

This is the solution for BCIS 5140-HW2 Part1 by Srivyshnavi Tangellapelli,Fall 2022.

Setup

```
# Python ≥3.5 is required
import sys
assert sys.version_info >= (3, 5)

# Scikit-Learn ≥0.20 is required
import sklearn
assert sklearn.__version__ >= "0.20"

# Common imports
import numpy as np
import os

# to make this notebook's output stable across runs
np.random.seed(42)

# To plot pretty figures
%matplotlib inline
import matplotlib as mpl
import matplotlib.pyplot as plt
mpl.rc('axes', labelsize=14)
mpl.rc('xtick', labelsize=12)
mpl.rc('ytick', labelsize=12)

# Where to save the figures
PROJECT_ROOT_DIR = "."
CHAPTER_ID = "training_linear_models"
IMAGES_PATH = os.path.join(PROJECT_ROOT_DIR, "images", CHAPTER_ID)
os.makedirs(IMAGES_PATH, exist_ok=True)

def save_fig(fig_id, tight_layout=True, fig_extension="png", resolution=300):
    path = os.path.join(IMAGES_PATH, fig_id + "." + fig_extension)
    print("Saving figure", fig_id)
    if tight_layout:
        plt.tight_layout()
    plt.savefig(path, format=fig_extension, dpi=resolution)
```

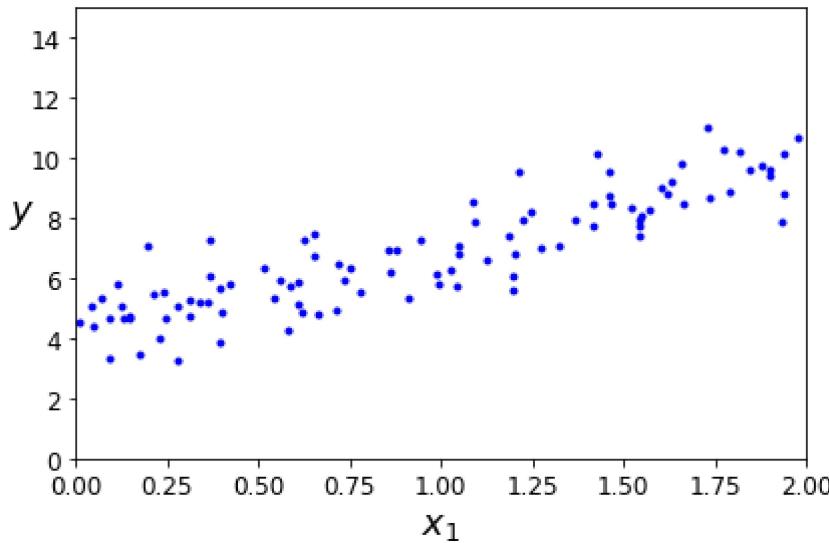
Linear Regression (The Normal Equation)

```
import numpy as np
```

```
X = 2 * np.random.rand(100, 1)
y = 4 + 3 * X + np.random.randn(100, 1)
```

```
plt.plot(X, y, "b.")
plt.xlabel("$x_1$", fontsize=18)
plt.ylabel("$y$", rotation=0, fontsize=18)
plt.axis([0, 2, 0, 15])
save_fig("generated_data_plot")
plt.show()
```

Saving figure generated_data_plot



```
X_b = np.c_[np.ones((100, 1)), X] # add x0 = 1 to each instance
theta_best = np.linalg.inv(X_b.T.dot(X_b)).dot(X_b.T).dot(y)
```

