

## ITEC 6720 - Final Project: Screenshots Log

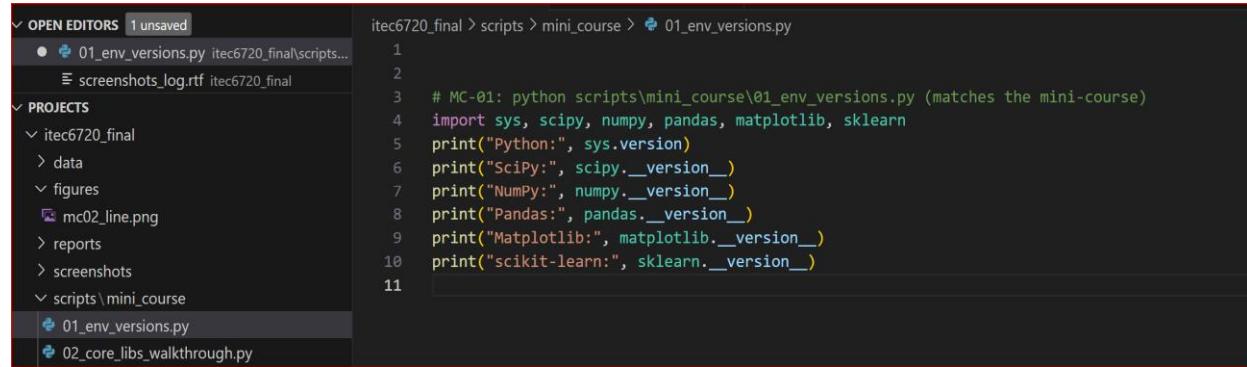
1. Mini-Course (MC-01 to MC-14)
2. Iris Step-by-Step Project

### 1) Mini-Course

#### MC-01 -- Environment Check

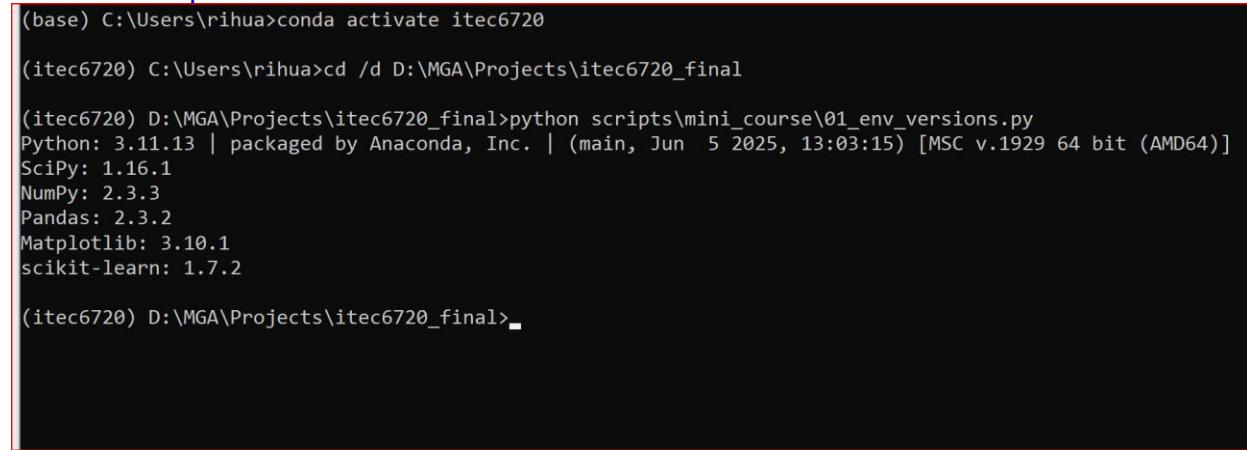
Caption: MC-01 -- Python, SciPy, NumPy, Pandas, Matplotlib, scikit-learn versions printed in the (itec6720) environment.

#### VS code screenshot



```
itec6720_final > scripts > mini_course > 01_env_versions.py
1
2
3  # MC-01: python scripts\mini_course\01_env_versions.py (matches the mini-course)
4  import sys, scipy, numpy, pandas, matplotlib, sklearn
5  print("Python:", sys.version)
6  print("SciPy:", scipy.__version__)
7  print("NumPy:", numpy.__version__)
8  print("Pandas:", pandas.__version__)
9  print("Matplotlib:", matplotlib.__version__)
10 print("scikit-learn:", sklearn.__version__)
11
```

#### Anaconda output screenshot



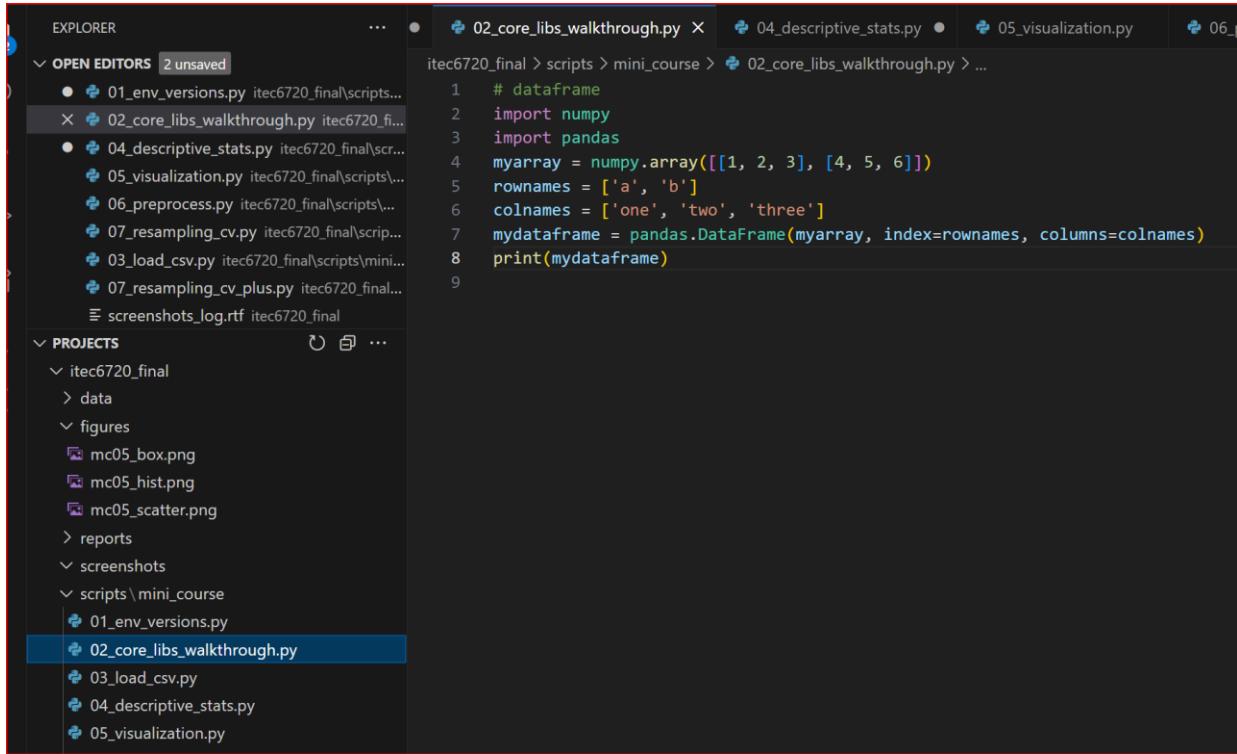
```
(base) C:\Users\rihua>conda activate itec6720
(itec6720) C:\Users\rihua>cd /d D:\MGA\Projects\itec6720_final
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\01_env_versions.py
Python: 3.11.13 | packaged by Anaconda, Inc. | (main, Jun 5 2025, 13:03:15) [MSC v.1929 64 bit (AMD64)]
SciPy: 1.16.1
NumPy: 2.3.3
Pandas: 2.3.2
Matplotlib: 3.10.1
scikit-learn: 1.7.2

(itec6720) D:\MGA\Projects\itec6720_final>
```

### MC-02 -- Core Libs Tour

Caption: MC-02 — Core libs tour: created a small NumPy array and printed a labeled 2×3 Pandas DataFrame (rows a,b; columns one,two,three).

