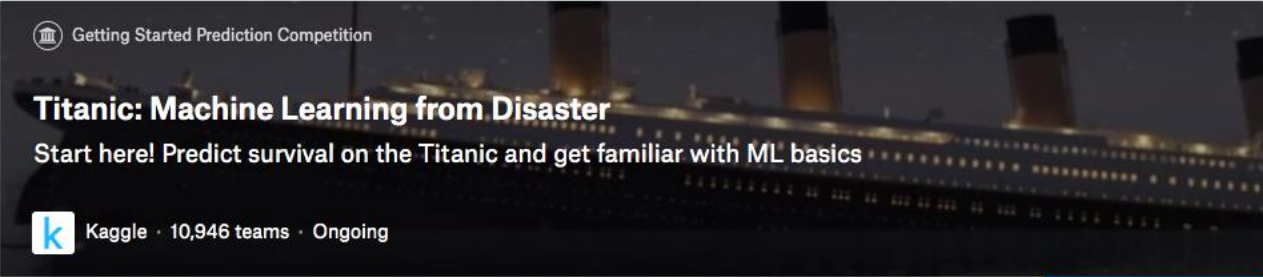# COMP 4331 Tut 3

Jiaxin Xie

# Outline

- Brief introduction to Kaggle
- Brief introduction to Pandas
- Demonstration of using pandas to preprocess the raw data
- Brief introduction to sklearn
- Review various plots

# What is Kaggle

Kaggle is a competition platform, where you can find data, examples, tutorials, even courses.

- Free
- Some of the projects even provide computing environments
- You can ask people there
- For more information, you can refer to: https://www.kaggle.com

# Titanic

# Experiment setup

- Start the Jupyter Notebook (jupyter notebook)

- Connect to it via your notebook via your browser.

- Create a new ipynb file (click 'new' button)

- In the first box, copy and paste all the codes in next slides to load used packages

- Press shift+enter to execute this box. (PS: you may get error saying that you cannot find a package, please install it in your server with "pip install [package name]")

# Experiment setup

```python
import numpy as np
import pandas as pd
from sklearn import tree
from sklearn.metrics import accuracy_score
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from IPython.display import Image as PImage
from subprocess import check_call
# import all the needed package
```

# Introduction to Pandas

- Powerful and productive python data analysis and management library.

- Short name of **Pan**el **Da**ta **S**ystem

- Open sourced by AQR capital management, LLC in late 2009

- Used by both academic and industry.

# Load data with pandas

Load the train and test data

```
In [26]:   1   # Loading the data
           2   train = pd.read_csv('train.csv')
           3   test = pd.read_csv('test.csv')
```

```
In [27]:   1   # Showing overview of the train dataset
           2   train.head(3)
```

Show the first three data example

Out[27]:

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| 2 | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |

# Preprocess data with Pandas (1)

```
In [28]:
    1  # Copy original dataset in case we need it later when digging into interesting features
    2  # WARNING: Beware of actually copying the dataframe instead of just referencing it
    3  # "original_train = train" will create a reference to the train variable (changes in 'train' will apply to 'origina
    4  original_train = train.copy() # Using 'copy()' allows to clone the dataset, creating a different object with
    5
    6  # Feature engineering steps taken from Sina and Anisotropic, with minor changes to avoid warn
    7  full_data = [train, test]
    8
    9  # Feature that tells whether a passenger had a cabin on the Titanic
   10  train['Has_Cabin'] = train["Cabin"].apply(lambda x: 0 if type(x) == float else 1)
   11  test['Has_Cabin'] = test["Cabin"].apply(lambda x: 0 if type(x) == float else 1)
   12
   13  # Create new feature FamilySize as a combination of SibSp and Parch
   14  for dataset in full_data:
   15      dataset['FamilySize'] = dataset['SibSp'] + dataset['Parch'] + 1
   16  # Create new feature IsAlone from FamilySize
   17  for dataset in full_data:
   18      dataset['IsAlone'] = 0
   19      dataset.loc[dataset['FamilySize'] == 1, 'IsAlone'] = 1
   20  # Remove all NULLS in the Embarked column
   21  for dataset in full_data:
   22      dataset['Embarked'] = dataset['Embarked'].fillna('S')
   23  # Remove all NULLS in the Fare column
   24  for dataset in full_data:
   25      dataset['Fare'] = dataset['Fare'].fillna(train['Fare'].median())
   26
   27  # Remove all NULLS in the Age column
   28  for dataset in full_data:
   29      age_avg = dataset['Age'].mean()
   30      age_std = dataset['Age'].std()
   31      age_null_count = dataset['Age'].isnull().sum()
   32      age_null_random_list = np.random.randint(age_avg - age_std, age_avg + age_std, size=age_null_count)
   33      # Next line has been improved to avoid warning
   34      dataset.loc[np.isnan(dataset['Age']), 'Age'] = age_null_random_list
   35      dataset['Age'] = dataset['Age'].astype(int)
```