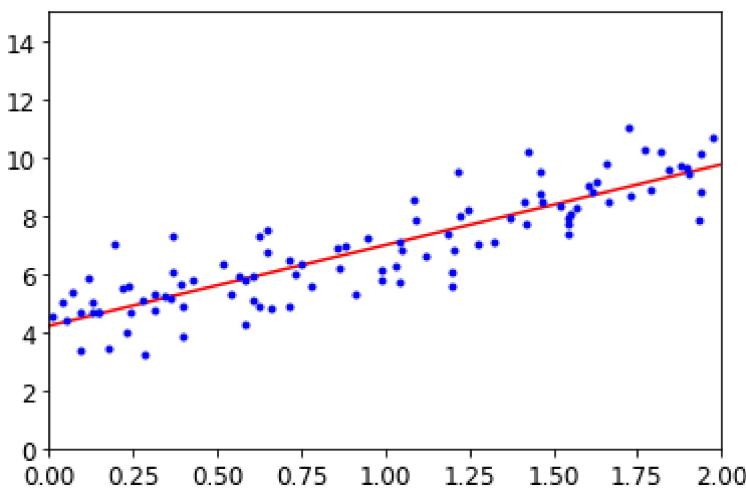
theta_best

```
array([[4.21509616],
       [2.77011339]])
```

```
X_new = np.array([[0], [2]])
X_new_b = np.c_[np.ones((2, 1)), X_new] # add x0 = 1 to each instance
y_predict = X_new_b.dot(theta_best)
y_predict
```

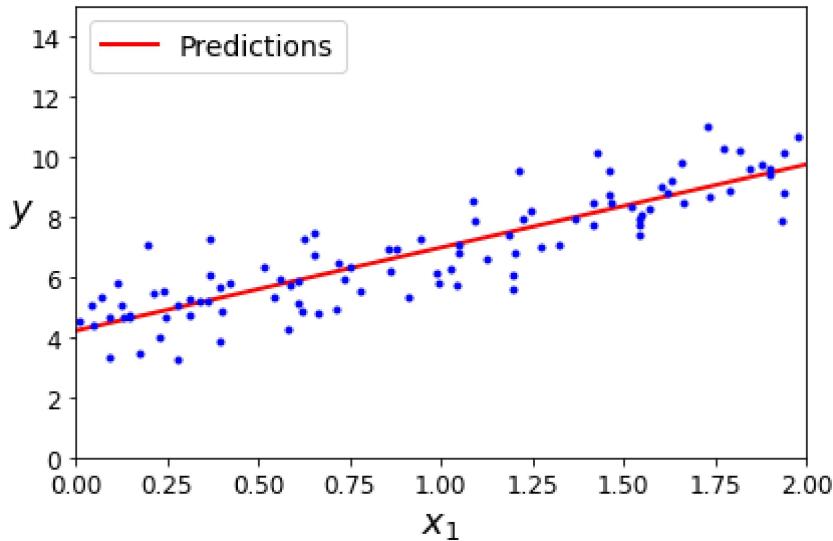
```
array([[4.21509616],
       [9.75532293]])
```

```
plt.plot(X_new, y_predict, "r-")
plt.plot(X, y, "b.")
plt.axis([0, 2, 0, 15])
plt.show()
```



```
plt.plot(X_new, y_predict, "r-", linewidth=2, label="Predictions")
plt.plot(X, y, "b.")
plt.xlabel("$x_1$", fontsize=18)
plt.ylabel("$y$", rotation=0, fontsize=18)
plt.legend(loc="upper left", fontsize=14)
plt.axis([0, 2, 0, 15])
save_fig("linear_model_predictions_plot")
plt.show()
```

Saving figure linear_model_predictions_plot



```
from sklearn.linear_model import LinearRegression

lin_reg = LinearRegression()
lin_reg.fit(X, y)
lin_reg.intercept_, lin_reg.coef_

(array([4.21509616]), array([[2.77011339]]))

lin_reg.predict(X_new)
```

```
array([[4.21509616],
       [9.75532293]])
```

```
theta_best_svd, residuals, rank, s = np.linalg.lstsq(X_b, y, rcond=1e-6)
theta_best_svd
```

```
array([[4.21509616],
       [2.77011339]])
```

```
np.linalg.pinv(X_b).dot(y)
```

```
array([[4.21509616],
       [2.77011339]])
```

