

An example machine learning notebook

Notebook by Randal S. Olson (<http://www.randalolson.com/>)

Supported by Jason H. Moore (<http://www.epistasis.org/>)

University of Pennsylvania Institute for Bioinformatics (<http://upibi.org/>)

It is recommended to view this notebook in nbviewer (<http://nbviewer.ipython.org/github/rhiever/Data-Analysis-and-Machine-Learning-Projects/blob/master/example-data-science-notebook/Example%20Machine%20Learning%20Notebook.ipynb>) for the best viewing experience.

Table of contents

1. [Introduction](#)
2. [License](#)
3. [Required libraries](#)
4. [The problem domain](#)
5. [Step 1: Answering the question](#)
6. [Step 2: Checking the data](#)
7. [Step 3: Tidying the data](#)
 - [Bonus: Testing our data](#)
8. [Step 4: Exploratory analysis](#)
9. [Step 5: Classification](#)
 - [Cross-validation](#)
 - [Parameter tuning](#)
10. [Step 6: Reproducibility](#)
11. [Conclusions](#)
12. [Further reading](#)
13. [Acknowledgements](#)

Introduction

[\[go back to the top \]](#)

In the time it took you to read this sentence, terabytes of data have been collectively generated across the world — more data than any of us could ever hope to process, much less make sense of, on the machines we're using to read this notebook.

In response to this massive influx of data, the field of Data Science has come to the forefront in the past decade. Cobbled together by people from a diverse array of fields — statistics, physics, computer science, design, and many more — the field of Data Science represents our collective desire to understand and harness the abundance of data around us to build a better world.

In this notebook, I'm going to go over a basic Python data analysis pipeline from start to finish to show you what a typical data science workflow looks like.

In addition to providing code examples, I also hope to imbue in you a sense of good practices so you can be a more effective — and more collaborative — data scientist.

I will be following along with the data analysis checklist from [The Elements of Data Analytic Style](https://leanpub.com/datastyle) (<https://leanpub.com/datastyle>), which I strongly recommend reading as a free and quick guidebook to performing outstanding data analysis.

This notebook is intended to be a public resource. As such, if you see any glaring inaccuracies or if a critical topic is missing, please feel free to point it out or (preferably) submit a pull request to improve the notebook.

License

[\[go back to the top \]](#)

Please see the [repository README file](https://github.com/rhiever/Data-Analysis-and-Machine-Learning-Projects#license) (<https://github.com/rhiever/Data-Analysis-and-Machine-Learning-Projects#license>) for the licenses and usage terms for the instructional material and code in this notebook. In general, I have licensed this material so that it is as widely usable and shareable as possible.

Required libraries

[\[go back to the top \]](#)

If you don't have Python on your computer, you can use the [Anaconda Python distribution \(http://continuum.io/downloads\)](http://continuum.io/downloads) to install most of the Python packages you need. Anaconda provides a simple double-click installer for your convenience.

This notebook uses several Python packages that come standard with the Anaconda Python distribution. The primary libraries that we'll be using are:

- **NumPy**: Provides a fast numerical array structure and helper functions.
- **pandas**: Provides a DataFrame structure to store data in memory and work with it easily and efficiently.
- **scikit-learn**: The essential Machine Learning package in Python.
- **matplotlib**: Basic plotting library in Python; most other Python plotting libraries are built on top of it.
- **Seaborn**: Advanced statistical plotting library.

To make sure you have all of the packages you need, install them with conda:

```
conda install numpy pandas scikit-learn matplotlib seaborn
```

conda may ask you to update some of them if you don't have the most recent version. Allow it to do so.

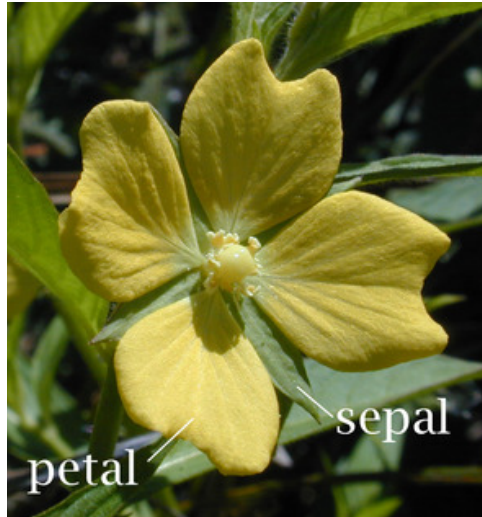
Note: I will not be providing support for people trying to run this notebook outside of the Anaconda Python distribution.

The problem domain

[\[go back to the top \]](#)

For the purposes of this exercise, let's pretend we're working for a startup that just got funded to create a smartphone app that automatically identifies species of flowers from pictures taken on the smartphone. We're working with a moderately-sized team of data scientists and will be building part of the data analysis pipeline for this app.

We've been tasked by our head of data science to create a demo machine learning model that takes four measurements from the flowers (sepal length, sepal width, petal length, and petal width) and identifies the species based on those measurements alone.



We've been given a [data set \(iris-data.csv\)](#) from our field researchers to develop the demo, which only includes measurements for three types of *Iris* flowers:

Iris setosa



Iris versicolor



Step 1: Answering the question

[\[go back to the top \]](#)

The first step to any data analysis project is to define the question or problem we're looking to solve, and to define a measure (or set of measures) for our success at solving that task. The data analysis checklist has us answer a handful of questions to accomplish that, so let's work through those questions.

Did you specify the type of data analytic question (e.g. exploration, association causality) before touching the data?

We're trying to classify the species (i.e., class) of the flower based on four measurements that we're provided: sepal length, sepal width, petal length, and petal width.

Did you define the metric for success before beginning?

Let's do that now. Since we're performing classification, we can use [accuracy \(https://en.wikipedia.org/wiki/Accuracy_and_precision\)](https://en.wikipedia.org/wiki/Accuracy_and_precision) — the fraction of correctly classified flowers — to quantify how well our model is performing. Our head of data has told us that we should achieve at least 90% accuracy.

Did you understand the context for the question and the scientific or business application?

We're building part of a data analysis pipeline for a smartphone app that will be able to classify the species of flowers from pictures taken on the smartphone. In the future, this pipeline will be connected to another pipeline that automatically measures from pictures the traits we're using to perform this classification.

Did you record the experimental design?

Our head of data has told us that the field researchers are hand-measuring 50 randomly-sampled flowers of each species using a standardized methodology. The field researchers take pictures of each flower they sample from pre-defined angles so the measurements and species can be confirmed by the other field researchers at a later point. At the end of each day, the data is compiled and stored on a private company GitHub repository.

Did you consider whether the question could be answered with the available data?

The data set we currently have is only for three types of *Iris* flowers. The model built off of this data set will only work for those *Iris* flowers, so we will need more data to create a general flower classifier.

Notice that we've spent a fair amount of time working on the problem without writing a line of code or even looking at the data.

Thinking about and documenting the problem we're working on is an important step to performing effective data analysis that often goes overlooked. Don't skip it.

Step 2: Checking the data

[\[go back to the top \]](#)

The next step is to look at the data we're working with. Even curated data sets from the government can have errors in them, and it's vital that we spot these errors before investing too much time in our analysis.

Generally, we're looking to answer the following questions:

- Is there anything wrong with the data?
- Are there any quirks with the data?
- Do I need to fix or remove any of the data?

Let's start by reading the data into a pandas DataFrame.

```
In [1]: import pandas as pd

iris_data = pd.read_csv('iris-data.csv')
iris_data.head()
```

```
Out[1]:
```

	sepal_length_cm	sepal_width_cm	petal_length_cm	petal_width_cm	class
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

We're in luck! The data seems to be in a usable format.

The first row in the data file defines the column headers, and the headers are descriptive enough for us to understand what each column represents. The headers even give us the units that the measurements were recorded in, just in case we needed to know at a later point in the project.

Each row following the first row represents an entry for a flower: four measurements and one class, which tells us the species of the flower.

One of the first things we should look for is missing data. Thankfully, the field researchers already told us that they put a 'NA' into the spreadsheet when they were missing a measurement.

We can tell pandas to automatically identify missing values if it knows our missing value marker.

```
In [2]: iris_data = pd.read_csv('iris-data.csv', na_values=['NA'])
```

Voilà! Now pandas knows to treat rows with 'NA' as missing values.

Next, it's always a good idea to look at the distribution of our data — especially the outliers.

Let's start by printing out some summary statistics about the data set.

In [3]: `iris_data.describe()`

```
/home/dima/anaconda2/lib/python2.7/site-packages/numpy/lib/function_base.py:3834: RuntimeWarning: Invalid value encountered in percentile
RuntimeWarning)
```

Out[3]:

	sepal_length_cm	sepal_width_cm	petal_length_cm	petal_width_cm
count	150.000000	150.000000	150.000000	145.000000
mean	5.644627	3.054667	3.758667	1.236552
std	1.312781	0.433123	1.764420	0.755058
min	0.055000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	NaN
50%	5.700000	3.000000	4.350000	NaN
75%	6.400000	3.300000	5.100000	NaN
max	7.900000	4.400000	6.900000	2.500000

We can see several useful values from this table. For example, we see that five `petal_width_cm` entries are missing.

If you ask me, though, tables like this are rarely useful unless we know that our data should fall in a particular range. It's usually better to visualize the data in some way. Visualization makes outliers and errors immediately stand out, whereas they might go unnoticed in a large table of numbers.

Since we know we're going to be plotting in this section, let's set up the notebook so we can plot inside of it.