Put the train and test data together

Convert "family size" feature to integer

Create a new feature "IsAlone"

Replace empty data with average

# Preprocess data with Pandas (2)

```python
36  # Define function to extract titles from passenger names
37  def get_title(name):
38      title_search = re.search(' ([A-Za-z]+)\.', name)
39      # If the title exists, extract and return it.
40      if title_search:
41          return title_search.group(1)
42      return ""
43
44  for dataset in full_data:
45      dataset['Title'] = dataset['Name'].apply(get_title)
46  # Group all non-common titles into one single grouping "Rare"
47  for dataset in full_data:
48      dataset['Title'] = dataset['Title'].replace(['Lady', 'Countess','Capt', 'Col','Don', 'Dr', 'Major', 'Rev', 'Sir
49
50      dataset['Title'] = dataset['Title'].replace('Mlle', 'Miss')
51      dataset['Title'] = dataset['Title'].replace('Ms', 'Miss')
52      dataset['Title'] = dataset['Title'].replace('Mme', 'Mrs')
53
54  for dataset in full_data:
55      # Mapping Sex
56      dataset['Sex'] = dataset['Sex'].map( {'female': 0, 'male': 1} )
57
58      # Mapping titles
59      title_mapping = {"Mr": 1, "Master": 2, "Mrs": 3, "Miss": 4, "Rare": 5}
60      dataset['Title'] = dataset['Title'].map(title_mapping)
61      dataset['Title'] = dataset['Title'].fillna(0)
62
63      # Mapping Embarked
64      dataset['Embarked'] = dataset['Embarked'].map( {'S': 0, 'C': 1, 'Q': 2} )
65
66      # Mapping Fare
67      dataset.loc[ dataset['Fare'] <= 7.91, 'Fare']                          = 0
68      dataset.loc[(dataset['Fare'] > 7.91) & (dataset['Fare'] <= 14.454), 'Fare'] = 1
69      dataset.loc[(dataset['Fare'] > 14.454) & (dataset['Fare'] <= 31), 'Fare']   = 2
70      dataset.loc[ dataset['Fare'] > 31, 'Fare']                             = 3
71      dataset['Fare'] = dataset['Fare'].astype(int)
72
73      # Mapping Age
74      dataset.loc[ dataset['Age'] <= 16, 'Age']                          = 0
75      dataset.loc[(dataset['Age'] > 16) & (dataset['Age'] <= 32), 'Age'] = 1
76      dataset.loc[(dataset['Age'] > 32) & (dataset['Age'] <= 48), 'Age'] = 2
77      dataset.loc[(dataset['Age'] > 48) & (dataset['Age'] <= 64), 'Age'] = 3
78      dataset.loc[ dataset['Age'] > 64, 'Age'] ;
79
80      del dataset['Name']
81      del dataset['Ticket']
82      del dataset['Cabin']
```

Define a function to clean the data

Clean all the title feature

Replace all the feature with number

# Introduction to scikit-learn (sklearn)

Extensions to Scipy (Scientific Python) are called Scikits. Scikie-learn provides machine learning algorithm.

- Algorithms for supervised and unsupervised learning

- Built on Scipy and Numpy

- Standard Python API interface

- Sits on top of c libraries, LAPACK, LibSVM, and Cython

- Open Source

For more information, you can refer to: https://scikit-learn.org/stable/

# Build a decision tree model with sklearn

sklearn. tree **.DecisionTreeClassifier**

*class* sklearn.tree. **DecisionTreeClassifier** (*criterion='gini', splitter='best', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, class_weight=None, presort=False*)

[source]

A decision tree classifier.