## Code



```
1 # dataframe
2 import numpy
3 import pandas
4 myarray = numpy.array([[1, 2, 3], [4, 5, 6]])
5 rownames = ['a', 'b']
6 colnames = ['one', 'two', 'three']
7 mydataframe = pandas.DataFrame(myarray, index=rownames, columns=colnames)
8 print(mydataframe)
9
```

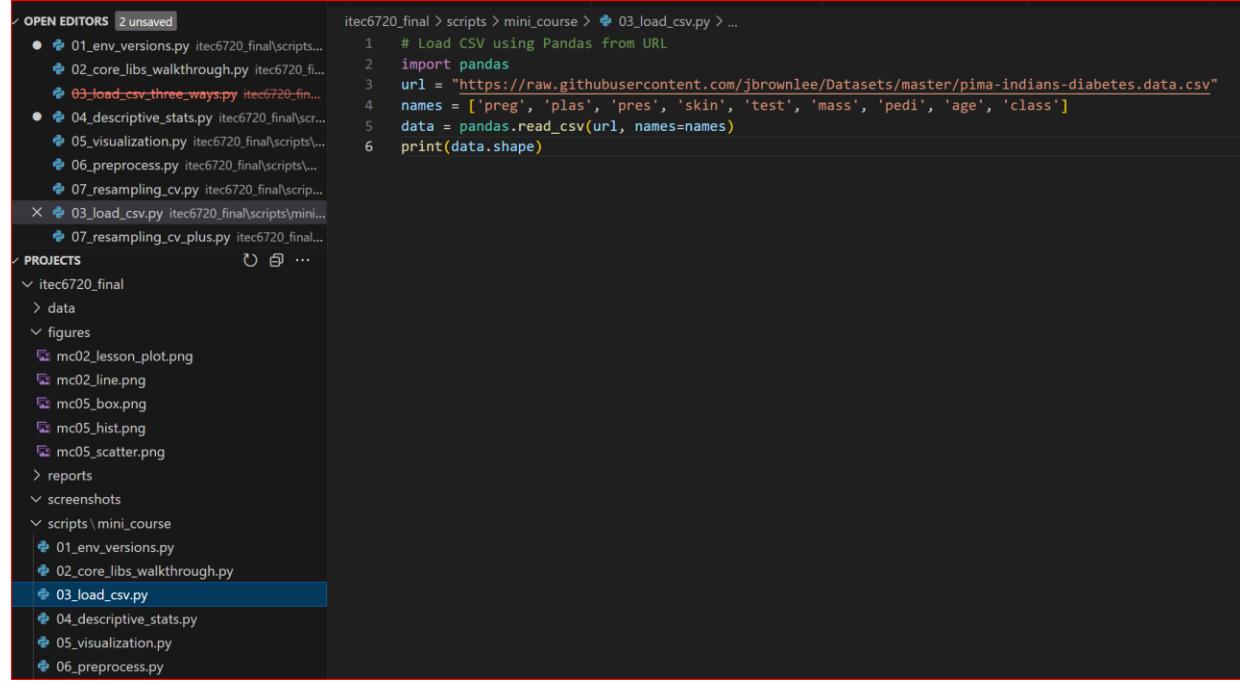
## Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\02_core_libs_walkthrough.py
   one  two  three
a    1    2    3
b    4    5    6
```

## MC-03 -- Load CSV

Caption: MC-03 — Loaded the Pima Indians Diabetes CSV from the URL with pandas.read\_csv (using column names) and verified the shape (768, 9).

### Code



The screenshot shows a Jupyter Notebook interface with the following details:

- OPEN EDITORS:** 2 unsaved
- PROJECTS:** itec6720\_final
- Contents of itec6720\_final:**
  - data
  - figures
    - mc02\_lesson\_plot.png
    - mc02\_line.png
    - mc05\_box.png
    - mc05\_hist.png
    - mc05\_scatter.png
  - reports
  - screenshots
  - scripts\mini\_course
    - 01\_env\_versions.py
    - 02\_core\_libs\_walkthrough.py
    - 03\_load\_csv.py
    - 04\_descriptive\_stats.py
    - 05\_visualization.py
    - 06\_preprocess.py

The code cell for `03_load_csv.py` contains the following Python code:

```
1 # Load CSV using Pandas from URL
2 import pandas
3 url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
4 names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
5 data = pandas.read_csv(url, names=names)
6 print(data.shape)
```

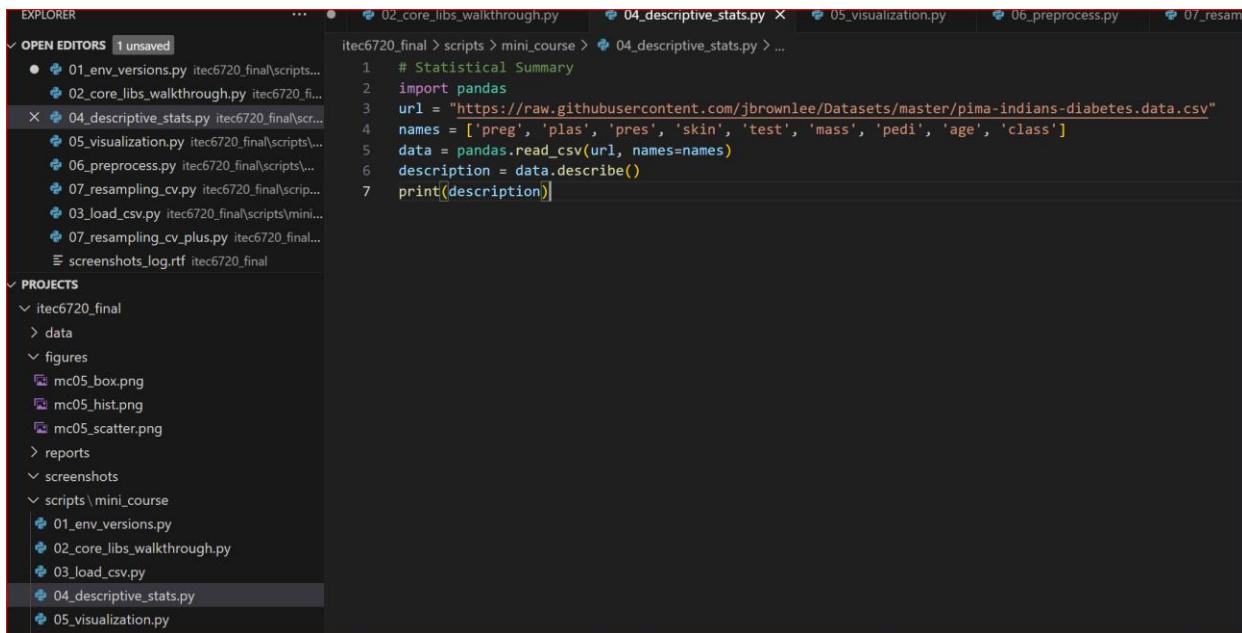
### Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\03_load_csv.py
(768, 9)
```

## MC-04 -- Descriptive Statistics

Caption: loaded the Pima Indians Diabetes dataset from the URL with pandas.read\_csv and printed the dataset's summary using data.describe().