In [4]: `# This line tells the notebook to show plots inside of the notebook`
`%matplotlib inline`

```
import matplotlib.pyplot as plt
import seaborn as sb
```

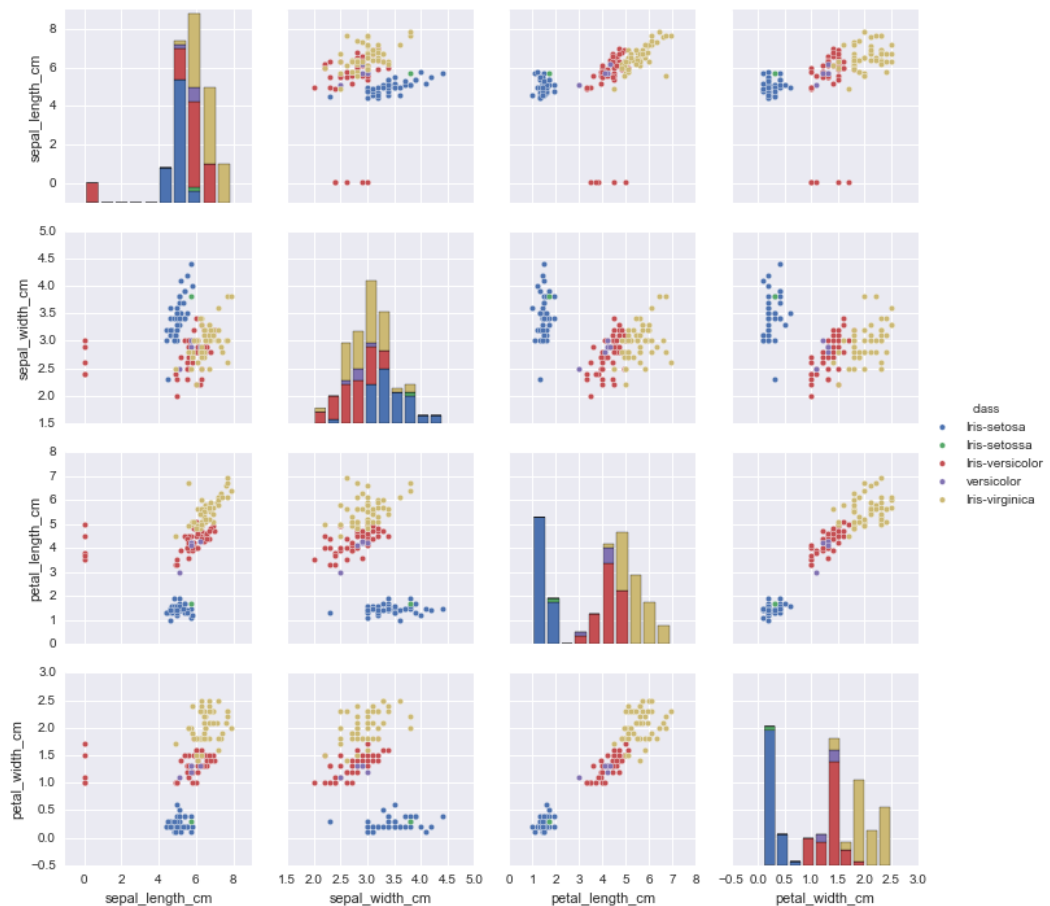
```
/home/dima/anaconda2/lib/python2.7/site-packages/matplotlib/font_manager.py:273: UserWarning: Matplotlib is building the font cache using fc-list.
This may take a moment.
warnings.warn('Matplotlib is building the font cache using fc-list. This
s may take a moment.')
```

Next, let's create a **scatterplot matrix**. Scatterplot matrices plot the distribution of each column along the diagonal, and then plot a scatterplot matrix for the combination of each variable. They make for an efficient tool to look for errors in our data.

We can even have the plotting package color each entry by its class to look for trends within the classes.


```
In [5]: # We have to temporarily drop the rows with 'NA' values
# because the Seaborn plotting function does not know
# what to do with them
sb.pairplot(iris_data.dropna(), hue='class')
```

```
Out[5]: <seaborn.axisgrid.PairGrid at 0x7fb3a4dccad0>
```



From the scatterplot matrix, we can already see some issues with the data set:

1. There are five classes when there should only be three, meaning there were some coding errors.
2. There are some clear outliers in the measurements that may be erroneous: one `sepal_width_cm` entry for `Iris-setosa` falls well outside its normal range, and several `sepal_length_cm` entries for `Iris-versicolor` are near-zero for some reason.
3. We had to drop those rows with missing values.

In all of these cases, we need to figure out what to do with the erroneous data. Which takes us to the next step...

Step 3: Tidying the data

[\[go back to the top \]](#)

Now that we've identified several errors in the data set, we need to fix them before we proceed with the analysis.

Let's walk through the issues one-by-one.

There are five classes when there should only be three, meaning there were some coding errors.

After talking with the field researchers, it sounds like one of them forgot to add *Iris* - before their *Iris-versicolor* entries. The other extraneous class, *Iris-setossa*, was simply a typo that they forgot to fix.

Let's use the `DataFrame` to fix these errors.

```
In [55]: iris_data.loc[iris_data['class'] ==  
                'versicolor', 'class'] = 'Iris-versicolor'  
iris_data.loc[iris_data['class'] ==  
                'Iris-setossa', 'class'] = 'Iris-setosa'  
  
iris_data['class'].unique()  
  
Out[55]: array(['Iris-setosa', 'Iris-versicolor', 'Iris-virginica'], dtype=object)
```

Much better! Now we only have three class types. Imagine how embarrassing it would've been to create a model that used the wrong classes.

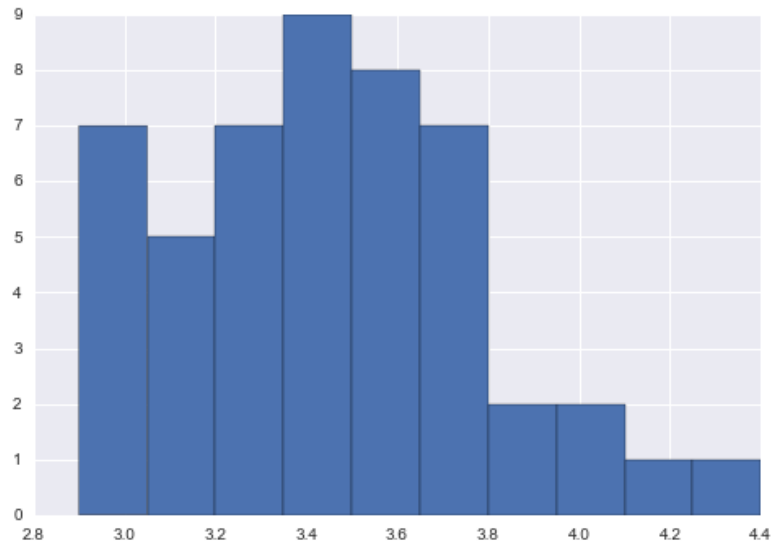
There are some clear outliers in the measurements that may be erroneous: one `sepal_width_cm` entry for *Iris-setosa* falls well outside its normal range, and several `sepal_length_cm` entries for *Iris-versicolor* are near-zero for some reason.

Fixing outliers can be tricky business. It's rarely clear whether the outlier was caused by measurement error, recording the data in improper units, or if the outlier is a real anomaly. For that reason, we should be judicious when working with outliers: if we decide to exclude any data, we need to make sure to document what data we excluded and provide solid reasoning for excluding that data. (i.e., "This data didn't fit my hypothesis" will not stand peer review.)

In the case of the one anomalous entry for *Iris-setosa*, let's say our field researchers know that it's impossible for *Iris-setosa* to have a sepal width below 2.5 cm. Clearly this entry was made in error, and we're better off just scrapping the entry than spending hours finding out what happened.

```
In [7]: # This line drops any 'Iris-setosa'
#rows with a sepal width less than 2.5 cm
iris_data = iris_data.loc[(iris_data['class'] != 'Iris-setosa') |
                          (iris_data['sepal_width_cm'] >= 2.5)]
iris_data.loc[iris_data['class'] == 'Iris-setosa',
              'sepal_width_cm'].hist()
```

```
Out[7]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb396cf9810>
```



Excellent! Now all of our *Iris-setosa* rows have a sepal width greater than 2.5.

The next data issue to address is the several near-zero sepal lengths for the *Iris-versicolor* rows. Let's take a look at those rows.

```
In [8]: iris_data.loc[(iris_data['class'] == 'Iris-versicolor') &
                      (iris_data['sepal_length_cm'] < 1.0)]
```

```
Out[8]:
```

	sepal_length_cm	sepal_width_cm	petal_length_cm	petal_width_cm	class
77	0.067	3.0	5.0	1.7	Iris-versicolor
78	0.060	2.9	4.5	1.5	Iris-versicolor
79	0.057	2.6	3.5	1.0	Iris-versicolor
80	0.055	2.4	3.8	1.1	Iris-versicolor
81	0.055	2.4	3.7	1.0	Iris-versicolor

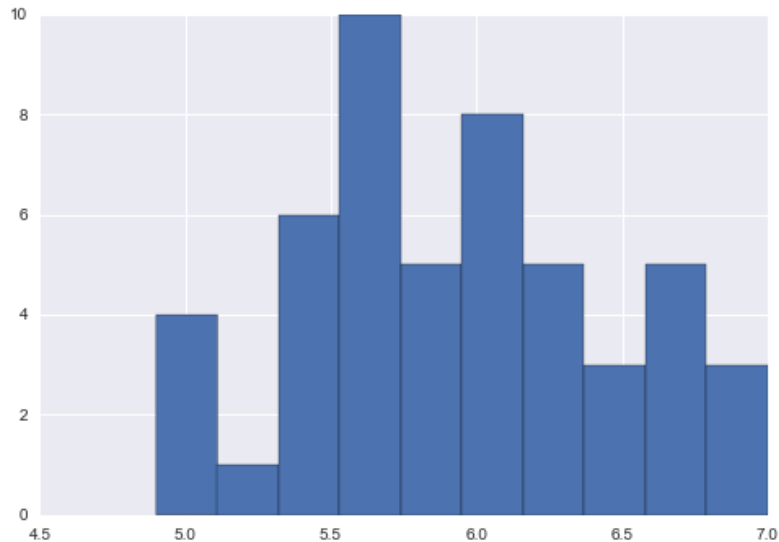
How about that? All of these near-zero `sepal_length_cm` entries seem to be off by two orders of magnitude, as if they had been recorded in meters instead of centimeters.