**Examples**

```
from sklearn.datasets import load_iris
from sklearn.model_selection import cross_val_score
from sklearn.tree import DecisionTreeClassifier
clf = DecisionTreeClassifier(random_state=0)
iris = load_iris()
cross_val_score(clf, iris.data, iris.target, cv=10)
```

# Build a decision tree model with sklearn

**Methods**

| | |
|---|---|
| **apply** (self, X[, check_input]) | Returns the index of the leaf that each sample is predicted as. |
| **decision_path** (self, X[, check_input]) | Return the decision path in the tree |
| **fit** (self, X, y[, sample_weight, ...]) | Build a decision tree classifier from the training set (X, y). |
| **get_depth** (self) | Returns the depth of the decision tree. |
| **get_n_leaves** (self) | Returns the number of leaves of the decision tree. |
| **get_params** (self[, deep]) | Get parameters for this estimator. |
| **predict** (self, X[, check_input]) | Predict class or regression value for X. |
| **predict_log_proba** (self, X) | Predict class log-probabilities of the input samples X. |
| **predict_proba** (self, X[, check_input]) | Predict class probabilities of the input samples X. |
| **score** (self, X, y[, sample_weight]) | Returns the mean accuracy on the given test data and labels. |
| **set_params** (self, \*\*params) | Set the parameters of this estimator. |

# Visualize Decision Tree

sklearn.tree.**export_graphviz**

sklearn.tree. **export_graphviz** (*decision_tree, out_file=None, max_depth=None, feature_names=None, class_names=None, label='all', filled=False, leaves_parallel=False, impurity=True, node_ids=False, proportion=False, rotate=False, rounded=False, special_characters=False, precision=3*) ¶                    [source]

**Examples**

```
from sklearn.datasets import load_iris
from sklearn import tree
```

```
clf = tree.DecisionTreeClassifier()
iris = load_iris()
```

```
clf = clf.fit(iris.data, iris.target)
tree.export_graphviz(clf)
```

# Visualize Decision Tree

Export a decision tree in DOT format.

This function generates a GraphViz representation of the decision tree, which is then written into `out_file`. Once exported, graphical renderings can be generated using, for example:

```
$ dot -Tps tree.dot -o tree.ps      (PostScript format)
$ dot -Tpng tree.dot -o tree.png    (PNG format)
```

# Implement PCA

sklearn.decomposition.**PCA**

*class* `sklearn.decomposition.` **PCA** (*n_components=None, copy=True, whiten=False, svd_solver='auto', tol=0.0, iterated_power='auto', random_state=None*)  [source]

```python
import numpy as np
from sklearn.decomposition import PCA
X = np.array([[-1, -1], [-2, -1], [-3, -2], [1, 1], [2, 1], [3, 2]])
pca = PCA(n_components=2)
pca.fit(X)

print(pca.explained_variance_ratio_)

print(pca.singular_values_)
```
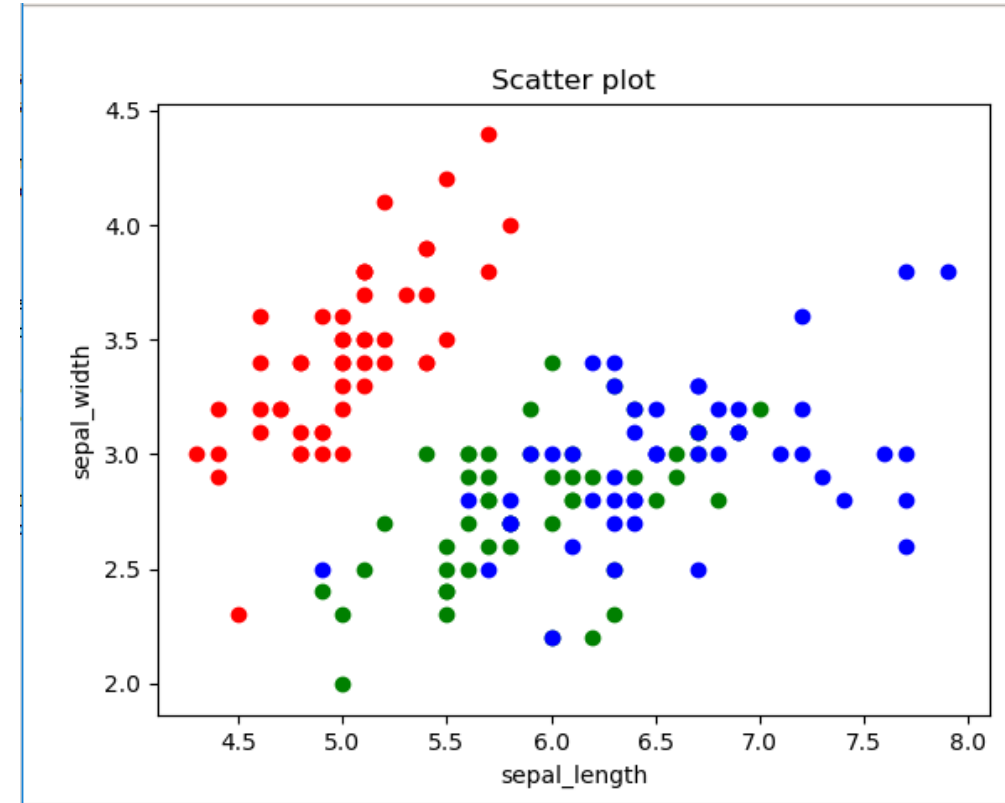
In tutorial 1,  we also release code to implement PCA.