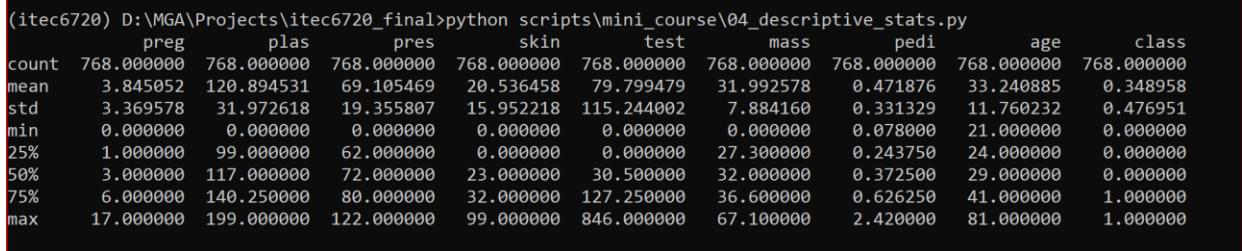
### Code



The screenshot shows a code editor with a dark theme. The left sidebar is titled 'EXPLORER' and shows a project structure for 'itec6720\_final'. The 'scripts' folder contains several Python files: '01\_env\_versions.py', '02\_core\_libs\_walkthrough.py', '04\_descriptive\_stats.py' (which is the active file), '05\_visualization.py', '06\_preprocess.py', '07\_resampling\_cv.py', '03\_load\_csv.py', '07\_resampling\_cv\_plus.py', and 'screenshots\_log.rtf'. The 'PROJECTS' section shows subfolders 'data', 'figures' (containing 'mc05\_box.png', 'mc05\_hist.png', 'mc05\_scatter.png'), 'reports', 'screenshots', and 'scripts\mini\_course' (containing '01\_env\_versions.py', '02\_core\_libs\_walkthrough.py', '03\_load\_csv.py', '04\_descriptive\_stats.py' (selected), and '05\_visualization.py'). The main editor area displays the following Python code:

```
1  # Statistical Summary
2  import pandas
3  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
4  names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
5  data = pandas.read_csv(url, names=names)
6  description = data.describe()
7  print(description)
```

### Output



The screenshot shows a terminal window with a red border. The command 'python scripts\mini\_course\04\_descriptive\_stats.py' is run, and the output is a table of descriptive statistics for the Pima Indians Diabetes dataset. The table includes columns for 'count', 'mean', 'std', 'min', '25%', '50%', '75%', and 'max' for each of the nine features: 'preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', and 'class'.

	preg	plas	pres	skin	test	mass	pedi	age	class
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	0.471876	33.240885	0.348958
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	21.000000	0.000000
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	24.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	29.000000	0.000000
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	41.000000	1.000000
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000	1.000000

## MC-05 -- Visualization

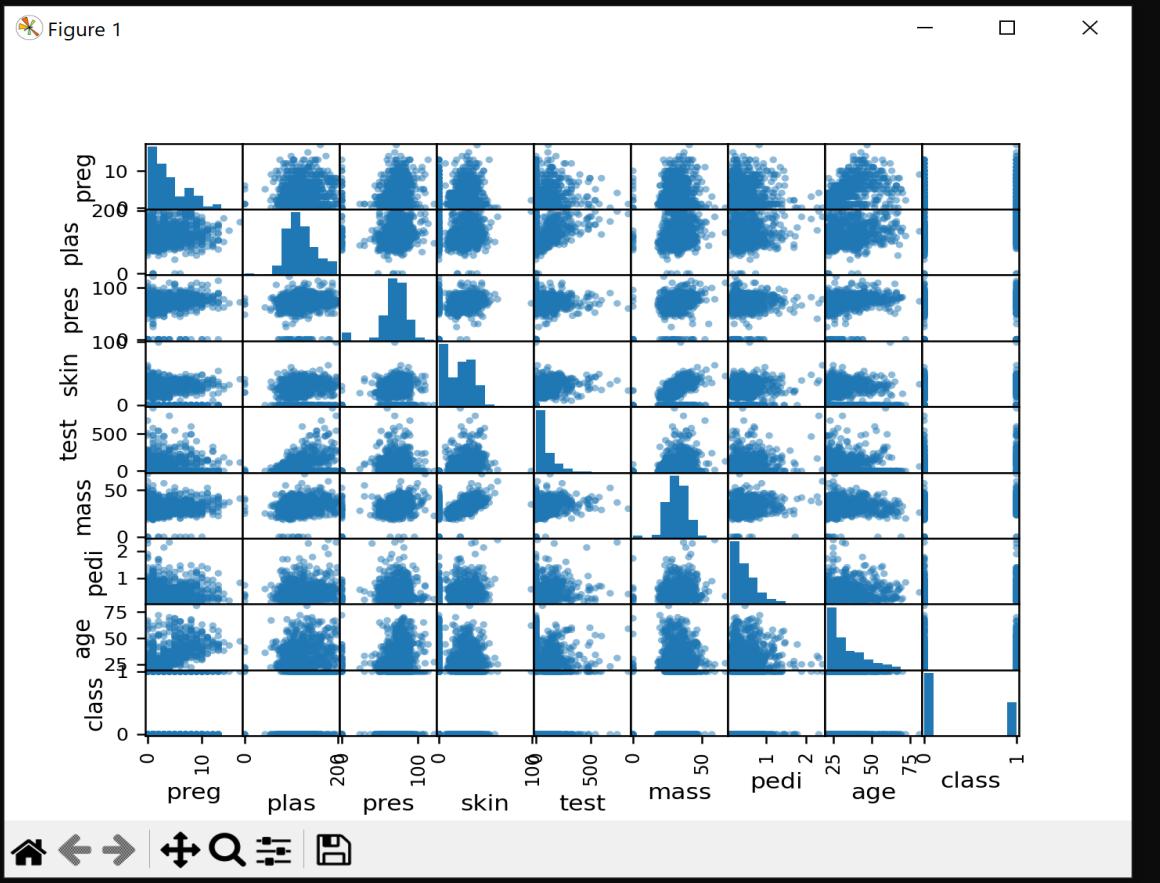
Caption: MC-05 —loaded the Pima dataset and produced a **scatter-plot matrix** of all attributes using `pandas.plotting.scatter_matrix` (displayed with Matplotlib).