After some brief correspondence with the field researchers, we find that one of them forgot to convert those measurements to centimeters. Let's do that for them.

```
In [9]: iris_data.loc[(iris_data['class'] == 'Iris-versicolor') &
                    (iris_data['sepal_length_cm'] < 1.0),
                    'sepal_length_cm'] *= 100.0

iris_data.loc[iris_data['class'] == 'Iris-versicolor',
              'sepal_length_cm'].hist()
```

```
Out[9]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb3974acd50>
```



Phew! Good thing we fixed those outliers. They could've really thrown our analysis off.

We had to drop those rows with missing values.

Let's take a look at the rows with missing values:

```
In [10]: iris_data.loc[(iris_data['sepal_length_cm'].isnull() |
                        (iris_data['sepal_width_cm'].isnull() |
                         (iris_data['petal_length_cm'].isnull() |
                          (iris_data['petal_width_cm'].isnull()))))]
```

```
Out[10]:
```

	sepal_length_cm	sepal_width_cm	petal_length_cm	petal_width_cm	class
7	5.0	3.4	1.5	NaN	Iris-setosa
8	4.4	2.9	1.4	NaN	Iris-setosa
9	4.9	3.1	1.5	NaN	Iris-setosa
10	5.4	3.7	1.5	NaN	Iris-setosa
11	4.8	3.4	1.6	NaN	Iris-setosa

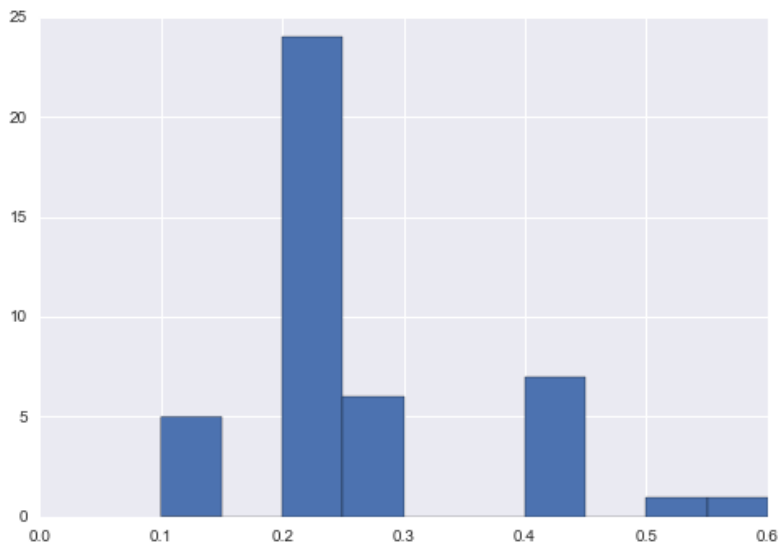
It's not ideal that we had to drop those rows, especially considering they're all *Iris-setosa* entries. Since it seems like the missing data is systematic — all of the missing values are in the same column for the same *Iris* type — this error could potentially bias our analysis.

One way to deal with missing data is **mean imputation**: If we know that the values for a measurement fall in a certain range, we can fill in empty values with the average of that measurement.

Let's see if we can do that here.

```
In [11]: iris_data.loc[iris_data['class'] == 'Iris-setosa',
                    'petal_width_cm'].hist()
```

```
Out[11]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb3974ac190>
```



Most of the petal widths for Iris-setosa fall within the 0.2-0.3 range, so let's fill in these entries with the average measured petal width.

```
In [12]: average_petal_width = iris_data.loc[iris_data['class'] ==
                    'Iris-setosa',
                    'petal_width_cm'].mean()

iris_data.loc[(iris_data['class'] == 'Iris-setosa') &
              (iris_data['petal_width_cm'].isnull()),
              'petal_width_cm'] = average_petal_width

iris_data.loc[(iris_data['class'] == 'Iris-setosa') &
              (iris_data['petal_width_cm'] == average_petal_width)]
```

```
Out[12]:
```

	sepal_length_cm	sepal_width_cm	petal_length_cm	petal_width_cm	class
7	5.0	3.4	1.5	0.25	Iris-setosa
8	4.4	2.9	1.4	0.25	Iris-setosa
9	4.9	3.1	1.5	0.25	Iris-setosa
10	5.4	3.7	1.5	0.25	Iris-setosa
11	4.8	3.4	1.6	0.25	Iris-setosa

```
In [13]: iris_data.loc[(iris_data['sepal_length_cm'].isnull() |
                    (iris_data['sepal_width_cm'].isnull() |
                     (iris_data['petal_length_cm'].isnull() |
                      (iris_data['petal_width_cm'].isnull()))))]
```

```
Out[13]:
```

	sepal_length_cm	sepal_width_cm	petal_length_cm	petal_width_cm	class
--	-----------------	----------------	-----------------	----------------	-------

Great! Now we've recovered those rows and no longer have missing data in our data set.

Note: If you don't feel comfortable imputing your data, you can drop all rows with missing data with the `dropna()` call:

```
iris_data.dropna(inplace=True)
```

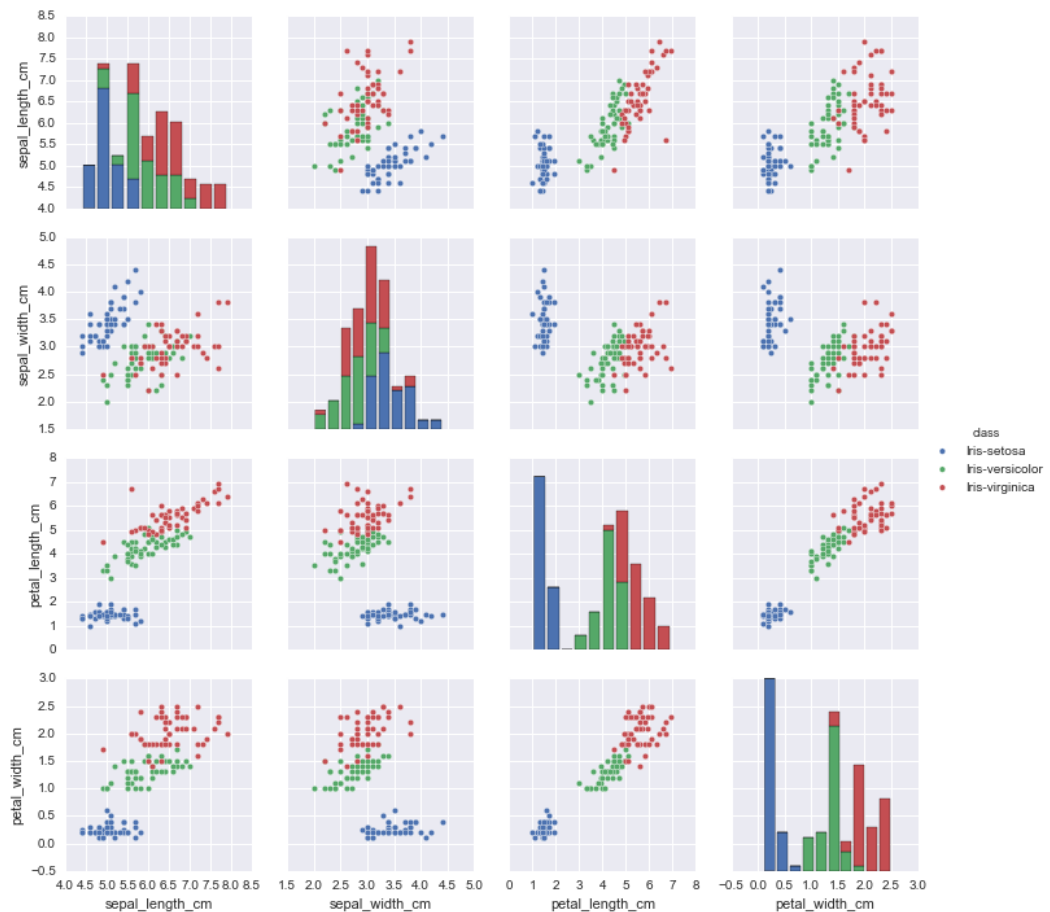
After all this hard work, we don't want to repeat this process every time we work with the data set. Let's save the tidied data file as a *separate file* and work directly with that data file from now on.

```
In [14]: iris_data.to_csv('iris-data-clean.csv', index=False)
iris_data_clean = pd.read_csv('iris-data-clean.csv')
```

Now, let's take a look at the scatterplot matrix now that we've tidied the data.

```
In [15]: sb.pairplot(iris_data_clean, hue='class')
```

```
Out[15]: <seaborn.axisgrid.PairGrid at 0x7fb396b06ad0>
```



Of course, I purposely inserted numerous errors into this data set to demonstrate some of the many possible scenarios you may face while tidying your data.

The general takeaways here should be:

- Make sure your data is encoded properly
- Make sure your data falls within the expected range, and use domain knowledge whenever possible to define that expected range
- Deal with missing data in one way or another: replace it if you can or drop it
- Never tidy your data manually because that is not easily reproducible
- Use code as a record of how you tidied your data
- Plot everything you can about the data at this stage of the analysis so you can *visually* confirm everything looks correct

Bonus: Testing our data

[\[go back to the top \]](#)

At SciPy 2015, I was exposed to a great idea: We should test our data. Just how we use unit tests to verify our expectations from code, we can similarly set up unit tests to verify our expectations about a data set.

We can quickly test our data using assert statements: We assert that something must be true, and if it is, then nothing happens and the notebook continues running. However, if our assertion is wrong, then the notebook stops running and brings it to our attention. For example:

```
In [16]: assert 1 == 2
```

```
AssertionErrorTraceback (most recent call last)
<ipython-input-16-a810b3a4aded> in <module>()
----> 1 assert 1 == 2

AssertionError:
```

Let's test a few things that we know about our data set now.

```
In [17]: # We know that we should only have three classes
assert len(iris_data_clean['class'].unique()) == 3
```

```
In [18]: # We know that sepal lengths for 'Iris-versicolor'
#should never be below 2.5 cm
assert iris_data_clean.loc[iris_data_clean['class'] ==
                          'Iris-versicolor', 'sepal_length_cm'].min() >
= 2.5
```

```
In [19]: # We know that our data set should have no missing measurements
assert len(iris_data_clean.loc[(iris_data_clean['sepal_length_cm'].isnull() |
                                (iris_data_clean['sepal_width_cm'].isnull() |
                                (iris_data_clean['petal_length_cm'].isnull() |
                                (iris_data_clean['petal_width_cm'].isnull()
                                ()))]) == 0
```

And so on. If any of these expectations are violated, then our analysis immediately stops and we have to return to the tidying stage.