# Scatter plot

```python
import matplotlib.pyplot as plt
#scatter plot
colors = {'Setosa':'r', 'Versicolor':'g', 'Virginica':'b'}
# create a figure and axis
fig, ax = plt.subplots()
# plot each data-point
for i in range(len(iris['sepal_length'])):
    ax.scatter(iris['sepal_length'][i], iris['sepal_width'][i],\
    color=colors[iris['class'][i]])
# set a title and labels
ax.set_title('Scatter Plot')
ax.set_xlabel('sepal_length')
ax.set_ylabel('sepal_width')
plt.show()
```
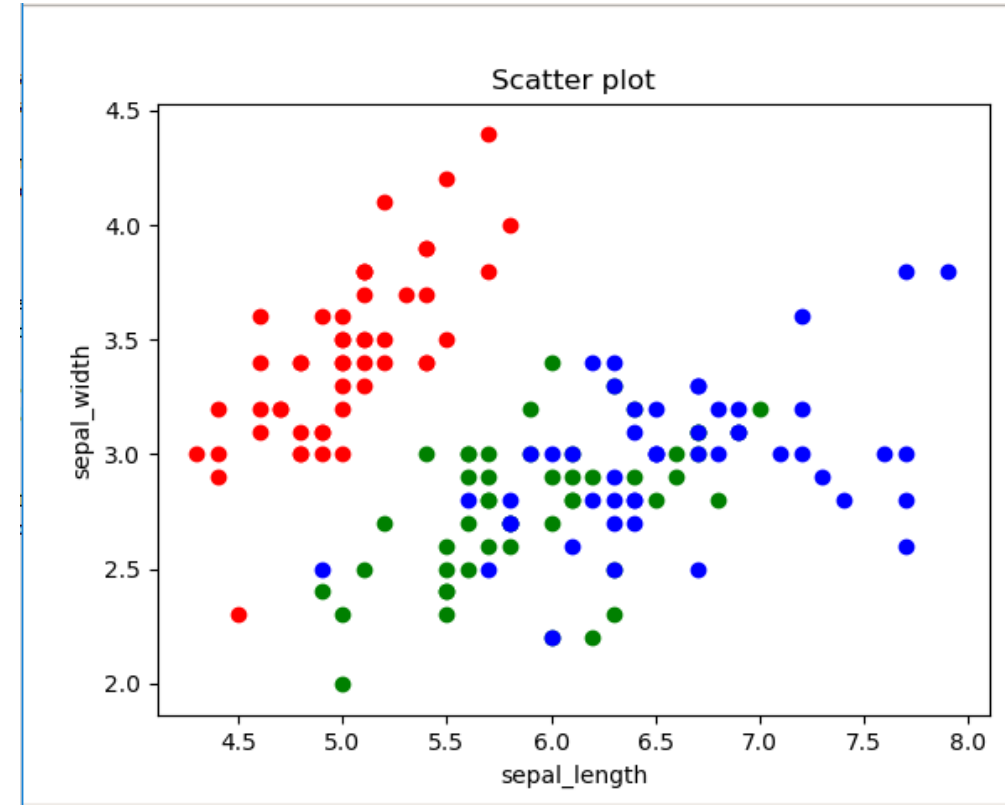