### Code

```
itec6720_final > scripts > mini_course > 05_visualization.py > ...
1  # Scatter Plot Matrix
2  import matplotlib.pyplot as plt
3  import pandas
4  from pandas.plotting import scatter_matrix
5  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
6  names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
7  data = pandas.read_csv(url, names=names)
8  scatter_matrix(data)
9  plt.show()
```

### Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\05_visualization.py
```



## MC-06 -- Preprocessing

Caption:MC-06 —standardized the Pima features  $X$  with scikit-learn's StandardScaler (mean 0, std 1) and printed the first 5 rows of the rescaled data.

### Code

```
itec6720_final > scripts > mini_course > 06_preprocess.py > ...
1  # Standardize data (0 mean, 1 stdev)
2  from sklearn.preprocessing import StandardScaler
3  import pandas
4  import numpy
5  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
6  names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
7  dataframe = pandas.read_csv(url, names=names)
8  array = dataframe.values
9  # separate array into input and output components
10 X = array[:,0:8]
11 Y = array[:,8]
12 scaler = StandardScaler().fit(X)
13 rescaledX = scaler.transform(X)
14 # summarize transformed data
15 numpy.set_printoptions(precision=3)
16 print(rescaledX[0:5,:])
```

### Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\06_preprocess.py
[[ 0.64  0.848  0.15  0.907 -0.693  0.204  0.468  1.426]
 [-0.845 -1.123 -0.161  0.531 -0.693 -0.684 -0.365 -0.191]
 [ 1.234  1.944 -0.264 -1.288 -0.693 -1.103  0.604 -0.106]
 [-0.845 -0.998 -0.161  0.155  0.123 -0.494 -0.921 -1.042]
 [-1.142  0.504 -1.505  0.907  0.766  1.41   5.485 -0.02 ]]
```

## MC-07 -- Resampling

Caption: MC-07 — evaluated Logistic Regression on the Pima dataset using 10-fold cross-validation and reported mean accuracy  $\pm$  standard deviation.

### Code

```
itec6720_final > scripts > mini_course > 07_resampling_cv.py > ...
1  # MC-07 (basic): Evaluate using 10-fold CV (matches the mini-course)
2  from pandas import read_csv
3  from sklearn.model_selection import KFold, cross_val_score
4  from sklearn.linear_model import LogisticRegression
5
6  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
7  names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
8  dataframe = read_csv(url, names=names)
9  array = dataframe.values
10 X, y = array[:, 0:8], array[:, 8]
11
12 kfold = KFold(n_splits=10, random_state=7, shuffle=True)
13 model = LogisticRegression(solver="liblinear", max_iter=1000, random_state=42)
14 results = cross_val_score(model, X, y, cv=kfold, scoring="accuracy")
15
16 print("Accuracy: %.3f%% (%.3f%%)" % (results.mean()*100.0, results.std()*100.0))
17 |
```

### Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\07_resampling_cv.py
Accuracy: 77.086% (5.091%)
```

## MC-08 -- Metrics

Caption: MC-08 — Accuracy via 10-fold CV and LogLoss (lower is better) for Logistic Regression on the Pima dataset; KFold uses shuffle=True, random\_state=7.

### Code

```
itec6720_final > scripts > mini_course > 08_metrics.py > ...
1  # MC-08: Algorithm evaluation metrics (accuracy + log loss)
2  from pandas import read_csv
3  from sklearn.model_selection import KFold, cross_val_score
4  from sklearn.linear_model import LogisticRegression
5  import numpy as np
6
7  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
8  names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
9
10 df = read_csv(url, names=names)
11 X = df.iloc[:, 0:8].values
12 y = df.iloc[:, 8].values
13
14 # FIX: enable shuffling if you want to use random_state
15 kfold = KFold(n_splits=10, shuffle=True, random_state=7)
16
17 model = LogisticRegression(solver='liblinear', max_iter=1000)
18
19 # Accuracy (higher is better)
20 acc = cross_val_score(model, X, y, cv=kfold, scoring="accuracy")
21 print("Accuracy: %.3f%% (%.3f%%)" % (acc.mean() * 100.0, acc.std() * 100.0))
22
23 # Log loss (lower is better) -> cross_val_score returns *negative* log loss
24 neg_ll = cross_val_score(model, X, y, cv=kfold, scoring="neg_log_loss")
25 print("LogLoss: %.3f (%.3f)" % (-neg_ll.mean(), neg_ll.std()))
26
```

### Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\08_metrics.py
Accuracy: 77.086% (5.091%)
LogLoss: 0.494 (%0.042)
```

## MC-09 -- Spot-Check Algorithms

Caption: MC-09 — Spot-check: KNN regressor on Boston Housing with 10-fold CV; report mean  $\pm$  std RMSE.

### Code

```
itec6720_final > scripts > mini_course > 09_spot_check.py > ...
1  from pandas import read_csv
2  from sklearn.model_selection import KFold, cross_val_score
3  from sklearn.neighbors import KNeighborsRegressor
4  import numpy as np
5
6  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/housing.csv"
7  names = ['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT', 'MEDV']
8  df = read_csv(url, sep=r"\s+", names=names)
9
10 X, y = df.iloc[:, :13].values, df.iloc[:, 13].values
11 cv = KFold(n_splits=10, shuffle=True, random_state=7)
12 model = KNeighborsRegressor()
13
14 scores = cross_val_score(model, X, y, cv=cv, scoring='neg_mean_squared_error')
15 rmse = np.sqrt(-scores)
16 print(f"RMSE: {rmse.mean():.3f} ({rmse.std():.3f})")
17
```

### Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\09_spot_check.py
RMSE: 6.130 (1.131)
```

## MC-10 -- Model Comparison

Caption: MC-10 -- Comparing Logistic Regression vs. LDA on the Pima Diabetes dataset using 10-fold cross-validation and reporting mean  $\pm$  std accuracy to choose the more reliable model.

### Code

```
itec6720_final > scripts > mini_course > 10_model_compare.py > ...
1  # Compare Algorithms
2  from pandas import read_csv
3  from sklearn.model_selection import KFold
4  from sklearn.model_selection import cross_val_score
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
7  # load dataset
8  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
9  names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
10 dataframe = read_csv(url, names=names)
11 array = dataframe.values
12 X = array[:,0:8]
13 Y = array[:,8]
14 # prepare models
15 models = []
16 models.append(('LR', LogisticRegression(solver='liblinear')))
17 models.append(('LDA', LinearDiscriminantAnalysis()))
18 # evaluate each model in turn
19 results = []
20 names = []
21 scoring = 'accuracy'
22 for name, model in models:
23
24     cv = KFold(n_splits=10, shuffle=True, random_state=7)
25     cv_results = cross_val_score(model, X, Y, cv=cv, scoring=scoring)
26     results.append(cv_results)
27     names.append(name)
28     msg = "%s: %f (%f)" % (name, cv_results.mean(), cv_results.std())
29     print(msg)
```

### Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\10_model_compare.py
LR: 0.770865 (0.050905)
LDA: 0.766969 (0.047966)
```

## MC-11 -- Tuning

Caption: Tuning model hyperparameters with GridSearchCV: try a set of alpha values for Ridge using cross-validation, report the best CV score, and select the alpha that generalizes best.

