# 4 Regression-Based Classification

Certain algorithms that were first developed for regression can be adapted for classification, and some classifiers can be modified to predict continuous values. Converting a linear regressor into a classifier by adding a logistic link, for instance, retains the original coefficient vector, which makes it easy to see how each input feature influences the decision. As these dual-purpose models expose many of the hyper-parameters found in their regression versions, they often provide more opportunities for fine-tuning and clearer interpretability than algorithms designed purely for classification.

## 4.1 Logistic Regression

Logistic Regression, sometimes called logit regression, models the probability that a sample belongs to a specific class, such as the chance that an email is spam. When this probability exceeds 50%, the model predicts the positive class (label "1"); otherwise, it predicts the negative class (label "0"). It therefore acts as a binary classifier.

The predicted probability is

$$\hat{p} = h_\theta(X) = \sigma(\theta^{\mathrm{T}} \cdot X). \tag{4.1}$$

Here $\hat{p}$ is the estimated probability, and $\sigma(\cdot)$ is the sigmoid function, an $S$-shaped curve that maps any real number to the interval $(0,1)$. The logistic function is defined in Equation 4.2 and illustrated in Figure 4.1. As before, $h_\theta(X)$ denotes the hypothesis that the input matrix $X$, augmented with a bias term, belongs to the positive class under the parameters $\theta$.
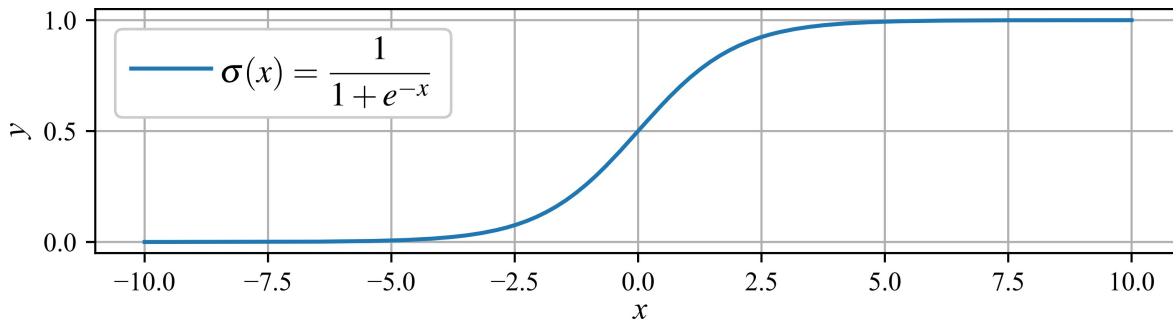
$$\sigma(x) = \frac{1}{1+e^{-x}}. \tag{4.2}$$



Figure 4.1: Sigmoid function that maps any real-valued input $x$ to a value between 0 and 1.

Once the probability $\hat{p} = h_\theta(X)$ that an instance $X$ belongs to the positive class has been estimated using Logistic Regression, the prediction ($\hat{y}$) can be made. $\hat{y}$ is calculated as

$$\hat{y} = \begin{cases} 0 & \text{if } \hat{p} < 0.5, \\ 1 & \text{if } \hat{p} \geq 0.5. \end{cases} \tag{4.3}$$

Note that $\sigma(x) < 0.5$ when $x < 0$, and $\sigma(x) \geq 0.5$ when $x \geq 0$. Therefore, the Logistic Regression model predicts 0 if $\theta^{\mathrm{T}} \cdot X$ is negative, and 0 if it is positive.

Now that we understand how logistic regression assigns probabilities and produces predictions, let's walk through a concise example that shows its training procedure and the associated cost function. Training seeks parameter values $\theta$ that give high predicted probabilities to positive examples ($y = 1$) and low probabilities to negative ones ($y = 0$). This aim is captured by the cost defined in equation 4.4, which evaluates a single training sample $x$. To achieve this, we require a cost function, such as

$$C(\theta) = \begin{cases} -\log(\hat{p}) & \text{if } y = 1, \\ -\log(1 - \hat{p}) & \text{if } y = 0. \end{cases} \tag{4.4}$$