How to draw a scatter plot with sepal_length
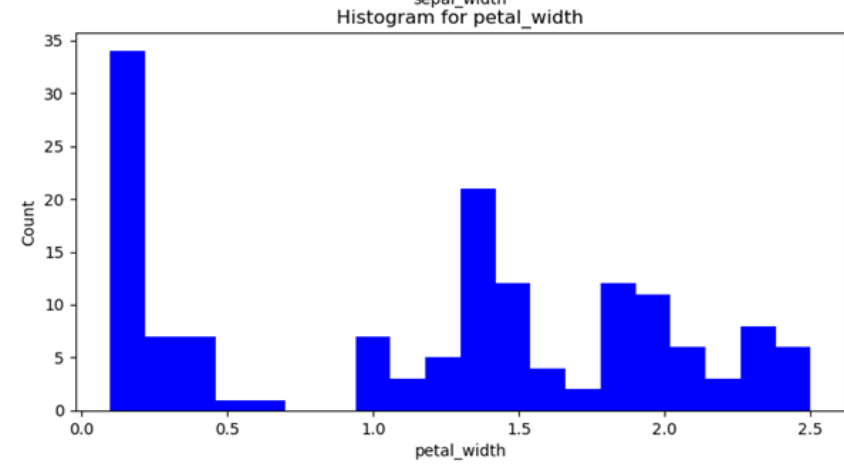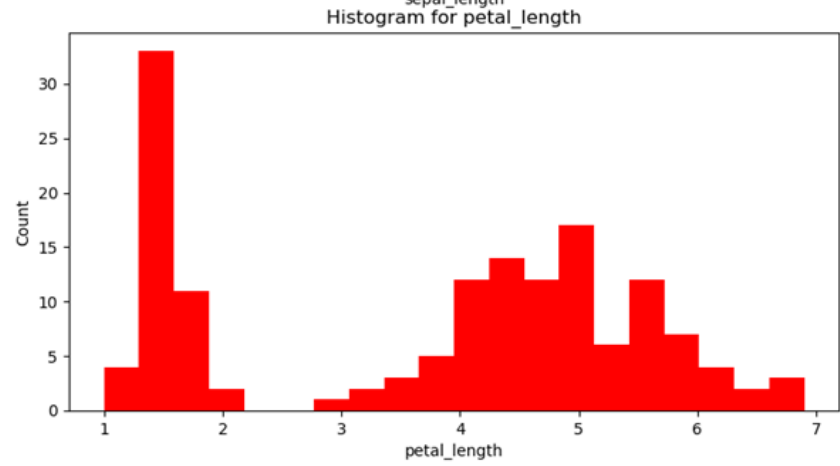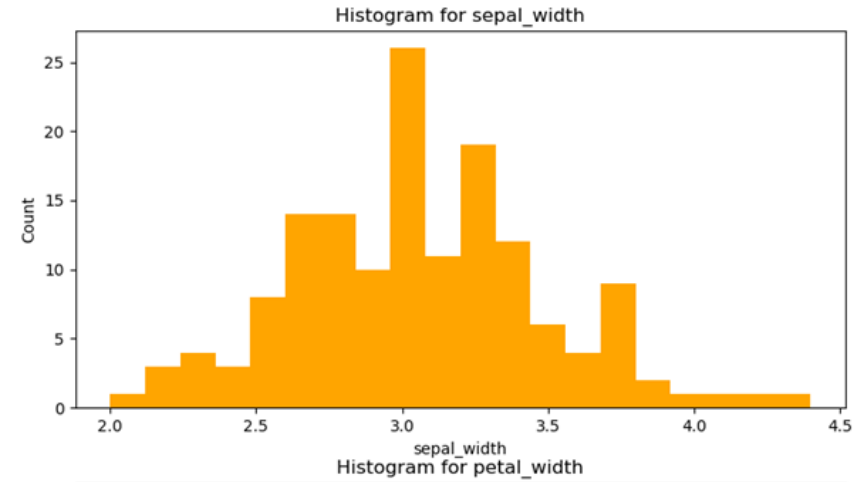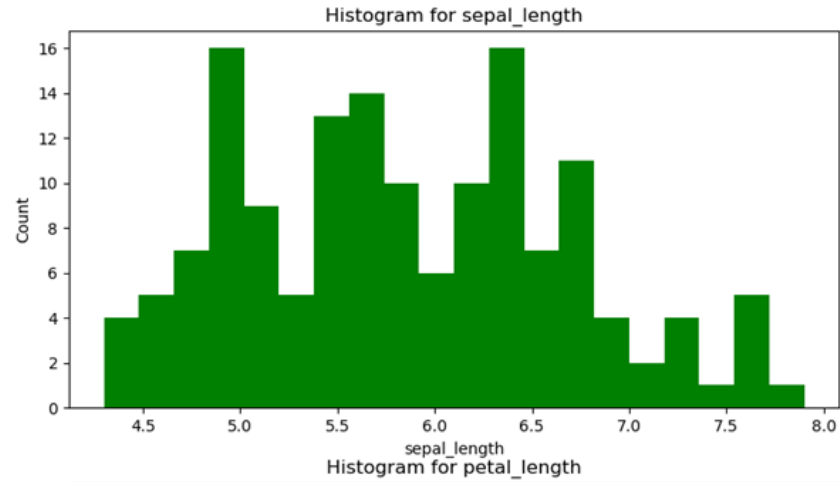for the x-axis and petal_length for the y_axis?