### Code

```
1  # Grid Search for Algorithm Tuning
2  from pandas import read_csv
3  import numpy
4  from sklearn.linear_model import Ridge
5  from sklearn.model_selection import GridSearchCV
6  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
7  names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
8  dataframe = read_csv(url, names=names)
9  array = dataframe.values
10 X = array[:,0:8]
11 Y = array[:,8]
12 alphas = numpy.array([1,0.1,0.01,0.001,0.0001,0])
13 param_grid = dict(alpha=alphas)
14 model = Ridge()
15 grid = GridSearchCV(estimator=model, param_grid=param_grid, cv=3)
16 grid.fit(X, Y)
17 print(grid.best_score_)
18 print(grid.best_estimator_.alpha)
19
20
```

### Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\11_tuning.py
0.27961755931297233
1.0
```

## MC-12 -- Ensembles

Caption: MC-12 -- Evaluate a 100-tree RandomForestClassifier on the Pima Diabetes dataset using shuffled 10-fold cross-validation (max\_features=3) and print the mean CV accuracy.

### Code

```
itec6720_final > scripts > mini_course > 12_ensembles.py > ...
1  # Random Forest Classification
2  from pandas import read_csv
3  from sklearn.model_selection import KFold
4  from sklearn.model_selection import cross_val_score
5  from sklearn.ensemble import RandomForestClassifier
6  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
7  names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
8  dataframe = read_csv(url, names=names)
9  array = dataframe.values
10 X = array[:,0:8]
11 Y = array[:,8]
12 num_trees = 100
13 max_features = 3
14 kfold = KFold(n_splits=10, shuffle=True, random_state=7)
15 model = RandomForestClassifier(n_estimators=num_trees, max_features=max_features)
16 results = cross_val_score(model, X, Y, cv=kfold)
17 print(results.mean())|
```

### Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\12_ensembles.py
0.7695488721804511
```

## MC-13 -- Finalize and Save

Caption: Hold out 33% for testing, train Logistic Regression on the 67% training split, pickle the fitted model to disk, reload it later, and report test accuracy—your “finalized” model workflow.

### Code

```
itec6720_final > scripts > mini_course > 13_finalize_save.py > ...
1  # Save Model Using Pickle
2  from pandas import read_csv
3  from sklearn.model_selection import train_test_split
4  from sklearn.linear_model import LogisticRegression
5  import pickle
6  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv"
7  names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
8  dataframe = read_csv(url, names=names)
9  array = dataframe.values
10 X = array[:,0:8]
11 Y = array[:,8]
12 test_size = 0.33
13 seed = 7
14 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=test_size, random_state=seed)
15 # Fit the model on 67%
16 model = LogisticRegression(solver='liblinear')
17 model.fit(X_train, Y_train)
18 # save the model to disk
19 filename = 'finalized_model.sav'
20 pickle.dump(model, open(filename, 'wb'))
21
22 # some time later...
23
24 # load the model from disk
25 loaded_model = pickle.load(open(filename, 'rb'))
26 result = loaded_model.score(X_test, Y_test)
27 print(result)
```

### Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\13_finalize_save.py
0.7559055118110236
```

## MC-14 – Hello World End-to-End Project

### 1.Understand data (EDA)

**Caption: “L14 – EDA code: load iris, print shape/classes/summary.”**

**Code**

```
itec6720_final > scripts > mini_course > 14_hello_world_project.py > ...
4  from sklearn.datasets import load_iris
5  from sklearn.model_selection import train_test_split, StratifiedKFold, cross_val_score, GridSearchCV
6  from sklearn.preprocessing import StandardScaler
7  from sklearn.pipeline import Pipeline
8  from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
9  from sklearn.linear_model import LogisticRegression
10 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
11 from sklearn.neighbors import KNeighborsClassifier
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.naive_bayes import GaussianNB
14 from sklearn.svm import SVC
15 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
16 from joblib import dump, load
17
18 # -----
19 # 1) Load & EDA (minimal)
20 # -----
21 iris = load_iris(as_frame=True)
22 df = iris.frame # features + target
23 df['target'] = iris.target
24 print("Shape:", df.shape)
25 print("Classes:", df['target'].value_counts().sort_index().to_dict())
26 print(df.describe().T[['mean', 'std', 'min', 'max']])
27
```

**Output**

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\mini_course\14_hello_world_project.py
Shape: (150, 5)
Classes: {0: 50, 1: 50, 2: 50}
      mean      std   min   max
sepal length (cm)  5.843333  0.828066  4.3  7.9
sepal width (cm)  3.057333  0.435866  2.0  4.4
petal length (cm) 3.758000  1.765298  1.0  6.9
petal width (cm)  1.199333  0.762238  0.1  2.5
target          1.000000  0.819232  0.0  2.0
```

### 2.Preprocess / Train–test split

**Caption: “L14 – Split code: stratified 75/25 train/test.”**

```

# -----
# 2) Train/Test split
# -----
X = iris.data
y = iris.target
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.25, stratify=y, random_state=7
)

print("Train/Test shapes:", X_train.shape, X_test.shape)
print("Train class counts:", np.bincount(y_train))
print("Test class counts:", np.bincount(y_test))

```

### Output L14 – Split output: shapes and class balance

```

Train/Test shapes: (112, 4) (38, 4)
Train class counts: [37 37 38]
Test class counts: [13 13 12]

```

## 3.Spot-check algorithms (10-fold CV)

Caption: “L14 – Spot-check code: LR, LDA, KNN, CART, NB, SVM with 10-fold CV.”

```

# -----
# 3) Spot-check algorithms
#     (scale where it matters via Pipeline)
# -----
candidates = [
    ("LR", Pipeline([
        ('scaler', StandardScaler()),
        ('clf', LogisticRegression(solver='lbfgs', max_iter=2000))])),  # <- no multi_class

    ("LDA", LDA := LinearDiscriminantAnalysis()),
    ("KNN", Pipeline([
        ('scaler', StandardScaler()),
        ('clf', KNeighborsClassifier())])), 
    ("CART", DecisionTreeClassifier(random_state=7)),
    ("NB", GaussianNB()),
    ("SVM", Pipeline([
        ('scaler', StandardScaler()),
        ('clf', SVC(kernel='rbf'))]))
]

cv = StratifiedKFold(n_splits=10, shuffle=True, random_state=7)
print("\nSpot-check (10-fold CV accuracy):")
scores_table = []
for name, model in candidates:
    scores = cross_val_score(model, X_train, y_train, cv=cv, scoring='accuracy')
    scores_table.append((name, scores.mean(), scores.std()))
    print(f"{name:>4}: {scores.mean():.3f} (+/- {scores.std():.3f})")

```

## Output: the CV accuracy table

Spot-check (10-fold CV accuracy):

LR: 0.948 (+/- 0.057)
LDA: 0.973 (+/- 0.041)
KNN: 0.957 (+/- 0.057)
CART: 0.948 (+/- 0.057)
NB: 0.939 (+/- 0.068)
SVM: 0.948 (+/- 0.057)

## 4. Tune best model (GridSearchCV)

**Caption: “L14 – Tuning code: SVM GridSearchCV (kernel, C, gamma).”**

```
# -----
# 4) Tune the top model(s)
#     Example: SVM (often strong on Iris)
# -----
svm_pipe = Pipeline([('scaler', StandardScaler()),
| | | | | | ('clf', SVC())])

param_grid = {
    'clf_kernel': ['rbf', 'linear'],
    'clf_C': [0.1, 1, 10, 100],
    'clf_gamma': ['scale', 0.1, 0.01, 0.001] # gamma used when rbf
}
grid = GridSearchCV(svm_pipe, param_grid, scoring='accuracy', cv=cv, n_jobs=-1, return_train_score=True)
grid.fit(X_train, y_train)
print("\nBest SVM params:", grid.best_params_)
print("Best CV accuracy:", grid.best_score_)
```

**Output: “L14 – Tuning results: best params & CV score.”**

```
Best SVM params: {'clf_C': 1, 'clf_gamma': 'scale', 'clf_kernel': 'linear'}  
Best CV accuracy: 0.9818181818181818
```