Gradient Descent (Batch Gradient Descent)

```
eta = 0.1 # learning rate #deciding the size of the steps to take for better convergence
n_iterations = 1000
m = 100

theta = np.random.randn(2,1) # random initialization # randomly initializing the theta before

for iteration in range(n_iterations):
    gradients = 2/m * X_b.T.dot(X_b.dot(theta) - y) # computing the gradient vector of the cost function
    theta = theta - eta * gradients # deciding the gradient descent step or your next step to take

theta
array([[4.21509616],
       [2.77011339]])

X_new_b.dot(theta)
array([[4.21509616],
       [9.75532293]])

theta_path_bgd = []

def plot_gradient_descent(theta, eta, theta_path=None):
    m = len(X_b)
    plt.plot(X, y, "b.")
    n_iterations = 1000
    for iteration in range(n_iterations):
        if iteration < 10:
            y_predict = X_new_b.dot(theta)
            style = "b-" if iteration > 0 else "r--"
```

```

        plt.plot(X_new, y_predict, style)
gradients = 2/m * X_b.T.dot(X_b.dot(theta) - y)
theta = theta - eta * gradients
if theta_path is not None:
    theta_path.append(theta)
plt.xlabel("$x_1$", fontsize=18)
plt.axis([0, 2, 0, 15])
plt.title(r"$\eta = {}$".format(eta), fontsize=16)

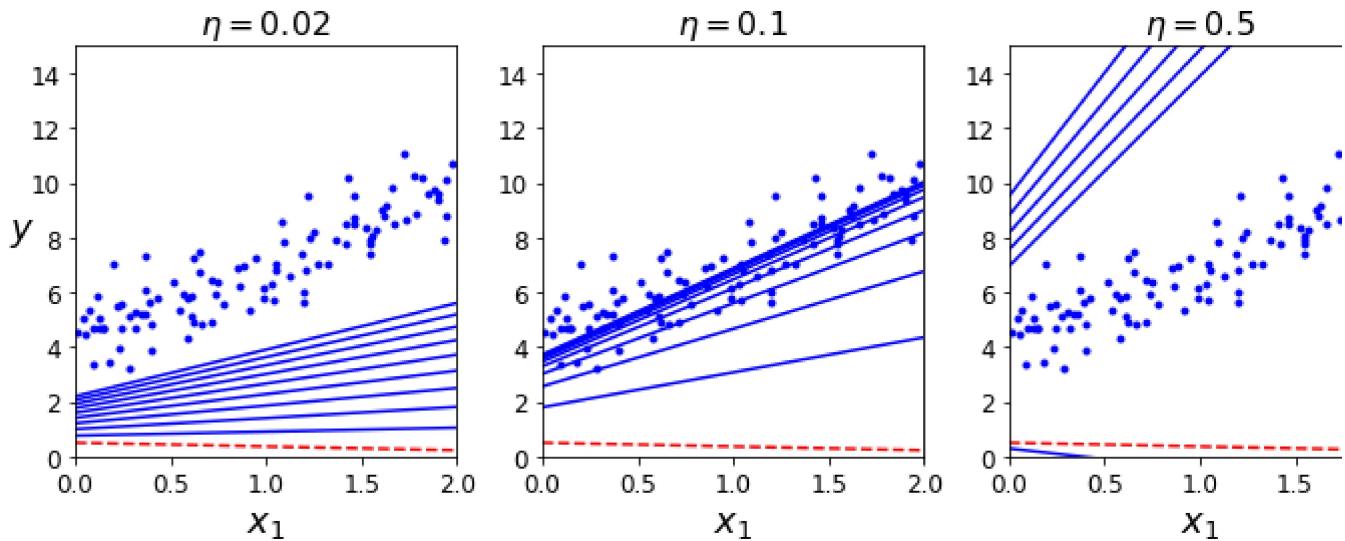
np.random.seed(42)
theta = np.random.randn(2,1) # random initialization

plt.figure(figsize=(10,4))
plt.subplot(131); plot_gradient_descent(theta, eta=0.02)
plt.ylabel("$y$", rotation=0, fontsize=18)
plt.subplot(132); plot_gradient_descent(theta, eta=0.1, theta_path=theta_path_bgd)
plt.subplot(133); plot_gradient_descent(theta, eta=0.5)

save_fig("gradient_descent_plot")
plt.show()

```

Saving figure gradient_descent_plot



Stochastic Gradient Descent

```

theta_path_sgd = []
m = len(X_b)
np.random.seed(42)

n_epochs = 50
t0, t1 = 5, 50 # learning schedule hyperparameters

def learning_schedule(t):
    return t0 / (t + t1)