# Scatter plot

```python
import matplotlib.pyplot as plt
#scatter plot
colors = {'Setosa':'r', 'Versicolor':'g', 'Virginica':'b'}
# create a figure and axis
fig, ax = plt.subplots()
# plot each data-point
for i in range(len(iris['sepal_length'])):
    ax.scatter(iris['sepal_length'][i], iris['sepal_width'][i],\
    color=colors[iris['class'][i]])
# set a title and labels
ax.set_title('Scatter Plot')
ax.set_xlabel('sepal_length')
ax.set_ylabel('sepal_width')
plt.show()
```



How to draw a scatter plot with sepal_length
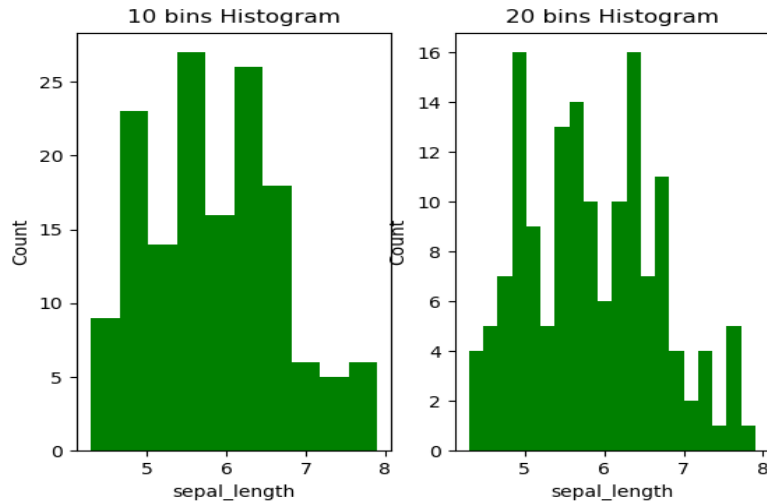for the x-axis and petal_length for the y_axis?

# Histogram

# Histogram

```
1    import matplotlib.pyplot as plt
2    plt.figure()
3    x = iris["sepal_length"]
4    plt.hist(x, bins = 20, color = "green")
5    plt.title("Histogram for sepal_length")
6    plt.xlabel("sepal_length")
7    plt.ylabel("Count")
```