## 5. Try an ensemble

**Caption: “L14 – Ensemble code: RF & GB CV evaluation.”**

```
# -----
# 5) Try an ensemble (comparison)
# -----
rf = RandomForestClassifier(n_estimators=200, random_state=7)
gb = GradientBoostingClassifier(random_state=7)
for name, model in [("RF", rf), ("GB", gb)]:
    scores = cross_val_score(model, X_train, y_train, cv=cv, scoring='accuracy')
    print(f"{name} CV accuracy: {scores.mean():.3f} (+/- {scores.std():.3f})")
```

## OutPut

RF CV accuracy: 0.965 (+/- 0.043)  
GB CV accuracy: 0.965 (+/- 0.043)

```

# -----
# 6) Finalize, evaluate on test, save/load
#   (Pick the best: grid.best_estimator_ here)
# -----
best_model = grid.best_estimator_
best_model.fit(X_train, y_train)
y_pred = best_model.predict(X_test)

print("\nTEST accuracy:", accuracy_score(y_test, y_pred))
print("Confusion matrix:\n", confusion_matrix(y_test, y_pred))
print("Classification report:\n", classification_report(y_test, y_pred, target_names=iris.target_names))

```

## 6.Finalize, evaluate, save & reload

**Caption: “L14 – Finalize code: fit best model, test metrics, save & reload.”**

```

# -----
# 6) Finalize, evaluate on test, save/load
#   (Pick the best: grid.best_estimator_ here)
# -----
best_model = grid.best_estimator_
best_model.fit(X_train, y_train)
y_pred = best_model.predict(X_test)

print("\nTEST accuracy:", accuracy_score(y_test, y_pred))
print("Confusion matrix:\n", confusion_matrix(y_test, y_pred))
print("Classification report:\n", classification_report(y_test, y_pred, target_names=iris.target_names))

dump(best_model, "iris_best_model.joblib")
loaded = load("iris_best_model.joblib")
print("Loaded model test accuracy (sanity check):",
     accuracy_score(y_test, loaded.predict(X_test)))

```

**Output: test accuracy, confusion matrix, classification report, and “Loaded model test accuracy”.**

```

TEST accuracy: 0.9736842105263158
Confusion matrix:
[[13  0  0]
 [ 0 12  1]
 [ 0  0 12]]
Classification report:
             precision    recall  f1-score   support
setosa       1.00     1.00     1.00      13
versicolor   1.00     0.92     0.96      13
virginica    0.92     1.00     0.96      12
accuracy          0.97     0.97     0.97      38
macro avg       0.97     0.97     0.97      38
weighted avg    0.98     0.97     0.97      38

```

Loaded model test accuracy (sanity check): 0.9736842105263158

## 2) Iris Step-by-Step Project

Step 1–6: code + output/plot per step.

### Step\_1. Install / versions

```
itec6720_final > scripts > step_by_step > 01_env_versions.py
1  # Check the versions of libraries
2
3  # Python version
4  import sys
5  print('Python: {}'.format(sys.version))
6  # scipy
7  import scipy
8  print('scipy: {}'.format(scipy.__version__))
9  # numpy
10 import numpy
11 print('numpy: {}'.format(numpy.__version__))
12 # matplotlib
13 import matplotlib
14 print('matplotlib: {}'.format(matplotlib.__version__))
15 # pandas
16 import pandas
17 print('pandas: {}'.format(pandas.__version__))
18 # scikit-learn
19 import sklearn
20 print('sklearn: {}'.format(sklearn.__version__))
```

### Output

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\step_by_step\01_env_versions.py
Python: 3.11.13 | packaged by Anaconda, Inc. | (main, Jun 5 2025, 13:03:15) [MSC v.1929 64 bit (AMD64)]
scipy: 1.16.1
numpy: 2.3.3
matplotlib: 3.10.1
pandas: 2.3.2
sklearn: 1.7.2
```

## Step\_2. Load the dataset

```
itec6720_final > scripts > step_by_step > 02_load_csv.py > ...
1  # Load libraries
2  from pandas import read_csv
3  from pandas.plotting import scatter_matrix
4  from matplotlib import pyplot as plt
5  from sklearn.model_selection import train_test_split
6  from sklearn.model_selection import cross_val_score
7  from sklearn.model_selection import StratifiedKFold
8  from sklearn.metrics import classification_report
9  from sklearn.metrics import confusion_matrix
10 from sklearn.metrics import accuracy_score
11 from sklearn.linear_model import LogisticRegression
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.neighbors import KNeighborsClassifier
14 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
15 from sklearn.naive_bayes import GaussianNB
16 from sklearn.svm import SVC
17
18
19 # Load dataset
20 url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/iris.csv"
21 names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
22 dataset = read_csv(url, names=names)
23
24
```

## Step\_3. Summarize the dataset

```
itec6720_final > scripts > step_by_step > 03_descriptive_stats.py > ...
1
2  # summarize the data
3  from pandas import read_csv
4  # Load dataset
5  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/iris.csv"
6  names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
7  dataset = read_csv(url, names=names)
8  # shape
9  print(dataset.shape)
10 # head
11 print(dataset.head(20))
12 # descriptions
13 print(dataset.describe())
14 # class distribution
15 print(dataset.groupby('class').size())
```

## Output

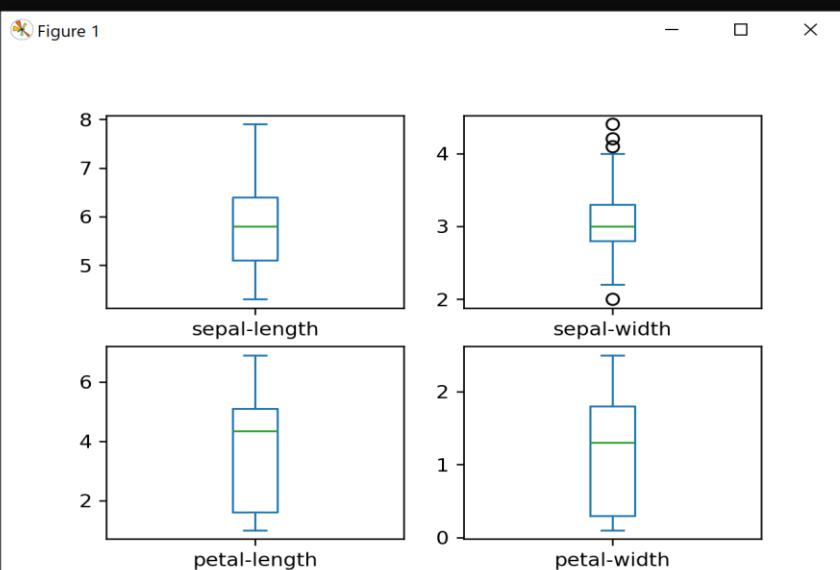
```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\step_by_step\03_descriptive_stats.py
(150, 5)
   sepal-length  sepal-width  petal-length  petal-width      class
0            5.1         3.5         1.4         0.2  Iris-setosa
1            4.9         3.0         1.4         0.2  Iris-setosa
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
5            5.4         3.9         1.7         0.4  Iris-setosa
6            4.6         3.4         1.4         0.3  Iris-setosa
7            5.0         3.4         1.5         0.2  Iris-setosa
8            4.4         2.9         1.4         0.2  Iris-setosa
9            4.9         3.1         1.5         0.1  Iris-setosa
10           5.4         3.7         1.5         0.2  Iris-setosa
11           4.8         3.4         1.6         0.2  Iris-setosa
12           4.8         3.0         1.4         0.1  Iris-setosa
13           4.3         3.0         1.1         0.1  Iris-setosa
14           5.8         4.0         1.2         0.2  Iris-setosa
15           5.7         4.4         1.5         0.4  Iris-setosa
16           5.4         3.9         1.3         0.4  Iris-setosa
17           5.1         3.5         1.4         0.3  Iris-setosa
18           5.7         3.8         1.7         0.3  Iris-setosa
19           5.1         3.8         1.5         0.3  Iris-setosa
   sepal-length  sepal-width  petal-length  petal-width
count    150.000000  150.000000  150.000000  150.000000
mean     5.843333    3.054000    3.758667    1.198667
std      0.828066    0.433594    1.764420    0.763161
min      4.300000    2.000000    1.000000    0.100000
25%     5.100000    2.800000    1.600000    0.300000
50%     5.800000    3.000000    4.350000    1.300000
75%     6.400000    3.300000    5.100000    1.800000
max     7.900000    4.400000    6.900000    2.500000
class
Iris-setosa      50
Iris-versicolor  50
Iris-virginica   50
dtype: int64
```

## Step\_4. Visualize the data