Step 4: Exploratory analysis

[\[go back to the top\]](#)

Now after spending entirely too much time tidying our data, we can start analyzing it!

Exploratory analysis is the step where we start delving deeper into the data set beyond the outliers and errors. We'll be looking to answer questions such as:

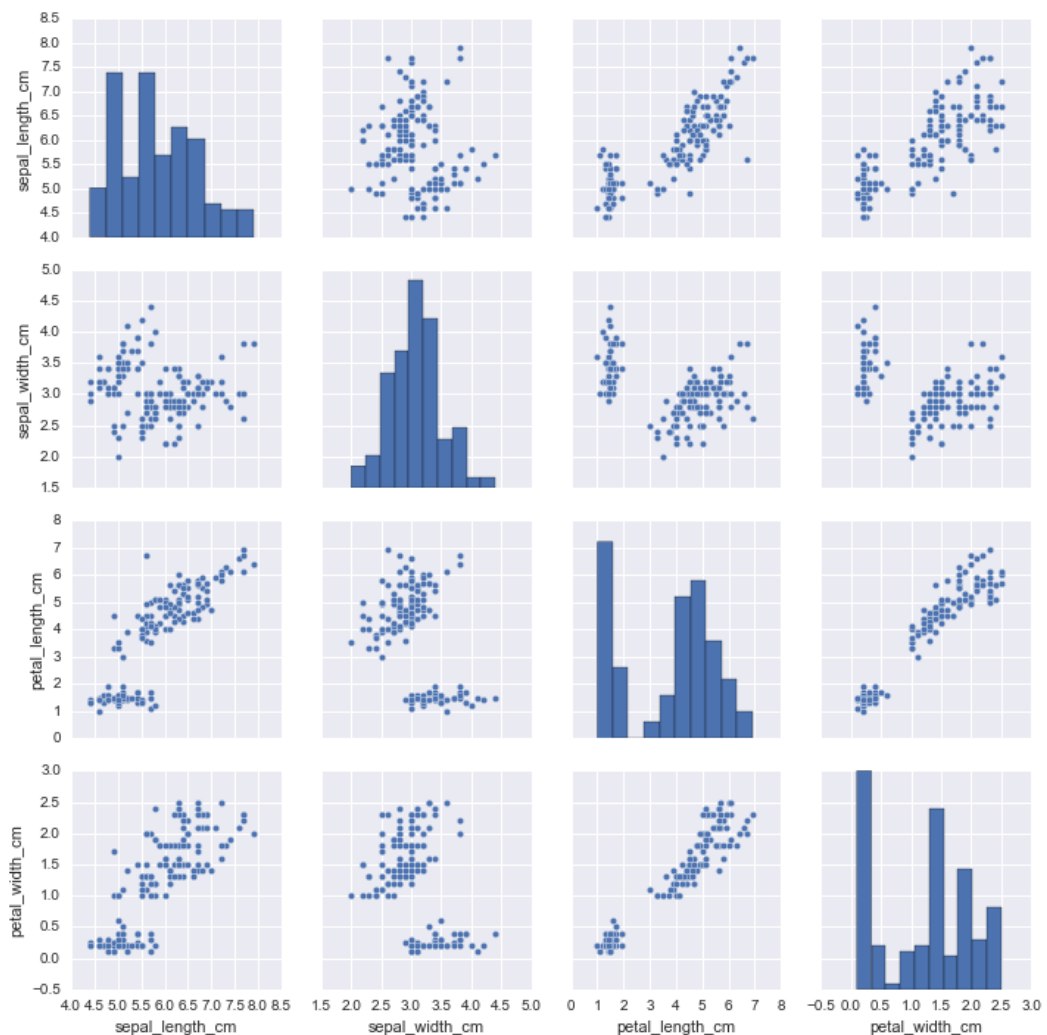
- How is my data distributed?
- Are there any correlations in my data?
- Are there any confounding factors that explain these correlations?

This is the stage where we plot all the data in as many ways as possible. Create many charts, but don't bother making them pretty — these charts are for internal use.

Let's return to that scatterplot matrix that we used earlier.

```
In [20]: sb.pairplot(iris_data_clean)
```

```
Out[20]: <seaborn.axisgrid.PairGrid at 0x7fb38fe11250>
```

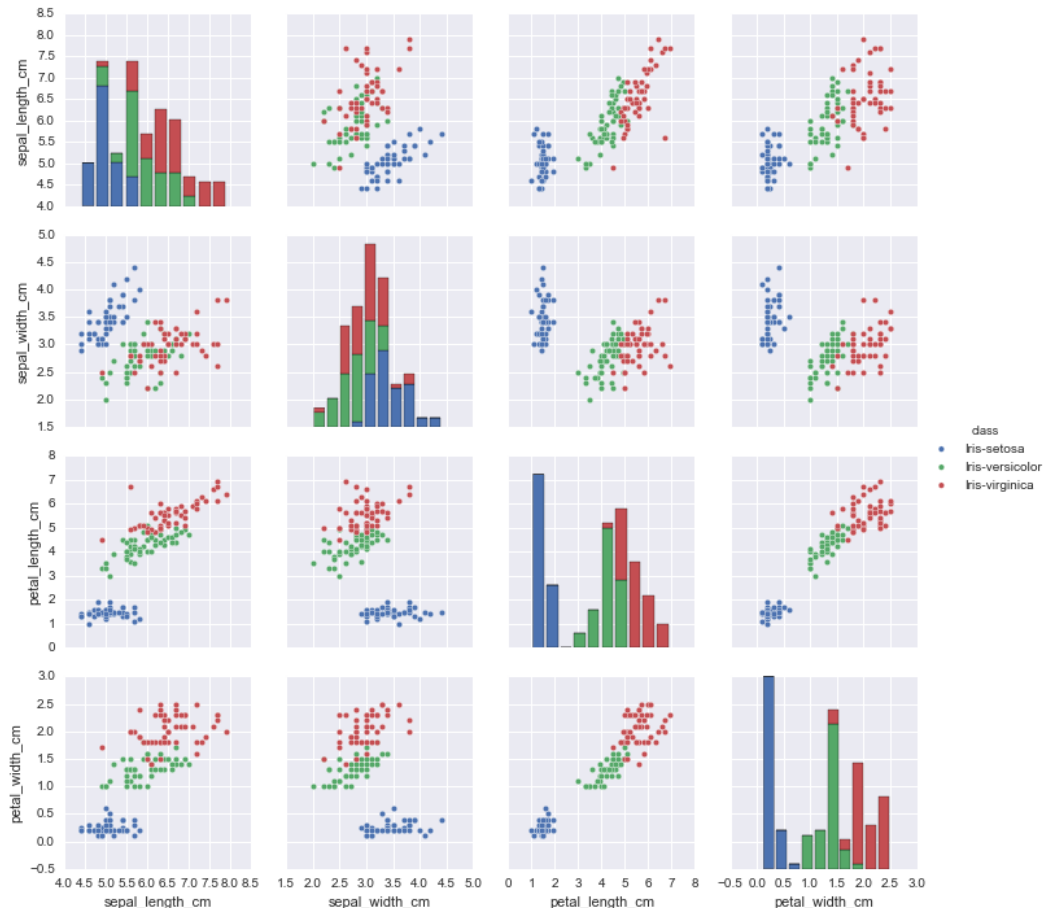


Our data is normally distributed for the most part, which is great news if we plan on using any modeling methods that assume the data is normally distributed.

There's something strange going with the petal measurements. Maybe it's something to do with the different Iris types. Let's color code the data by the class again to see if that clears things up.

```
In [21]: sb.pairplot(iris_data_clean, hue='class')
```

```
Out[21]: <seaborn.axisgrid.PairGrid at 0x7fb38ef94550>
```



Sure enough, the strange distribution of the petal measurements exist because of the different species. This is actually great news for our classification task since it means that the petal measurements will make it easy to distinguish between Iris-setosa and the other Iris types.

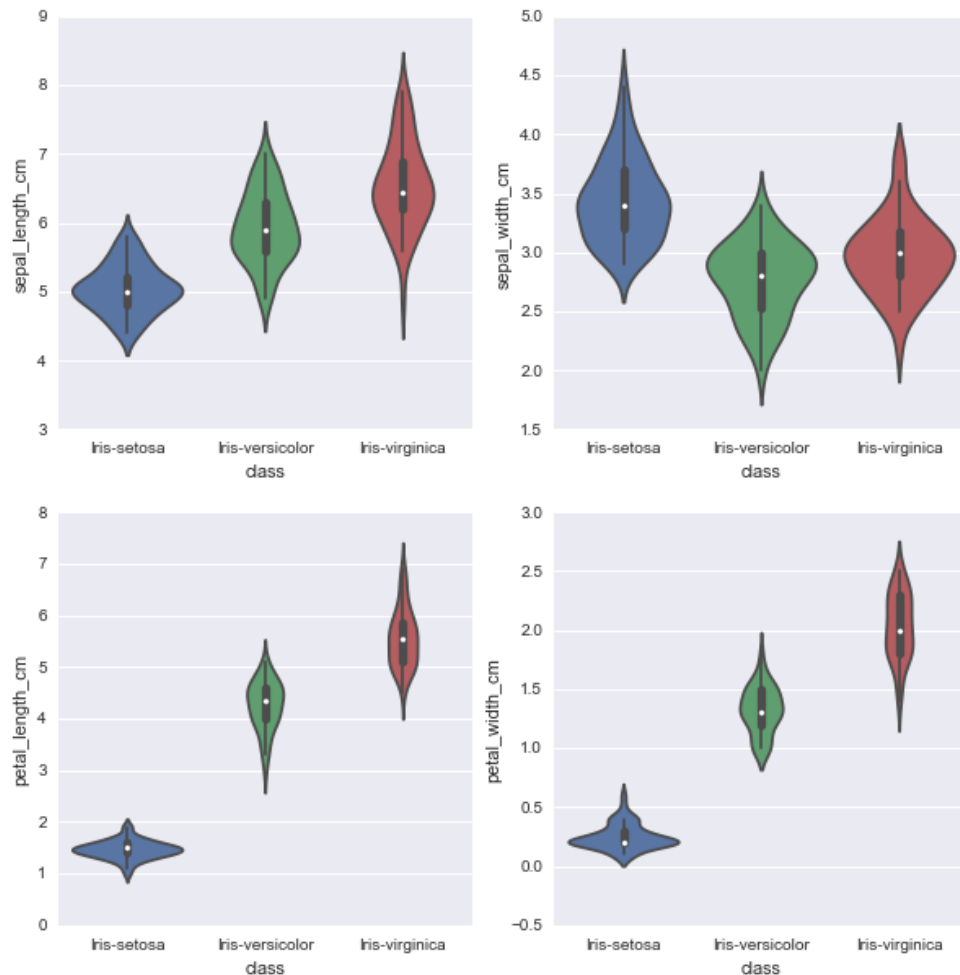
Distinguishing Iris-versicolor and Iris-virginica will prove more difficult given how much their measurements overlap.