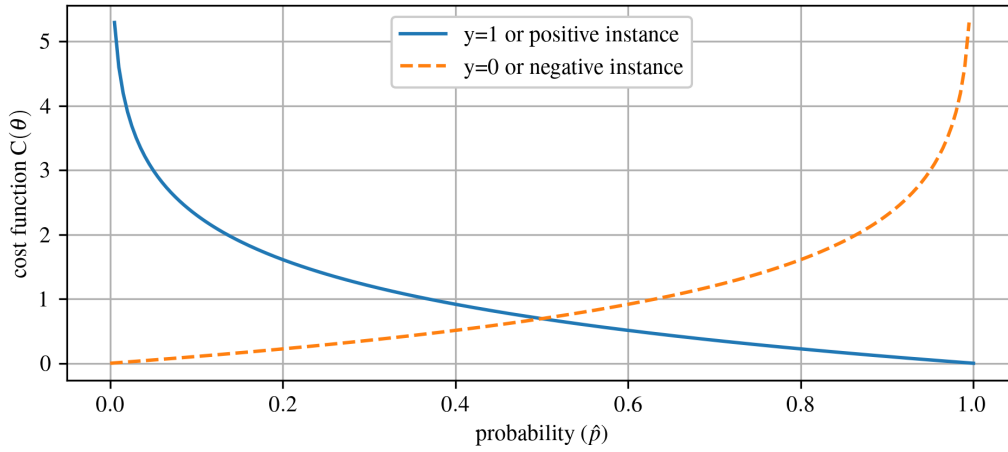The considered cost function is plotted in figure 4.2.



Figure 4.2: Cost function behavior for classification that heavily penalizes incorrect predictions.

The cost function in equation 4.4 behaves intuitively: $-\log(\hat{p})$ increases sharply as $\hat{p} \to 0$, so the loss is large when the model assigns a probability near 0 to a positive example. It likewise produces a high loss when the model predicts a probability close to 1 for a negative example. Conversely, because $-\log(\hat{p}) \to 0$ as $\hat{p} \to 1$, the loss becomes negligible when the predicted probability is near 1 for a positive instance or near 0 for a negative one, aligning with our expectations.

The overall cost, denoted $J(\theta)$, is the mean loss across all $m$ training examples. This metric is commonly called the log-loss and written as

$$J(\theta); =; -\frac{1}{m} \sum_{i=1}^{m} \left[ y^{(i)} \log\left(\hat{p}^{(i)}\right); +; \left(1 - y^{(i)}\right) \log\left(1 - \hat{p}^{(i)}\right) \right]. \tag{4.5}$$

Because there is no closed-form solution for the parameters $\theta$ (embedded within $\hat{p}$), the minimum must be found with an iterative optimizer. The cost surface is convex, so gradient descent or another suitable algorithm will reach the global minimum provided the learning rate is reasonable and enough iterations are allowed. The gradient of the cost with respect to the $j^{\text{th}}$ parameter $\theta_j$ is given as

$$\frac{\partial J}{\partial \theta_j}; =; \frac{1}{m} \sum_{i=1}^{m} \left[ \sigma\left(\theta^{\text{T}} X^{(i)}\right) - y^{(i)} \right], x_j^{(i)}. \tag{4.6}$$

2

This equation resembles the partial derivative used in gradient descent. For every training example, it finds the prediction error, multiplies it by the $j^{\text{th}}$ feature value, and then averages these products over all $m$ samples. With the resulting gradient vector of partial derivatives, you can update the parameters using the Batch Gradient Descent algorithm, completing the training of a Logistic Regression model. Stochastic Gradient Descent performs the update after each single example, while Mini-batch Gradient Descent updates the parameters after processing each mini-batch.

**Review 4.1  Iris Flower Dataset**

The Iris dataset, collected by botanist Edgar Anderson in 1935, lists sepal length, sepal width, petal length, and petal width for 150 flowers drawn equally from three species: Iris-Setosa, Iris-Versicolor, and Iris-Virginic (refer to Figure 4.3)a. Statistician Ronald Fisher analysed these measurements in a 1936 study on linear discriminant analysis, and the data have since become a standard benchmark in statistics and machine learning because they are small, clean, and perfectly balanced, making them ideal for testing classification algorithms and visualisation techniques.
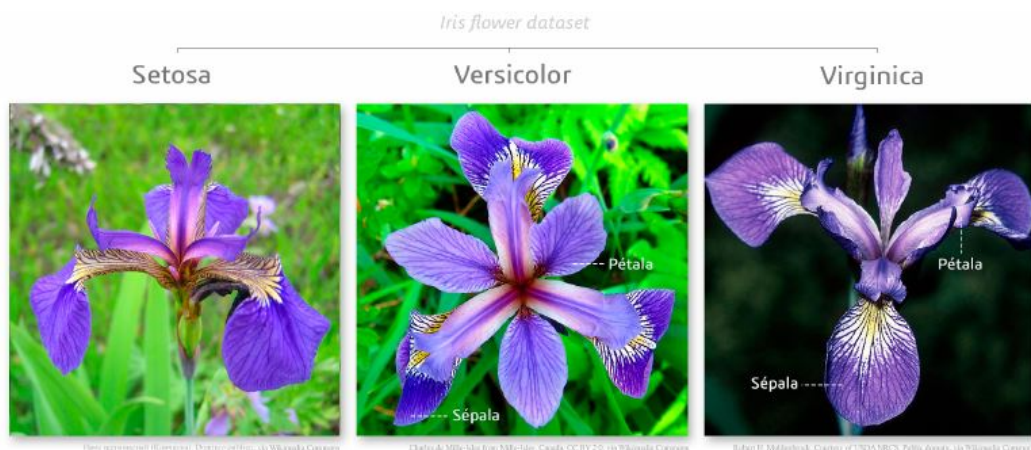