Function: matplotlib.pyplot.hist  More details plz see https://matplotlib.org/3.1.1/api/_as_gen/matplotlib.pyplot.hist.html
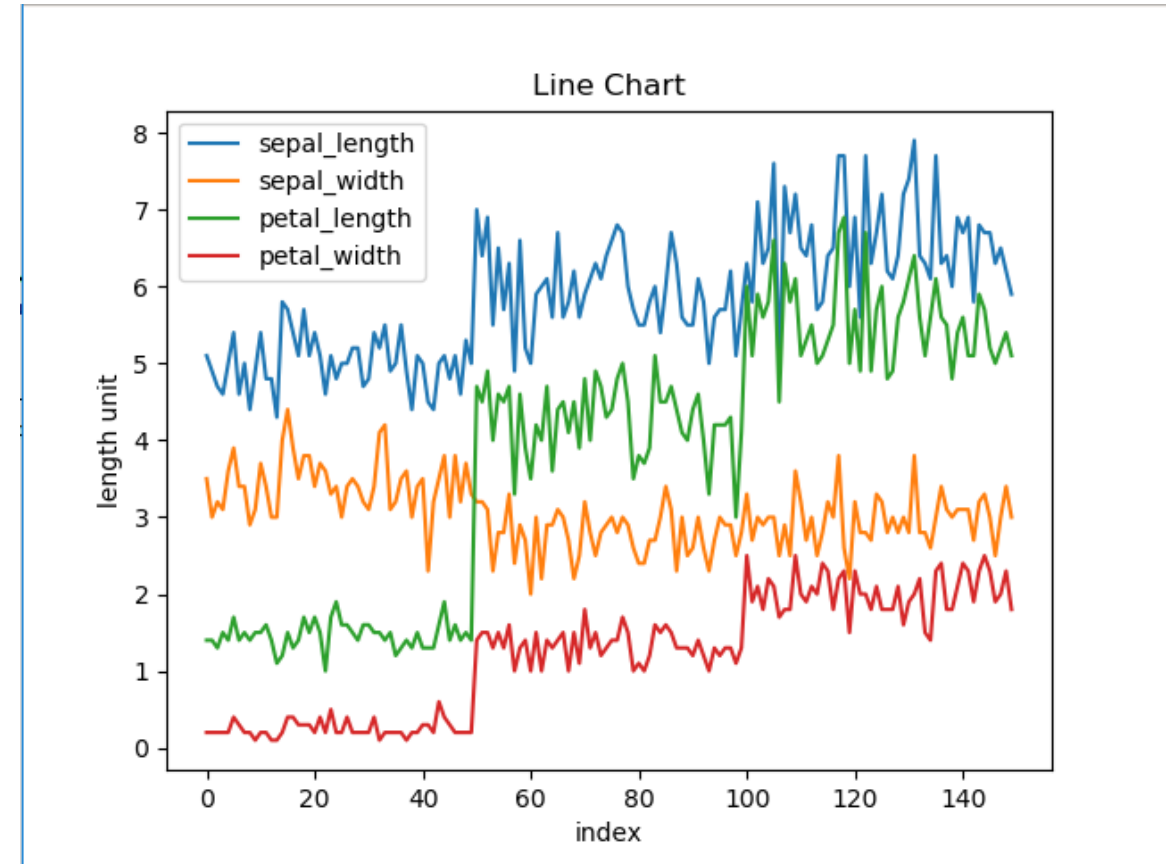
Change bins? plt.hist(x, bins = 10, color = "green")



Try to change other parameters, such as range.

# Line Chart

```
1   import matplotlib.pyplot as plt
2   #Line Chart
3   columns = iris.columns.drop(['class'])
4   # create x data
5   x_data = range(0, iris.shape[0])
6   # create figure and axis
7   fig, ax = plt.subplots()
8   # plot each column
9   for column in columns:
10      ax.plot(x_data, iris[column])
11  # set title and legend
12  ax.set_title('Line Chart')
13  plt.xlabel('index')
14  plt.ylabel('length unit')
15  ax.legend(['sepal_length', 'sepal_width', 'petal_length', 'petal_width'])
16  plt.show()
```



How to draw a Line Chart if we only want to plot sepal_length and sepal_width?

# Box Plot

```
1    import matplotlib.pyplot as plt
2    #Box plot
3    plt.figure()
4    new_iris=iris[["sepal_length", "sepal_width",\
5      "petal_length", "petal_width"]]
6    new_iris.boxplot()
7    plt.title('Box Plot')
8    plt.show()
```

Same question as Line Chart. if we only want to include sepal_length and sepal_width?