There are also correlations between petal length and petal width, as well as sepal length and sepal width. The field biologists assure us that this is to be expected: Longer flower petals also tend to be wider, and the same applies for sepals.

We can also make **violin plots** of the data to compare the measurement distributions of the classes. Violin plots contain the same information as [box plots](https://en.wikipedia.org/wiki/Box_plot) (https://en.wikipedia.org/wiki/Box_plot), but also scales the box according to the density of the data.

```
In [22]: plt.figure(figsize=(10, 10))

for column_index, column in enumerate(iris_data_clean.columns):
    if column == 'class':
        continue
    plt.subplot(2, 2, column_index + 1)
    sb.violinplot(x='class', y=column, data=iris_data_clean)
```



Enough flirting with the data. Let's get to modeling.

Step 5: Classification

[\[go back to the top \]](#)

Wow, all this work and we *still* haven't modeled the data!

As tiresome as it can be, tidying and exploring our data is a vital component to any data analysis. If we had jumped straight to the modeling step, we would have created a faulty classification model.

Remember: **Bad data leads to bad models**. Always check your data first.

Assured that our data is now as clean as we can make it — and armed with some cursory knowledge of the distributions and relationships in our data set — it's time to make the next big step in our analysis: Splitting the data into training and testing sets.

A **training set** is a random subset of the data that we use to train our models.

A **testing set** is a random subset of the data (mutually exclusive from the training set) that we use to validate our models on unforeseen data.

Especially in sparse data sets like ours, it's easy for models to **overfit** the data: The model will learn the training set so well that it won't be able to handle most of the cases it's never seen before. This is why it's important for us to build the model with the training set, but score it with the testing set.

Note that once we split the data into a training and testing set, we should treat the testing set like it no longer exists: We cannot use any information from the testing set to build our model or else we're cheating.

Let's set up our data first.

```
In [23]: iris_data_clean = pd.read_csv('iris-data-clean.csv')

# We're using all four measurements as inputs
# Note that scikit-learn expects each entry to be a list of values, e.g.
# [ [val1, val2, val3],
#   [val1, val2, val3],
#   ... ]
# such that our input data set is represented as a list of lists

# We can extract the data in this format from pandas like this:
all_inputs = iris_data_clean[['sepal_length_cm', 'sepal_width_cm',
                              'petal_length_cm', 'petal_width_cm']].values

# Similarly, we can extract the classes
all_classes = iris_data_clean['class'].values

# Make sure that you don't mix up the order of the entries
# all_inputs[5] inputs should correspond to the class in all_classes[5]

# Here's what a subset of our inputs looks like:
all_inputs[:5]
```

```
Out[23]: array([[ 5.1,  3.5,  1.4,  0.2],
                [ 4.9,  3. ,  1.4,  0.2],
                [ 4.7,  3.2,  1.3,  0.2],
                [ 4.6,  3.1,  1.5,  0.2],
                [ 5. ,  3.6,  1.4,  0.2]])
```

Now our data is ready to be split.

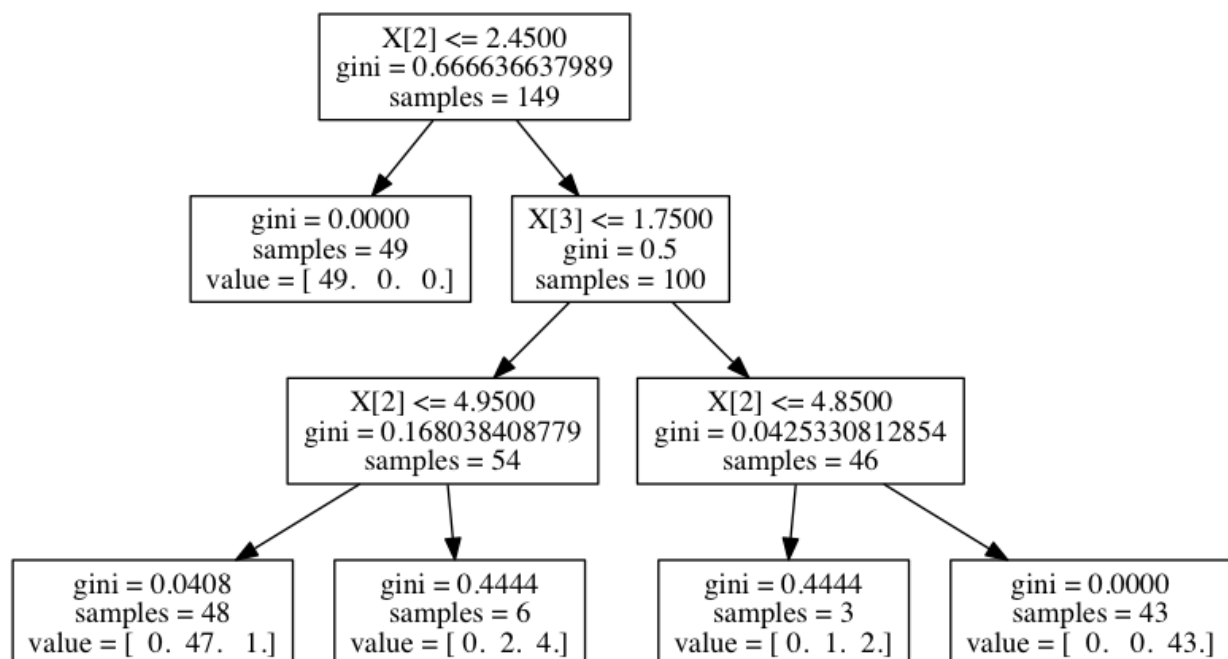
```
In [24]: from sklearn.cross_validation import train_test_split

(training_inputs,
 testing_inputs,
 training_classes,
 testing_classes) = train_test_split(all_inputs, all_classes,
                                     train_size=0.75, random_state=1)
```

With our data split, we can start fitting models to our data. Our head of data is all about decision tree classifiers, so let's start with one of those.

Decision tree classifiers are incredibly simple in theory. In their simplest form, decision tree classifiers ask a series of Yes/No questions about the data — each time getting closer to finding out the class of each entry — until they either classify the data set perfectly or simply can't differentiate a set of entries. Think of it like a game of Twenty Questions (https://en.wikipedia.org/wiki/Twenty_Questions), except the computer is *much, much* better at it.

Here's an example decision tree classifier:



Notice how the classifier asks Yes/No questions about the data — whether a certain feature is ≤ 1.75 , for example — so it can differentiate the records. This is the essence of every decision tree.

The nice part about decision tree classifiers is that they are **scale-invariant**, i.e., the scale of the features does not affect their performance, unlike many Machine Learning models. In other words, it doesn't matter if our features range from 0 to 1 or 0 to 1,000; decision tree classifiers will work with them just the same.

There are several parameters (<http://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>) that we can tune for decision tree classifiers, but for now let's use a basic decision tree classifier.

```
In [25]: from sklearn.tree import DecisionTreeClassifier

# Create the classifier
decision_tree_classifier = DecisionTreeClassifier()

# Train the classifier on the training set
decision_tree_classifier.fit(training_inputs,
                             training_classes)

# Validate the classifier on the testing set using classification accuracy
decision_tree_classifier.score(testing_inputs,
                               testing_classes)
```

Out[25]: 0.97368421052631582

Heck yeah! Our model achieves 97% classification accuracy without much effort.

However, there's a catch: Depending on how our training and testing set was sampled, our model can achieve anywhere from 80% to 100% accuracy:

```
In [26]: model_accuracies = []

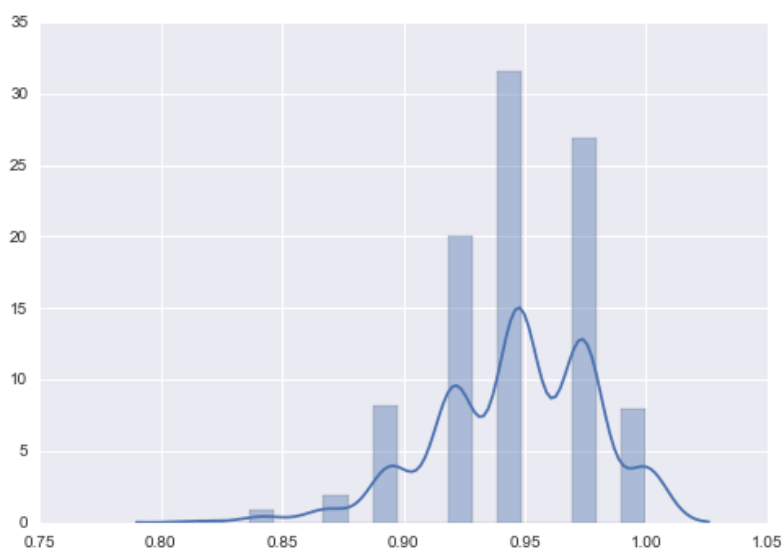
for repetition in range(1000):
    (training_inputs,
     testing_inputs,
     training_classes,
     testing_classes) = train_test_split(all_inputs,
                                          all_classes, train_size=0.75)

    decision_tree_classifier = DecisionTreeClassifier()
    decision_tree_classifier.fit(training_inputs, training_classes)
    classifier_accuracy = decision_tree_classifier.score(testing_inputs,
                                                         testing_classes)

    model_accuracies.append(classifier_accuracy)

sb.distplot(model_accuracies)
```

Out[26]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb38c643110>



It's obviously a problem that our model performs quite differently depending on the subset of the data it's trained on. This phenomenon is known as **overfitting**: The model is learning to classify the training set so well that it doesn't generalize and perform well on data it hasn't seen before.