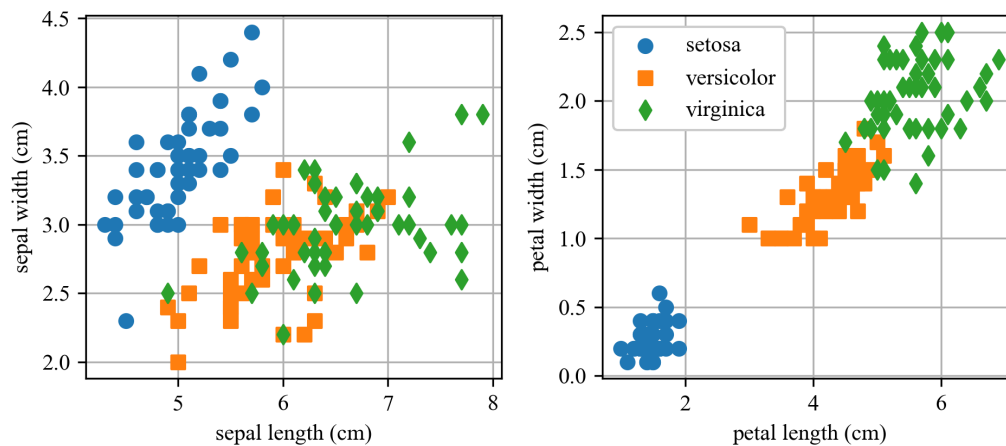Figure 4.3: Flowers representing three species of iris plants [a]

Figure 4.4: Iris dataset scatterplot showing sepal length vs. sepal width (left) and petal length vs. petal width (right).

**Example 4.1  Iris Dataset Exploration**
This example introduces the Iris dataset (Figure 4.3), originally compiled by Ronald Fisher. It demonstrates how to load the dataset, access feature and label information, and visualize the relationships between sepal and petal measurements across the three iris species.

### 4.1.1    1-D Decision Boundaries

The petal width of Iris-virginica flowers usually lies between 1.4 cm and 2.5 cm, whereas the other two species range from 0.1 cm to 1.8 cm. Because these intervals overlap, the classifier's certainty changes across the scale. Above roughly 2 cm, the model is confident that any given sample is Iris-virginica. Moreover, below 1 cm it is equally sure the flower is not Iris-virginica (This means the model redicts high probability for the class "Not Iris-Virginica"). In between the limits of the classes, the model is unsure. When you call `predict()` instead of `predict_proba()`, it outputs the most likely class, creating a decision boundary around 1.6 cm where both class probabilities reach 50 percent. Thus, flowers with petal widths greater than 1.6 cm lead to a prediction of Iris-virginica, while flowers with smaller petal widths yield the opposite prediction even if confidence is low.
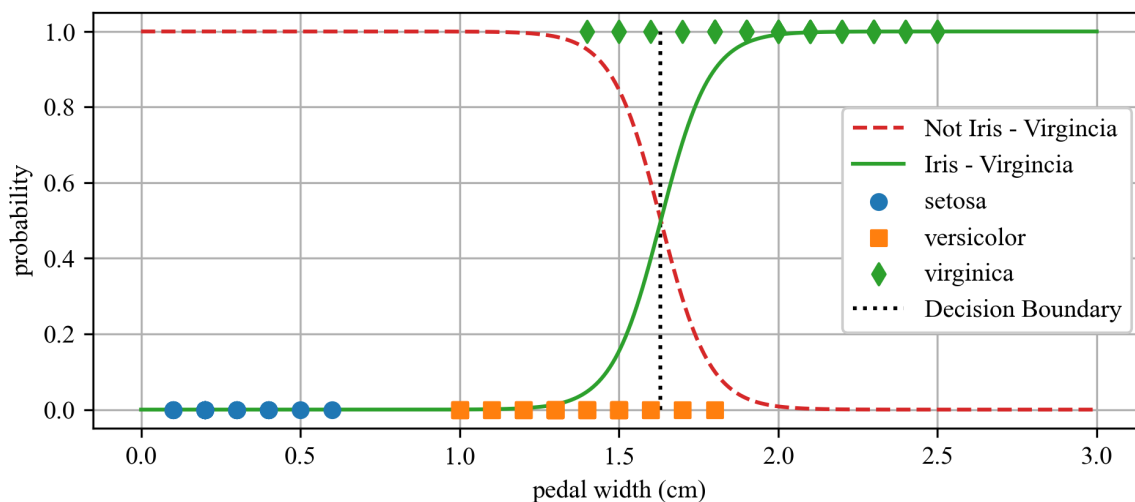
Figure 4.5: Decision boundary for the flowers of three Iris plant species with C set to $C = 10^{10}$ .

