designed to relay this information to you immediately. If you are new to working with Bash, you may need to train yourself to being aware of the prompt when typing. Listing 3.8 shows the default AWS Bash prompt. The information contained is the username, **ubuntu**, and the private IP address of the AWS instance. **This is not the public address you use to connect**. What is useful about this, is that we can immediately see that the user is **ubuntu**. This tells us we are connected to AWS.

Listing 3.8: The default AWS Bash prompt

```
ubuntu@ip-172-31-21-89:~$
```

Your local system will no doubt display something different (See Listing 3.9). Again, the important thing is to take note of what is displayed by the prompt

<sup>9</sup><https://help.github.com/articles/adding-a-new-ssh-key-to-your-github-account/>

and to learn to associate that prompt with the correct system. As you become a more advanced Bash user, you may wish to personalize your prompt, but for now it is imperative that you learn to read the prompt in order to always know to which system you are connected.

Listing 3.9: A local Bash prompt

```
joshuas-macbook-pro:~$
```

## 3.8 Test your SSH Connection to Github

Having added your AWS Public Key to your Github account, you should verify your SSH connection from your AWS instance. In Listing 3.10, we attempt to connect to Github via SSH. As before, we receive a message about the authenticity of the connection. Again, type **yes**, and continue. If successful, you will see a message telling you have successfully authenticated but that Github does not provide shell access.

Listing 3.10: Verify Github SSH Key

```
ubuntu@ip-172-31-21-89:~$ ssh -T git@github.com
The authenticity of host 'github.com (IP ADDRESS)' can't be
established.
RSA key fingerprint is
16:27:ac:a5:76:28:2d:36:63:1b:56:4d:eb:df:a6:48.
Are you sure you want to continue connecting (yes/no)? yes
Hi username! You've successfully authenticated, but GitHub does
not provide shell access.
```

## 3.9 Docker

Having configured our SSH connections and provisioned a new AWS EC2 instance, it is time to get to the business of building your data science platform. To do this you will use the containerization platform Docker and its Docker Compose tool. While Docker is very easy to use, it can be difficult to understand for the uninitiated. In an earlier work, *Docker for Data Science*<sup>1</sup>, I wrote:

dockerfordatascience: <https://www.apress.com/us/book/9781484230114>

<sup>1</sup> [Using Docker] we add a layer of complexity to our software, but in doing so gain the advantage of ensuring that our local development environment will be identical to any possible environment into which we would deploy the application.

It may be simpler, however, to simply think about using Docker as a way to manage a running process. Your system will be running two processes: an IPython shell and a PostgreSQL server. Were you to not use Docker, you would need to ensure that the AWS instance had all of the libraries required to run both of those processes (and keep those libraries up to date).

Instead, you will let Docker manage the processes using a container for each process. Each respective container will be run using a predefined image built using best practices and ready to run their respective process. The exchange is this: you will take on the cognitive burden of *understanding* what Docker is doing and Docker (and the Docker community) will take over the burden of making sure that your processes run.

## Docker Compose

Docker Compose is a tool built for managing an application consisting of multiple containers. Using Docker Compose, it is possible to completely define an application using a simple text file. To make this conversation less abstract, let's have a look at the `docker-compose.yml` file you will use to define your first application (See Listing 3.11).

Listing 3.11: Your Data Science Application

```
version: "3"
services:
  ipython_shell:
    image: jupyter/scipy-notebook
  database:
    image: postgres
    volumes:
      - postgres_data:/var/lib/postgresql/data
volumes:
  postgres_data
```

That's it. This simple file completely defines a fully-functioning Data Science Application. In it, we define the two services we need: `ipython_shell` and `database`. These two services are defined using the `jupyter/scipy-notebook` and `postgres` images. When we launch the application, the images will be pulled from Docker Hub into our local memory and then launched. The one other thing we do is create a data volume `postgres_data`. We will use this as the data volume for our database server so that if for some reason we have to shut our system down, we do not lose our data. The data will exist on this volume independent of the services.