Cross-validation

[\[go back to the top \]](#)

This problem is the main reason that most data scientists perform **k-fold cross-validation** on their models: Split the original data set into k subsets, use one of the subsets as the testing set, and the rest of the subsets are used as the training set. This process is then repeated k times such that each subset is used as the testing set exactly once.

10-fold cross-validation is the most common choice, so let's use that here. Performing 10-fold cross-validation on our data set looks something like this:

(each square is an entry in our data set)

```
In [27]: import numpy as np
from sklearn.cross_validation import StratifiedKFold

def plot_cv(cv, n_samples):
    masks = []
    for train, test in cv:
        mask = np.zeros(n_samples, dtype=bool)
        mask[test] = 1
        masks.append(mask)

    plt.figure(figsize=(15, 15))
    plt.imshow(masks, interpolation='none')
    plt.ylabel('Fold')
    plt.xlabel('Row #')

plot_cv(StratifiedKFold(all_classes, n_folds=10),
        len(all_classes))
```



You'll notice that we used **Stratified k-fold cross-validation** in the code above. Stratified k -fold keeps the class proportions the same across all of the folds, which is vital for maintaining a representative subset of our data set. (e.g., so we don't have 100% Iris setosa entries in one of the folds.)

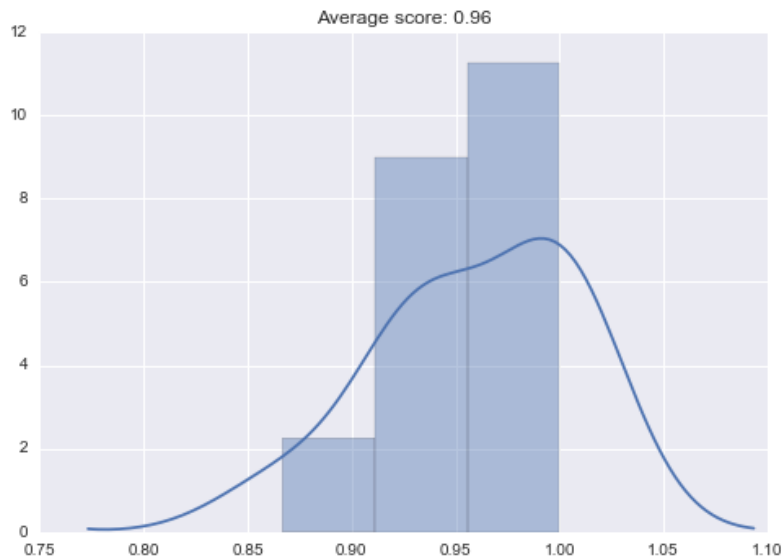
We can perform 10-fold cross-validation on our model with the following code:

```
In [28]: from sklearn.cross_validation import cross_val_score

decision_tree_classifier = DecisionTreeClassifier()

# cross_val_score returns a list of the scores, which we can visualize
# to get a reasonable estimate of our classifier's performance
cv_scores = cross_val_score(decision_tree_classifier,
                             all_inputs, all_classes, cv=10)
sb.distplot(cv_scores)
plt.title('Average score: {}'.format(np.mean(cv_scores)))
```

Out[28]: <matplotlib.text.Text at 0x7fb387d86f50>



Now we have a much more consistent rating of our classifier's general classification accuracy.

Parameter tuning

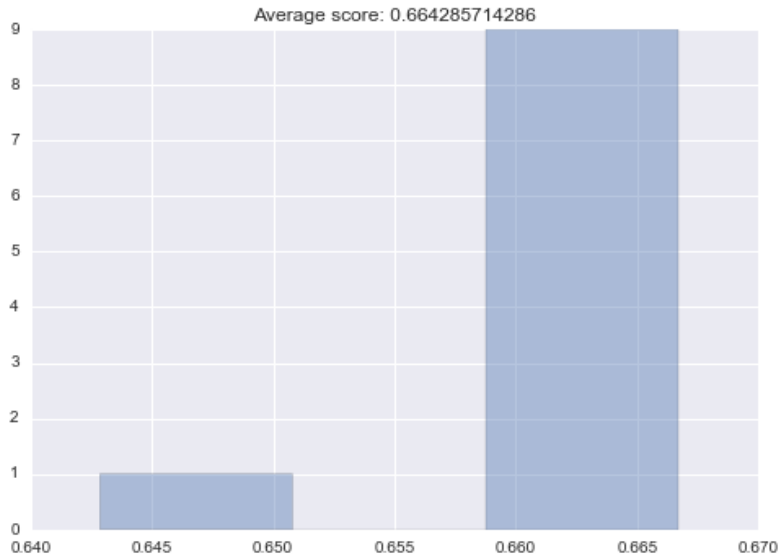
[\[go back to the top\]](#)

Every Machine Learning model comes with a variety of parameters to tune, and these parameters can be vitally important to the performance of our classifier. For example, if we severely limit the depth of our decision tree classifier:

```
In [29]: decision_tree_classifier = DecisionTreeClassifier(max_depth=1)

cv_scores = cross_val_score(decision_tree_classifier,
                             all_inputs, all_classes, cv=10)
sb.distplot(cv_scores, kde=False)
plt.title('Average score: {}'.format(np.mean(cv_scores)))
```

```
Out[29]: <matplotlib.text.Text at 0x7fb387c61850>
```



the classification accuracy falls tremendously.

Therefore, we need to find a systematic method to discover the best parameters for our model and data set.

The most common method for model parameter tuning is **Grid Search**. The idea behind Grid Search is simple: explore a range of parameters and find the best-performing parameter combination. Focus your search on the best range of parameters, then repeat this process several times until the best parameters are discovered.

Let's tune our decision tree classifier. We'll stick to only two parameters for now, but it's possible to simultaneously explore dozens of parameters if we want.

```
In [30]: from sklearn.grid_search import GridSearchCV

decision_tree_classifier = DecisionTreeClassifier()

parameter_grid = {'max_depth': [1, 2, 3, 4, 5],
                  'max_features': [1, 2, 3, 4]}

cross_validation = StratifiedKFold(all_classes, n_folds=10)

grid_search = GridSearchCV(decision_tree_classifier,
                           param_grid=parameter_grid,
                           cv=cross_validation)

grid_search.fit(all_inputs, all_classes)
print('Best score: {}'.format(grid_search.best_score_))
print('Best parameters: {}'.format(grid_search.best_params_))
```

```
Best score: 0.959731543624
```

```
Best parameters: {'max_features': 3, 'max_depth': 3}
```

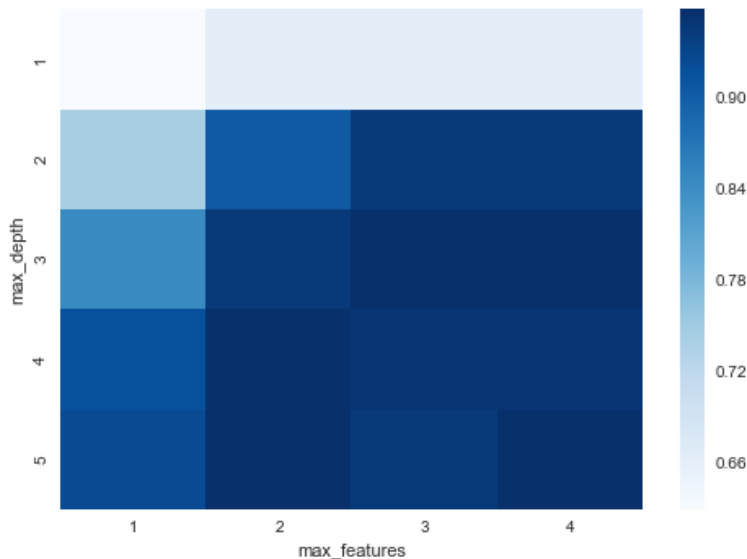
Now let's visualize the grid search to see how the parameters interact.


```
In [31]: grid_visualization = []

        for grid_pair in grid_search.grid_scores_:
            grid_visualization.append(grid_pair.mean_validation_score)

        grid_visualization = np.array(grid_visualization)
        grid_visualization.shape = (5, 4)
        sb.heatmap(grid_visualization, cmap='Blues')
        plt.xticks(np.arange(4) + 0.5, grid_search.param_grid['max_features'])
        plt.yticks(np.arange(5) + 0.5, grid_search.param_grid['max_depth'][::-1])
        plt.xlabel('max_features')
        plt.ylabel('max_depth')
```

Out[31]: <matplotlib.text.Text at 0x7fb387bb7c50>



Now we have a better sense of the parameter space: We know that we need a max_depth of at least 2 to allow the decision tree to make more than a one-off decision.

max_features doesn't really seem to make a big difference here as long as we have 2 of them, which makes sense since our data set has only 4 features and is relatively easy to classify. (Remember, one of our data set's classes was easily separable from the rest based on a single feature.)