```

```

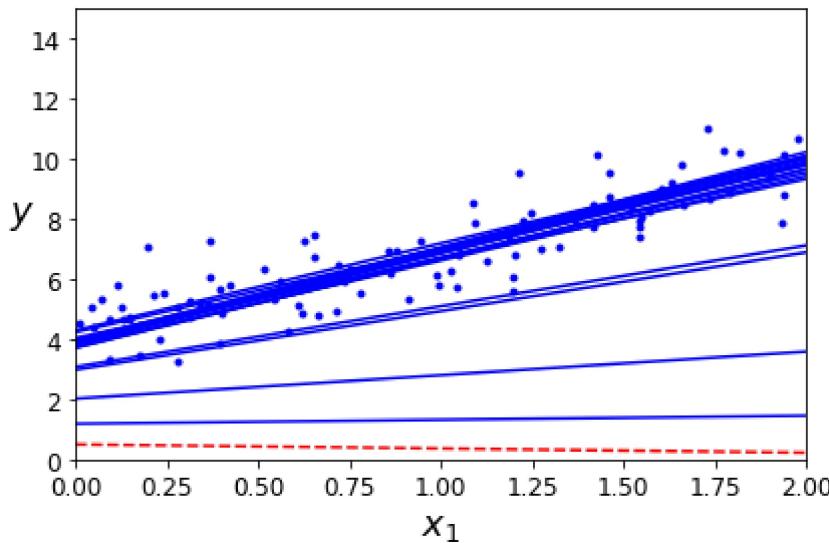
theta = np.random.randn(2,1) # random initialization

for epoch in range(n_epochs):
    for i in range(m):
        if epoch == 0 and i < 20:                      # not shown in the book
            y_predict = X_new_b.dot(theta)               # not shown
            style = "b-" if i > 0 else "r--"           # not shown
            plt.plot(X_new, y_predict, style)            # not shown
        random_index = np.random.randint(m)
        xi = X_b[random_index:random_index+1]
        yi = y[random_index:random_index+1]
        gradients = 2 * xi.T.dot(xi.dot(theta) - yi)
        eta = learning_schedule(epoch * m + i)
        theta = theta - eta * gradients
        theta_path_sgd.append(theta)                  # not shown

plt.plot(X, y, "b.")                                     # not shown
plt.xlabel("$x_1$", fontsize=18)                         # not shown
plt.ylabel("$y$", rotation=0, fontsize=18)                # not shown
plt.axis([0, 2, 0, 15])                                  # not shown
save_fig("sgd_plot")                                    # not shown
plt.show()                                              # not shown

```

Saving figure sgd_plot



theta

```
array([[4.21076011],
       [2.74856079]])
```

```
from sklearn.linear_model import SGDRegressor
```

```
sgd_reg = SGDRegressor(max_iter=1000, tol=1e-3, penalty=None, eta0=0.1, random_state=42)
sgd_reg.fit(X, y.ravel())
```

```
SGDRegressor(eta0=0.1, penalty=None, random_state=42)
```

```
sgd_reg.intercept_, sgd_reg.coef_
```

```
(array([4.24365286]), array([2.8250878]))
```

Mini-batch gradient descent

```
theta_path_mgd = []
```

```
n_iterations = 50
```

```
minibatch_size = 20
```

```
np.random.seed(42)
```

```
theta = np.random.randn(2,1) # random initialization
```

```
t0, t1 = 200, 1000
```

```
def learning_schedule(t):
```

```
    return t0 / (t + t1)
```

```
t = 0
```

```
for epoch in range(n_iterations):
```

```
    shuffled_indices = np.random.permutation(m)
```

```
    X_b_shuffled = X_b[shuffled_indices]
```

```
    y_shuffled = y[shuffled_indices]
```

```
    for i in range(0, m, minibatch_size):
```

```
        t += 1
```

```
        xi = X_b_shuffled[i:i+minibatch_size]
```

```
        yi = y_shuffled[i:i+minibatch_size]
```

```
        gradients = 2/minibatch_size * xi.T.dot(xi.dot(theta) - yi)
```

```
        eta = learning_schedule(t)
```

```
        theta = theta - eta * gradients
```

```
        theta_path_mgd.append(theta)
```

```
theta
```

```
array([[4.25214635],  
       [2.7896408 ]])
```

```
theta_path_bgd = np.array(theta_path_bgd)
```

```
theta_path_sgd = np.array(theta_path_sgd)
```

```
theta_path_mgd = np.array(theta_path_mgd)
```

```
plt.figure(figsize=(7,4))
```

```
plt.plot(theta_path_sgd[:, 0], theta_path_sgd[:, 1], "r-s", linewidth=1, label="Stochastic")
```

```
plt.plot(theta_path_mgd[:, 0], theta_path_mgd[:, 1], "g+-", linewidth=2, label="Mini-batch")
```