```
itec6720_final > scripts > step_by_step > 04_visualization.py > ...
1
2  # visualize the data
3  from pandas import read_csv
4  from pandas.plotting import scatter_matrix
5  from matplotlib import pyplot as plt
6  # Load dataset
7  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/iris.csv"
8  names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
9  dataset = read_csv(url, names=names)
10 # box and whisker plots
11 dataset.plot(kind='box', subplots=True, layout=(2,2), sharex=False, sharey=False)
12 plt.show()
13 # histograms
14 dataset.hist()
15 plt.show()
16 # scatter plot matrix
17 scatter_matrix(dataset)
18 plt.show()
```

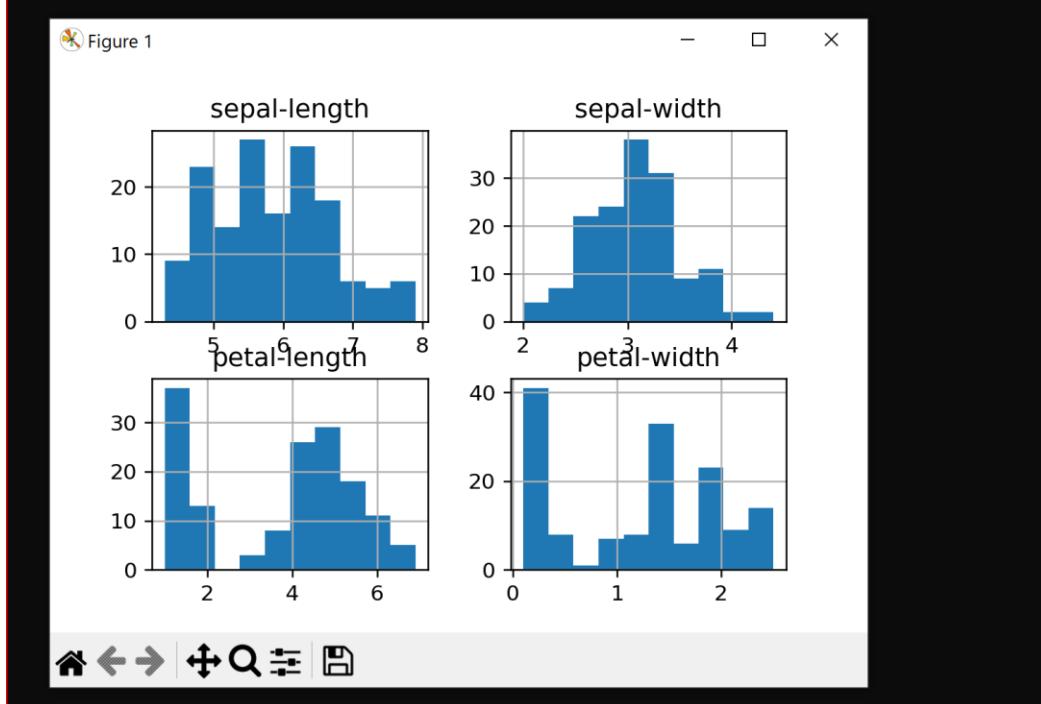
**Boxplots (figures)** — Box-and-whisker plots for the four Iris features (cm).

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\step_by_step\04_visualization.py
```



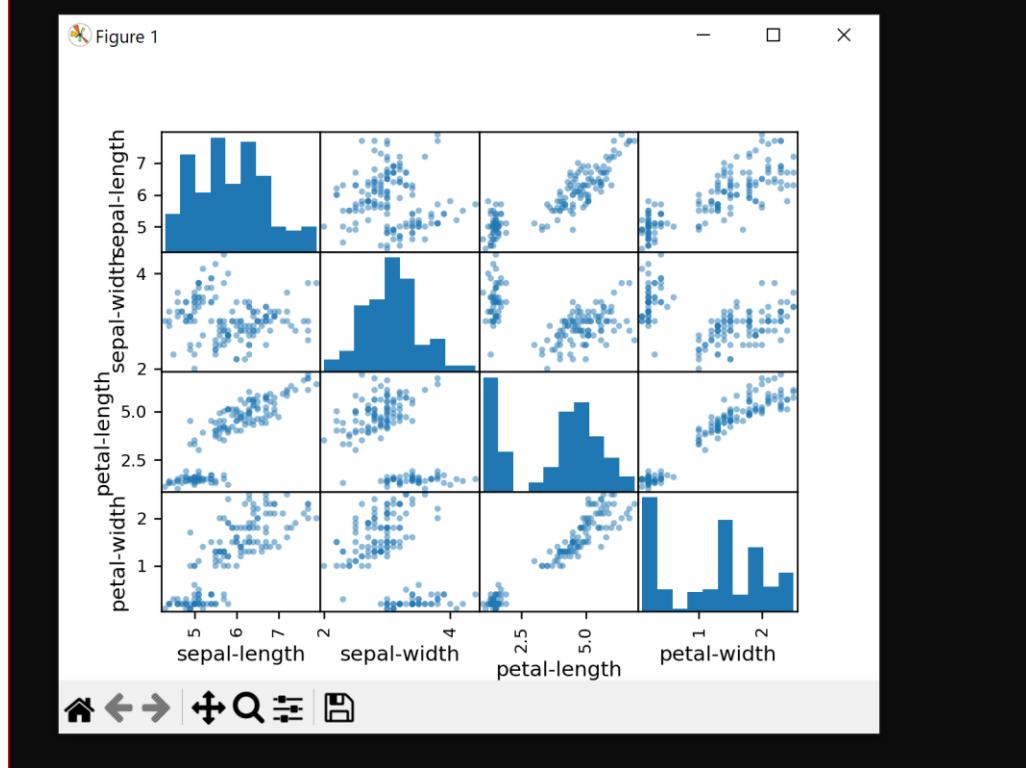
## Histograms (figures) — Univariate distributions.

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\step_by_step\04_visualization.py
```



## Scatter-matrix (figures) — Pairwise relationships.

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\step_by_step\04_visualization.py
```



## Step\_5.Evaluate some algorithms (Spot-check).

**Step-by-Step — 05 Evaluate Algorithms:** 10-fold stratified CV accuracy for LR, LDA, KNN, CART, NB, SVM (with pipelines to scale KNN/SVM); prints mean  $\pm$  std and compares models with a boxplot.