Let's go ahead and use a broad grid search to find the best settings for a handful of parameters.

```
In [32]: decision_tree_classifier = DecisionTreeClassifier()

parameter_grid = {'criterion': ['gini', 'entropy'],
                  'splitter': ['best', 'random'],
                  'max_depth': [1, 2, 3, 4, 5],
                  'max_features': [1, 2, 3, 4]}

cross_validation = StratifiedKFold(all_classes, n_folds=10)

grid_search = GridSearchCV(decision_tree_classifier,
                           param_grid=parameter_grid,
                           cv=cross_validation)

grid_search.fit(all_inputs, all_classes)
print('Best score: {}'.format(grid_search.best_score_))
print('Best parameters: {}'.format(grid_search.best_params_))

Best score: 0.96644295302
Best parameters: {'max_features': 3, 'splitter': 'best', 'criterion': 'gini', 'max_depth': 3}
```

Now we can take the best classifier from the Grid Search and use that:

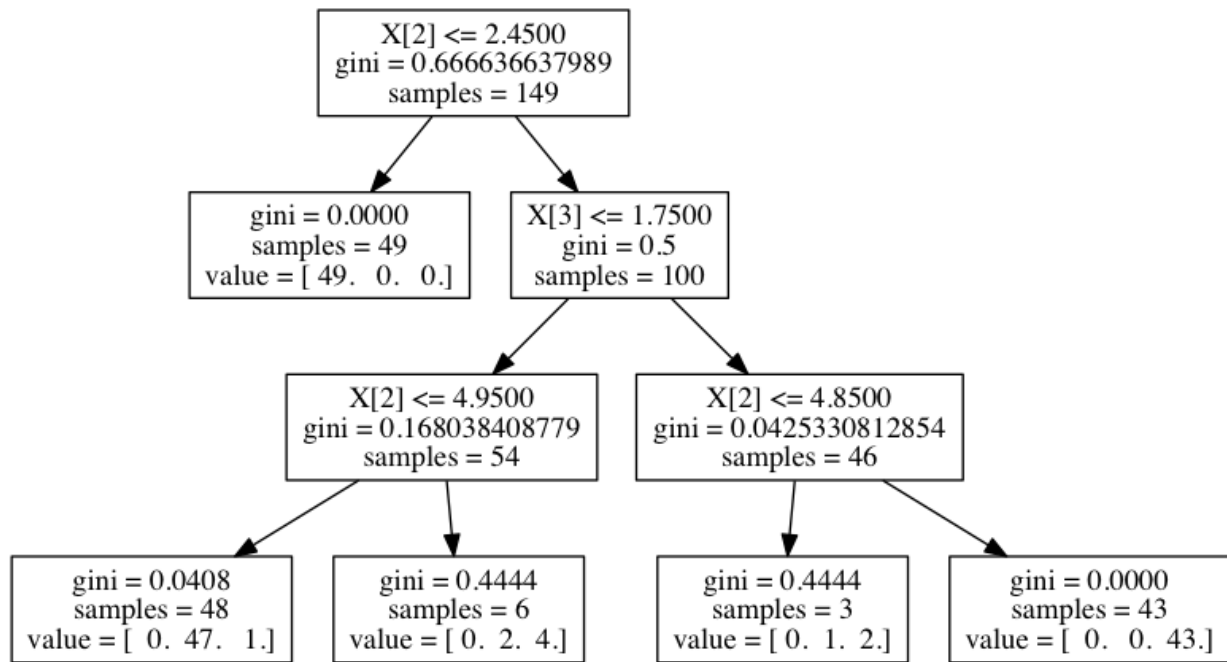
```
In [33]: decision_tree_classifier = grid_search.best_estimator_
decision_tree_classifier

Out[33]: DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=3,
                               max_features=3, max_leaf_nodes=None, min_samples_leaf=1,
                               min_samples_split=2, min_weight_fraction_leaf=0.0,
                               presort=False, random_state=None, splitter='best')
```

We can even visualize the decision tree with [GraphViz \(http://www.graphviz.org/\)](http://www.graphviz.org/) to see how it's making the classifications:

```
In [34]: import sklearn.tree as tree
from sklearn.externals.six import StringIO

with open('iris_dtc.dot', 'w') as out_file:
    out_file = tree.export_graphviz(decision_tree_classifier,
                                    out_file=out_file)
```



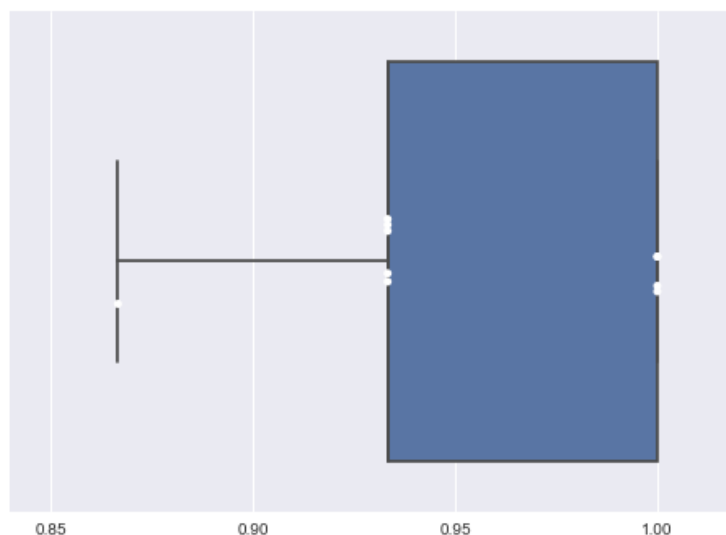
(This classifier may look familiar from earlier in the notebook.)

Alright! We finally have our demo classifier. Let's create some visuals of its performance so we have something to show our head of data.

```
In [35]: rf_scores = cross_val_score(decision_tree_classifier,
                                     all_inputs, all_classes, cv=10)

sb.boxplot(rf_scores)
sb.stripplot(rf_scores, jitter=True, color='white')
```

```
Out[35]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb3878b2e10>
```



Hmmm... that's a little boring by itself though. How about we compare another classifier to see how they perform?

We already know from previous projects that Random Forest classifiers usually work better than individual decision trees. A common problem that decision trees face is that they're prone to overfitting: They complexify to the point that they classify the training set near-perfectly, but fail to generalize to data they have not seen before.

Random Forest classifiers work around that limitation by creating a whole bunch of decision trees (hence "forest") — each trained on random subsets of training samples (drawn with replacement) and features (drawn without replacement) — and have the decision trees work together to make a more accurate classification.

Let that be a lesson for us: **Even in Machine Learning, we get better results when we work together!**

Let's see if a Random Forest classifier works better here.

The great part about scikit-learn is that the training, testing, parameter tuning, etc. process is the same for all models, so we only need to plug in the new classifier.

```
In [36]: from sklearn.ensemble import RandomForestClassifier

random_forest_classifier = RandomForestClassifier()

parameter_grid = {'n_estimators': [5, 10, 25, 50],
                  'criterion': ['gini', 'entropy'],
                  'max_features': [1, 2, 3, 4],
                  'warm_start': [True, False]}

cross_validation = StratifiedKFold(all_classes, n_folds=10)

grid_search = GridSearchCV(random_forest_classifier,
                           param_grid=parameter_grid,
                           cv=cross_validation)

grid_search.fit(all_inputs, all_classes)
print('Best score: {}'.format(grid_search.best_score_))
print('Best parameters: {}'.format(grid_search.best_params_))

grid_search.best_estimator_

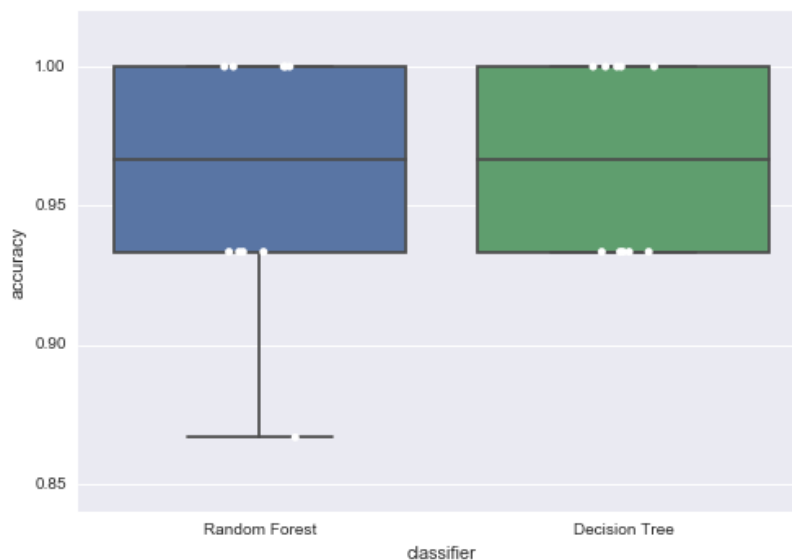
Best score: 0.96644295302
Best parameters: {'max_features': 1, 'n_estimators': 10, 'criterion': 'gini', 'warm_start': True}

Out[36]: RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                                max_depth=None, max_features=1, max_leaf_nodes=None,
                                min_samples_leaf=1, min_samples_split=2,
                                min_weight_fraction_leaf=0.0, n_estimators=10, n_jobs=1,
                                oob_score=False, random_state=None, verbose=0, warm_start=True)
```

Now we can compare their performance:

```
In [37]: random_forest_classifier = grid_search.best_estimator_  
  
rf_df = pd.DataFrame({'accuracy': cross_val_score(random_forest_classifi  
er,  
                                                    all_inputs, all_classe  
s, cv=10),  
                      'classifier': ['Random Forest'] * 10})  
dt_df = pd.DataFrame({'accuracy': cross_val_score(decision_tree_classifi  
er,  
                                                    all_inputs, all_classe  
s, cv=10),  
                      'classifier': ['Decision Tree'] * 10})  
both_df = rf_df.append(dt_df)  
  
sb.boxplot(x='classifier', y='accuracy', data=both_df)  
sb.stripplot(x='classifier', y='accuracy', data=both_df, jitter=True, co  
lor='white')
```

Out[37]: <matplotlib.axes._subplots.AxesSubplot at 0x7fb38c0b3450>



How about that? They both seem to perform about the same on this data set. This is probably because of the limitations of our data set: We have only 4 features to make the classification, and Random Forest classifiers excel when there's hundreds of possible features to look at. In other words, there wasn't much room for improvement with this data set.