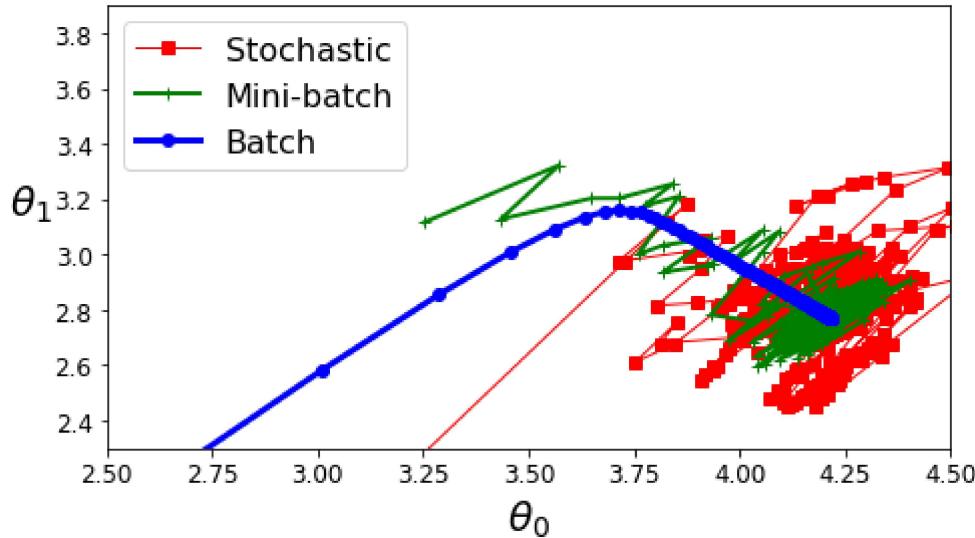
```
plt.plot(theta_path_bgd[:, 0], theta_path_bgd[:, 1], "b-o", linewidth=3, label="Batch")
```

```

plt.legend(loc="upper left", fontsize=16)
plt.xlabel(r"$\theta_0$", fontsize=20)
plt.ylabel(r"$\theta_1$ ", fontsize=20, rotation=0)
plt.axis([2.5, 4.5, 2.3, 3.9])
save_fig("gradient_descent_paths_plot")
plt.show()

```

Saving figure gradient_descent_paths_plot



Logistic Regression (Decision Boundaries)

```

t = np.linspace(-10, 10, 100)
sig = 1 / (1 + np.exp(-t))
plt.figure(figsize=(9, 3))
plt.plot([-10, 10], [0, 0], "k-")
plt.plot([-10, 10], [0.5, 0.5], "k:")
plt.plot([-10, 10], [1, 1], "k:")
plt.plot([0, 0], [-1.1, 1.1], "k-")
plt.plot(t, sig, "b-", linewidth=2, label=r"$\sigma(t) = \frac{1}{1 + e^{-t}}$")
plt.xlabel("t")
plt.legend(loc="upper left", fontsize=20)
plt.axis([-10, 10, -0.1, 1.1])
save_fig("logistic_function_plot")
plt.show()

```

Saving figure logistic_function_plot



```
from sklearn import datasets
iris = datasets.load_iris()
list(iris.keys())
```

```
['data',
 'target',
 'frame',
 'target_names',
 'DESCR',
 'feature_names',
 'filename',
 'data_module']
```

```
print(iris.DESCR)
```

```
.. _iris_dataset:
```

Iris plants dataset

Data Set Characteristics:

:Number of Instances: 150 (50 in each of three classes)
:Number of Attributes: 4 numeric, predictive attributes and the class
:Attribute Information:
 - sepal length in cm
 - sepal width in cm
 - petal length in cm
 - petal width in cm
 - class:
 - Iris-Setosa
 - Iris-Versicolour
 - Iris-Virginica

:Summary Statistics:

| | Min | Max | Mean | SD | Class Correlation |
|---------------|-----|-----|------|------|-------------------|
| sepal length: | 4.3 | 7.9 | 5.84 | 0.83 | 0.7826 |
| sepal width: | 2.0 | 4.4 | 3.05 | 0.43 | -0.4194 |
| petal length: | 1.0 | 6.9 | 3.76 | 1.76 | 0.9490 (high!) |
| petal width: | 0.1 | 2.5 | 1.20 | 0.76 | 0.9565 (high!) |

:Missing Attribute Values: None

:Class Distribution: 33.3% for each of 3 classes.

:Creator: R.A. Fisher

:Donor: Michael Marshall (MARSHALL%PLU@io.arc.nasa.gov)

:Date: July, 1988

The famous Iris database, first used by Sir R.A. Fisher. The dataset is taken from Fisher's paper. Note that it's the same as in R, but not as in the UCI Machine Learning Repository, which has two wrong data points.

This is perhaps the best known database to be found in the pattern recognition literature. Fisher's paper is a classic in the field and is referenced frequently to this day. (See Duda & Hart, for example.) The data set contains 3 classes of 50 instances each, where each class refers to a type of iris plant. One class is linearly separable from the other 2; the latter are NOT linearly separable from each other.

.. topic:: References

- Fisher, R.A. "The use of multiple measurements in taxonomic problems" Annual Eugenics, 7, Part II, 179-188 (1936); also in "Contributions to Mathematical Statistics" (John Wiley, NY, 1950).
- Duda, R.O., & Hart, P.E. (1973) Pattern Classification and Scene Analysis. (Q327.D83) John Wiley & Sons. ISBN 0-471-22361-1. See page 218.
- Dasarathy, B.V. (1980) "Nosing Around the Neighborhood: A New System Structure and Classification Rule for Recognition in Partially Exposed Environments". IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. PAMI-2, No. 1, 67-71.

```
X = iris["data"][:, 3:] # petal width
y = (iris["target"] == 2).astype(np.int) # 1 if Iris virginica, else 0

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:2: DeprecationWarning: `np
Deprecated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/rele
```

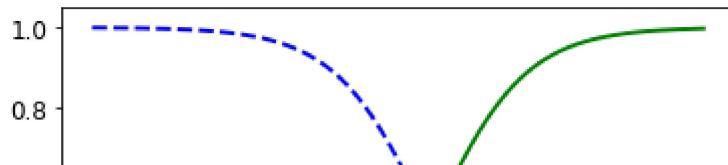
```
from sklearn.linear_model import LogisticRegression
log_reg = LogisticRegression(solver="lbfgs", random_state=42)
log_reg.fit(X, y)
```

LogisticRegression(random_state=42)

```
X_new = np.linspace(0, 3, 1000).reshape(-1, 1)
y_proba = log_reg.predict_proba(X_new)

plt.plot(X_new, y_proba[:, 1], "g-", linewidth=2, label="Iris virginica")
plt.plot(X_new, y_proba[:, 0], "b--", linewidth=2, label="Not Iris virginica")
```