*square* **Note:** Throughout this text, when discussing infrastructure, I may

casually refer to containers, services, and processes. At the risk of annoying your local site reliability engineer, you may treat these as terms as synonymous. Care should be taken, however, not to confuse services/containers/processes and images. An image defines a service, but a service should be thought of as a living and active thing. You may loosely compare the service-image relationship to the object-class relationship in Object-Oriented Programming. A service is a running container defined by an image, just like an object is an instance of a class that exists in memory.

## Installing and Configuring Docker

Installing Docker on your AWS instance is a downright trivial process. It consists of running an install script that can be obtained from Docker and then adding your user to the Docker group. In Listing 3.12, we run these two commands. First, we download the install script from <https://get.docker.com>, then immediately pipe the script into the shell (`| sh`).

*square* **Note:** It is generally considered to be a significant security vulnerability to execute arbitrary code obtained from an unknown, or untrusted source. For our purposes, the source (<https://get.docker.com>) is considered trustworthy, we are using SSL to perform the curl, and in practice this is the method I use to install Docker. Still, it may make the security minded more comfortable to `curl` the script, inspect, and then run it.



Listing 3.12: Install Docker via a Shell Script

```
$ curl -sSL https://get.docker.com/ | sh
# Executing docker install script, commit: 1d31602
+ sudo -E sh -c apt-get update -qq >/dev/null
...

Client:
 Version:      18.02.0-ce
 API version:  1.36
 Go version:   go1.9.3
 Git commit:   fc4de44
 Built: Wed Feb  7 21:16:33 2018
 OS/Arch:     linux/amd64
 Experimental: false
 Orchestrator: swarm

Server:
 Engine:
  Version:      18.02.0-ce
  API version:  1.36 (minimum version 1.12)
  Go version:   go1.9.3
  Git commit:   fc4de44
  Built:       Wed Feb  7 21:15:05 2018
  OS/Arch:     linux/amd64
  Experimental: false
...
```

When the script completes there is one last thing to be done. In Listing 3.13, you add the `ubuntu` user to the `docker` group. By default, the command line docker client will require `sudo` access in order to issue commands to the docker daemon. You can add the `ubuntu` user to the `docker` group in order to allow the `ubuntu` user to issue commands to docker without `sudo`.

Listing 3.13: Add the Ubuntu User to the Docker Group

```
$ sudo usermod -aG docker ubuntu
```

Finally, in order to force the changes to take effect, you should disconnect and reconnect to their remote system. You can achieve this by typing `exit` or `ctrl-d` and then reconnecting via `ssh` to your EC2 instance.

## Installing and Configuring Docker

Recall that regardless of your local operating system, you are working on an AWS EC2 Instance running the Linux variant, Ubuntu. As such, `docker-compose` can be installed using the instructions provided here: <https://github.com/docker/compose/releases>,

which are written specifically for Linux machines. As of the writing of this book, this consists of two steps.

In Listing 3.14, you use `curl` to retrieve the `docker-compose` binary from Github. As of the writing of this book, the latest version of `docker-compose` was 1.19.0. You should retrieve the latest version from the above url.

Listing 3.14: Retrieve `docker-compose` binary from Github

```
$ sudo curl -L
    https://github.com/docker/compose/releases/download/1.19.0/docker-compose -u
    -s'-'uname -m' -o /usr/local/bin/docker-compose
```

In Listing 3.15, we use the `chmod` utility to allow `docker-compose` to be executed (+x).

`chmod`: The unix "change mode" utility. I pronounce it "shmod".

Listing 3.15: Enable Docker Compose to be Executed

```
$ sudo chmod +x /usr/local/bin/docker-compose
```

Finally, in Listing 3.16, we check the version of `docker-compose` against what we expect to have installed.

Listing 3.16: Check Docker-Compose Version

```
$ docker-compose -v
docker-compose version 1.19.0, build 9e633ef
```