Step 6: Reproducibility

[\[go back to the top \]](#)

Ensuring that our work is reproducible is the last and — arguably — most important step in any analysis. **As a rule, we shouldn't place much weight on a discovery that can't be reproduced.** As such, if our analysis isn't reproducible, we might as well not have done it.

Notebooks like this one go a long way toward making our work reproducible. Since we documented every step as we moved along, we have a written record of what we did and why we did it — both in text and code.

Beyond recording what we did, we should also document what software and hardware we used to perform our analysis. This typically goes at the top of our notebooks so our readers know what tools to use.

Sebastian Raschka (<http://sebastianraschka.com/>) created a handy [notebook tool \(https://github.com/rasbt/watermark\)](https://github.com/rasbt/watermark) for this:

```
In [49]: %install_ext https://raw.githubusercontent.com/rasbt/watermark/master/watermark/watermark.py
```

```
Installed watermark.py. To use it, type:  
%load_ext watermark
```

```
In [ ]: %load_ext watermark
```

```
In [54]: %watermark -a 'Randal S. Olson' -nmv --packages numpy,pandas,scikit-learn,matplotlib,Seaborn
```

```
ERROR: Line magic function `%watermark` not found.
```

Finally, let's extract the core of our work from Steps 1-5 and turn it into a single pipeline.

```
In [56]: %matplotlib inline
import pandas as pd
import seaborn as sb
from sklearn.ensemble import RandomForestClassifier
from sklearn.cross_validation import train_test_split
from sklearn.cross_validation import cross_val_score

# We can jump directly to working with the clean data
#because we saved our cleaned data set
iris_data_clean = pd.read_csv('iris-data-clean.csv')

# Testing our data: Our analysis will stop here if any
#of these assertions are wrong

# We know that we should only have three classes
assert len(iris_data_clean['class'].unique()) == 3

# We know that sepal lengths for 'Iris-versicolor'
#should never be below 2.5 cm
assert iris_data_clean.loc[iris_data_clean['class'] == 'Iris-versicolor',
                           'sepal_length_cm'].min() >= 2.5

# We know that our data set should have no missing measurements
assert len(iris_data_clean.loc[(iris_data_clean['sepal_length_cm'].isnull() |
                                (iris_data_clean['sepal_width_cm'].isnull() |
                                (iris_data_clean['petal_length_cm'].isnull() |
                                (iris_data_clean['petal_width_cm'].isnull()
                                )))) == 0

all_inputs = iris_data_clean[['sepal_length_cm', 'sepal_width_cm',
                              'petal_length_cm', 'petal_width_cm']].values
s
all_classes = iris_data_clean['class'].values
```

```
In [57]: # This is the classifier that came out of Grid Search
random_forest_classifier = RandomForestClassifier(bootstrap=True,
                                                class_weight=None, criterion='gini',
                                                max_depth=None, max_features=3,
                                                max_leaf_nodes=None,
                                                min_samples_leaf=1, min_samples_split=2,
                                                min_weight_fraction_leaf=0.0,
                                                n_estimators=5, n_jobs=1,
                                                oob_score=False, random_state=None,
                                                verbose=0, warm_start=True)

# All that's left to do now is plot the cross-validation scores
rf_classifier_scores = cross_val_score(random_forest_classifier,
                                      all_inputs, all_classes, cv=10)

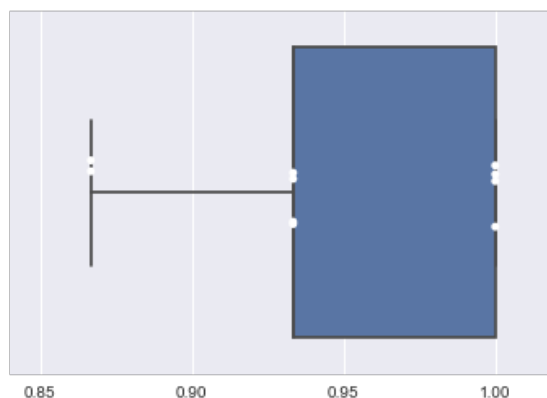
sb.boxplot(rf_classifier_scores)
sb.stripplot(rf_classifier_scores, jitter=True, color='white')

# ...and show some of the predictions from the classifier
(training_inputs,
 testing_inputs,
 training_classes,
 testing_classes) = train_test_split(all_inputs, all_classes, train_size
                                     =0.75)

random_forest_classifier.fit(training_inputs, training_classes)

for input_features, prediction, actual in zip(testing_inputs[:10],
                                             random_forest_classifier.predict(testing_inputs[
:10]),
                                             testing_classes[:10]):
    print('{}\t->\t{}\t(Actual: {})'.format(input_features, prediction,
                                             actual))

[ 6.9  3.1  4.9  1.5] --> Iris-versicolor (Actual: Iris-versicolor)
[ 5.6  3.   4.1  1.3] --> Iris-versicolor (Actual: Iris-versicolor)
[ 5.7  4.4  1.5  0.4] --> Iris-setosa   (Actual: Iris-setosa)
[ 6.9  3.1  5.4  2.1] --> Iris-virginica (Actual: Iris-virginica)
[ 5.6  2.9  3.6  1.3] --> Iris-versicolor (Actual: Iris-versicolor)
[ 4.9  3.1  1.5  0.1] --> Iris-setosa   (Actual: Iris-setosa)
[ 4.7  3.2  1.6  0.2] --> Iris-setosa   (Actual: Iris-setosa)
[ 4.8  3.4  1.6  0.25] --> Iris-setosa   (Actual: Iris-setosa)
[ 7.4  2.8  6.1  1.9] --> Iris-virginica (Actual: Iris-virginica)
[ 6.9  3.1  5.1  2.3] --> Iris-virginica (Actual: Iris-virginica)
```



There we have it: We have a complete and reproducible Machine Learning pipeline to demo to our head of data. We've met the success criteria that we set from the beginning (>90% accuracy), and our pipeline is flexible enough to handle new inputs or flowers when that data set is ready. Not bad for our first week on the job!

Conclusions

[\[go back to the top \]](#)

I hope you found this example notebook useful for your own work and learned at least one new trick by reading through it.

If you've spotted any errors or would like to contribute to this notebook, please don't hesitate to get in touch. I can be reached in the following ways:

- [Email \(http://www.randalolson.com/contact/\)](http://www.randalolson.com/contact/) me
- [Tweet \(https://twitter.com/randal_olson/\)](https://twitter.com/randal_olson/) at me
- [Submit an issue \(https://github.com/rhiever/Data-Analysis-and-Machine-Learning-Projects/issues\)](https://github.com/rhiever/Data-Analysis-and-Machine-Learning-Projects/issues) on GitHub
- Fork the [notebook repository \(https://github.com/rhiever/Data-Analysis-and-Machine-Learning-Projects/\)](https://github.com/rhiever/Data-Analysis-and-Machine-Learning-Projects/), make the fix/addition yourself, then send over a pull request

Further reading

[\[go back to the top \]](#)

This notebook covers a broad variety of topics but skips over many of the specifics. If you're looking to dive deeper into a particular topic, here's some recommended reading.

Data Science: William Chen compiled a [list of free books \(http://www.wzchen.com/data-science-books/\)](http://www.wzchen.com/data-science-books/) for newcomers to Data Science, ranging from the basics of R & Python to Machine Learning to interviews and advice from prominent data scientists.

Machine Learning: [/r/MachineLearning](#) has a useful [Wiki page \(https://www.reddit.com/r/MachineLearning/wiki/index\)](https://www.reddit.com/r/MachineLearning/wiki/index) containing links to online courses, books, data sets, etc. for Machine Learning. There's also a [curated list \(https://github.com/josephmisiti/awesome-machine-learning\)](https://github.com/josephmisiti/awesome-machine-learning) of Machine Learning frameworks, libraries, and software sorted by language.

Unit testing: Dive Into Python 3 has a [great walkthrough \(http://www.diveintopython3.net/unit-testing.html\)](http://www.diveintopython3.net/unit-testing.html) of unit testing in Python, how it works, and how it should be used

pandas has [several tutorials \(http://pandas.pydata.org/pandas-docs/stable/tutorials.html\)](http://pandas.pydata.org/pandas-docs/stable/tutorials.html) covering its myriad features.

scikit-learn has a [bunch of tutorials \(http://scikit-learn.org/stable/tutorial/index.html\)](http://scikit-learn.org/stable/tutorial/index.html) for those looking to learn Machine Learning in Python. Andreas Mueller's [scikit-learn workshop materials \(https://github.com/amueller/scipy_2015_sklearn_tutorial\)](https://github.com/amueller/scipy_2015_sklearn_tutorial) are top-notch and freely available.

matplotlib has many [books, videos, and tutorials \(http://matplotlib.org/resources/index.html\)](http://matplotlib.org/resources/index.html) to teach plotting in Python.

Seaborn has a [basic tutorial \(http://stanford.edu/~mwaskom/software/seaborn/tutorial.html\)](http://stanford.edu/~mwaskom/software/seaborn/tutorial.html) covering most of the statistical plotting features.

Acknowledgements

[\[go back to the top \]](#)

Many thanks to [Andreas Mueller](http://amueller.github.io/) (<http://amueller.github.io/>) for some of his [examples](https://github.com/amueller/scipy_2015_sklearn_tutorial) (https://github.com/amueller/scipy_2015_sklearn_tutorial) in the Machine Learning section. I drew inspiration from several of his excellent examples.

The photo of a flower with annotations of the petal and sepal was taken by [Eric Guinther](https://commons.wikimedia.org/wiki/File:Petal-sepal.jpg) (<https://commons.wikimedia.org/wiki/File:Petal-sepal.jpg>).

The photos of the various *Iris* flower types were taken by [Ken Walker](http://www.signa.org/index.pl?Display+Iris-setosa+2) (<http://www.signa.org/index.pl?Display+Iris-setosa+2>) and [Barry Glick](http://www.signa.org/index.pl?Display+Iris-virginica+3) (<http://www.signa.org/index.pl?Display+Iris-virginica+3>).