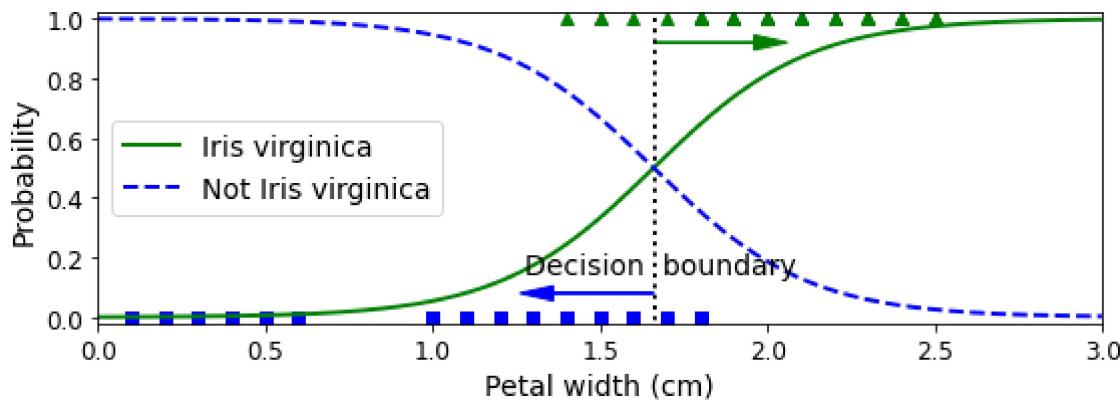
```
[<matplotlib.lines.Line2D at 0x7f8c60c9a250>]
```



```
X_new = np.linspace(0, 3, 1000).reshape(-1, 1)
y_proba = log_reg.predict_proba(X_new)
decision_boundary = X_new[y_proba[:, 1] >= 0.5][0]
```

```
plt.figure(figsize=(8, 3))
plt.plot(X[y==0], y[y==0], "bs")
plt.plot(X[y==1], y[y==1], "g^")
plt.plot([decision_boundary, decision_boundary], [-1, 2], "k:", linewidth=2)
plt.plot(X_new, y_proba[:, 1], "g-", linewidth=2, label="Iris virginica")
plt.plot(X_new, y_proba[:, 0], "b--", linewidth=2, label="Not Iris virginica")
plt.text(decision_boundary+0.02, 0.15, "Decision boundary", fontsize=14, color="k", ha="center")
plt.arrow(decision_boundary, 0.08, -0.3, 0, head_width=0.05, head_length=0.1, fc='b', ec='b')
plt.arrow(decision_boundary, 0.92, 0.3, 0, head_width=0.05, head_length=0.1, fc='g', ec='g')
plt.xlabel("Petal width (cm)", fontsize=14)
plt.ylabel("Probability", fontsize=14)
plt.legend(loc="center left", fontsize=14)
plt.axis([0, 3, -0.02, 1.02])
save_fig("logistic_regression_plot")
plt.show()
```

```
/usr/local/lib/python3.7/dist-packages/matplotlib/patches.py:1327: VisibleDeprecationWarning
    verts = np.dot(coords, M) + (x + dx, y + dy)
Saving figure logistic_regression_plot
```



```
decision_boundary
```

```
array([1.66066066])
```

```
log_reg.predict([[1.7], [1.5]])
```

```
array([1, 0])
```

▼ Question and Answers

1. Why would you add a vector of ones to the X matrix before performing linear regression using Normal equation (line 4 of Geron's Ch. 4 notebook) (1 points)

The X_b is the instance's feature vector, containing x_0 to x_n , here the value of $n = 1$ and the x_0 always equal to 1. From X matrix we have values for x_1 and hence we are adding a column vector with value just 1 to the X matrix to get the instance's feature vector.

2. Explain the code for running linear regression using LinearRegression estimator. What is the difference between `fit(X, y)` and `lin_reg.predict(X_new)` (ln. 9 and 10 of Geron's Ch. 4 notebook) (2 points)

The `.fit(X,y)` function of Scikit learn generally takes two arguments for the supervised learning and trains the model and estimates the weight parameters. The `.predict(X_new)` function of Scikit learn generally takes one argument, which is an test input data and predicts or estimates the output or target value using the trained model from the `fit(X,y)` function.

3. What are the key steps of the batch gradient descent? Add comments to the line gradient descent code (ln. 13 of Geron's Ch. 4 notebook) (2 points)

1. random initialization of the theta value
2. Deciding the learning rate, which decides the size of the steps
3. computing the gradient vector or partial derivative vector of the cost function (here it is MSE for linear regression) with respective to theta
4. deciding the gradient descent step or next step i.e., θ_{next} from subtracting learning rate multiplied with gradient descent from θ_{prev}

added the comments to the code snippet in the above

4. Why do we apply the logistic function to the linear combination of features in linear regression? (1 point)

Using the Logistic function we estimate the probability value between 0 to 1 of the linear combination of features to get the output class instead of computing the output value directly like in Linear Regression.

[Colab paid products](#) - [Cancel contracts here](#)

✓ 0s completed at 9:19 PM