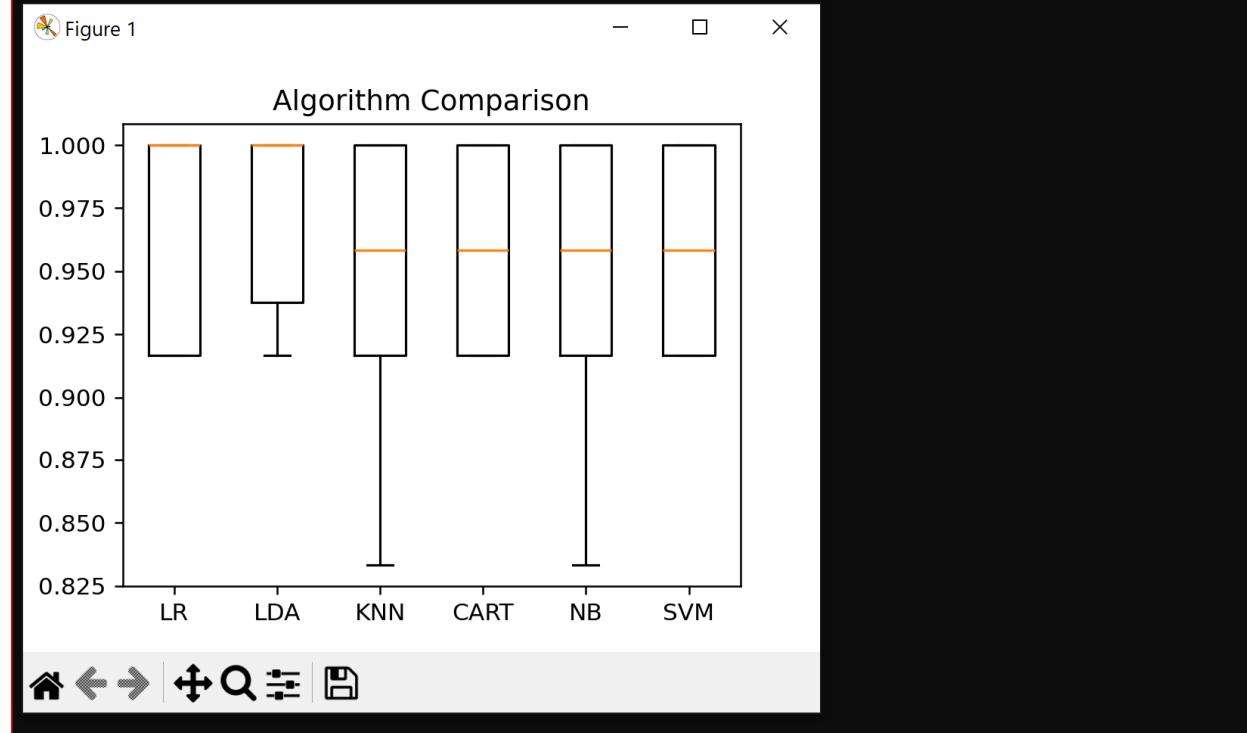
```

scripts > step_by_step > 05_evaluate_algorithms.py > ...
 1  from pandas import read_csv
 2  from matplotlib import pyplot as plt
 3  from sklearn.model_selection import train_test_split, StratifiedKFold, cross_val_score
 4  from sklearn.pipeline import make_pipeline
 5  from sklearn.preprocessing import StandardScaler
 6  from sklearn.linear_model import LogisticRegression
 7  from sklearn.tree import DecisionTreeClassifier
 8  from sklearn.neighbors import KNeighborsClassifier
 9  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
10  from sklearn.naive_bayes import GaussianNB
11  from sklearn.svm import SVC
12
13  # Load dataset
14  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/iris.csv"
15  names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
16  dataset = read_csv(url, names=names)
17
18  # Split-out validation dataset
19  array = dataset.values
20  X, y = array[:,0:4], array[:,4]
21  X_train, X_validation, Y_train, Y_validation = train_test_split(
22      X, y, test_size=0.20, random_state=1, shuffle=True
23  )
24
25  # Spot Check Algorithms (updated)
26  models = [
27      ('LR', make_pipeline(StandardScaler(), LogisticRegression(max_iter=1000))), # lbfgs by default; no deprecation
28      ('LDA', LinearDiscriminantAnalysis()),
29      ('KNN', make_pipeline(StandardScaler(), KNeighborsClassifier())),
30      ('CART', DecisionTreeClassifier(random_state=1)),
31      ('NB', GaussianNB()),
32      ('SVM', make_pipeline(StandardScaler(), SVC())) # default gamma='scale'
33  ]
34
35  # evaluate each model in turn
36  results, names = [], []
37  cv = StratifiedKFold(n_splits=10, random_state=1, shuffle=True)
38  for name, model in models:
39      cv_results = cross_val_score(model, X_train, Y_train, cv=cv, scoring='accuracy')
40      results.append(cv_results); names.append(name)
41      print(f'{name}: {cv_results.mean():.3f} ({cv_results.std():.3f})')
42
43  plt.boxplot(results, tick_labels=names)
44  plt.title('Algorithm Comparison')

```

Mean  $\pm$  std accuracy for LR, LDA, KNN, CART, NB, SVM printed at top; boxplot shows SVM/LDA near the top with tight spread  $\Rightarrow$  strong, consistent performance.

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\step_by_step\05_evaluate_algorithms.py
LR: 0.967 (0.041)
LDA: 0.975 (0.038)
KNN: 0.950 (0.055)
CART: 0.958 (0.042)
NB: 0.950 (0.055)
SVM: 0.958 (0.042)
```



## Step\_6. Make predictions

Train an SVC on the 80% training split, predict the 20% validation set, then print accuracy, the confusion matrix, and a full classification report.

```
itec6720_final > scripts > step_by_step > 06_predictions.py > ...
1  # make predictions
2  from pandas import read_csv
3  from sklearn.model_selection import train_test_split
4  from sklearn.metrics import classification_report
5  from sklearn.metrics import confusion_matrix
6  from sklearn.metrics import accuracy_score
7  from sklearn.svm import SVC
8  # Load dataset
9  url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/iris.csv"
10 names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
11 dataset = read_csv(url, names=names)
12 # Split-out validation dataset
13 array = dataset.values
14 X = array[:,0:4]
15 y = array[:,4]
16 X_train, X_validation, Y_train, Y_validation = train_test_split(X, y, test_size=0.20, random_state=1)
17 # Make predictions on validation dataset
18 model = SVC(gamma='auto')
19 model.fit(X_train, Y_train)
20 predictions = model.predict(X_validation)
21 # Evaluate predictions
22 print(accuracy_score(Y_validation, predictions))
23 print(confusion_matrix(Y_validation, predictions))
24 print(classification_report(Y_validation, predictions))
```

Accuracy  $\approx 0.967$  with 29/30 correct (one *versicolor* predicted as *virginica*). Macro F1  $\approx 0.96$ ; per-class scores are strong across all three classes.

```
(itec6720) D:\MGA\Projects\itec6720_final>python scripts\step_by_step\06_predictions.py
0.9666666666666667
[[11  0  0]
 [ 0 12  1]
 [ 0  0  6]]
      precision    recall  f1-score   support
Iris-setosa      1.00      1.00      1.00       11
Iris-versicolor   1.00      0.92      0.96       13
Iris-virginica    0.86      1.00      0.92        6

      accuracy         0.97       30
     macro avg       0.95      0.97      0.96       30
  weighted avg     0.97      0.97      0.97       30
```