---

**Example 4.2  1D Logistic Regression**

This example builds a logistic regression model to classify Iris-Virginica based on a single feature: petal width. It visualizes class probabilities and shows the decision boundary at the 50% threshold.

---

### 4.1.2   2-D Decision Boundaries

Figure 4.6 plots petal width against petal length for the Iris samples. After training, the logistic regression model assigns each point the probability that the flower is *Irisvirginica*. The dashed line marks where this probability equals 50 %, forming the decision boundary, which is linear. The parallel contour lines show equal-probability levels from 10% to 90%. Points be beyond the top-right line have a greater than 90% probability of being classified as Iris-Virginica by the model.
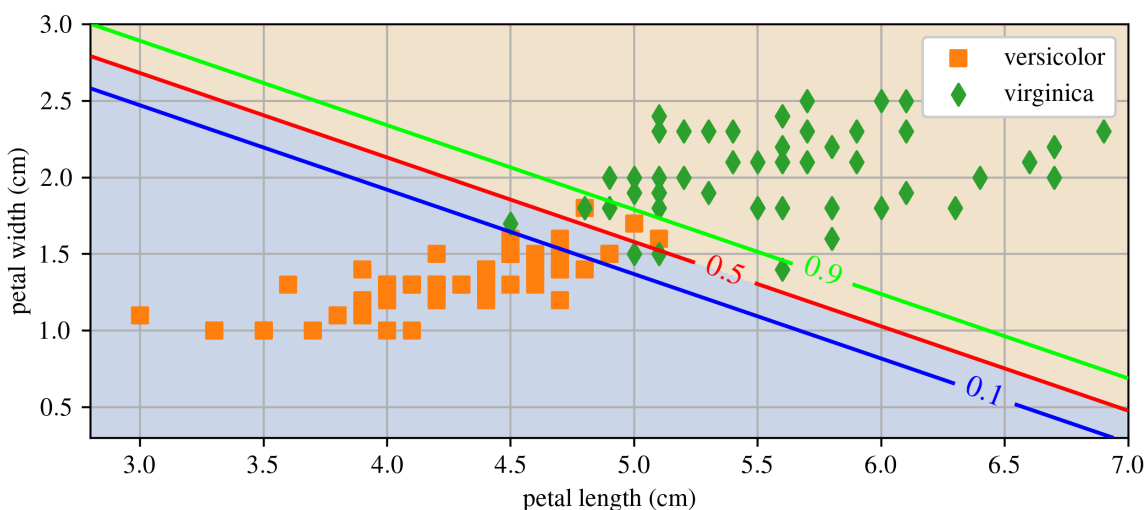


Figure 4.6: A 2D decision boundary for the Iris dataset.

---

**Example 4.3  2D Decision Boundary**

This example uses logistic regression to classify the Iris-Virginica species based on petal length and width. A 2D decision boundary is visualized in the petal feature space, with prediction regions and probability contours illustrating classifier confidence.

---

## 4.2   Softmax Regression

Logistic Regression can be generalised to handle many classes in a single model, so there is no need to train and merge several binary classifiers. This multiclass version is called Softmax Regression, or Multinomial Logistic Regression. Softmax Regression provides a straightforward way to perform regression-based classification when more than two classes are present. For an input vector $\mathbf{x}$, the model first computes a score $s_k(\mathbf{x})$ for every class $k$, then turns these scores into class probabilities with the softmax (normalised exponential) function. The score is obtained exactly as in linear regression:

$$s_k(\mathbf{x}) = \mathbf{x}^{\mathrm{T}} \boldsymbol{\theta}^{(k)}. \tag{4.7}$$

Here each class has its own parameter vector $\boldsymbol{\theta}^{(k)}$, and the collection of these vectors is usually stored as the rows of a parameter matrix $\Theta$.

After the model has calculated a score for each class given an input $\mathbf{x}$, it converts these scores to probabilities with the softmax function. For class $k$ the predicted probability is

$$\hat{p}_k = \sigma\big(s(\mathbf{x})\big)_k = \frac{\exp\big(s_k(\mathbf{x})\big)}{\sum_{j=1}^{K} \exp\big(s_j(\mathbf{x})\big)}. \tag{4.8}$$

In this expression, $\mathbf{s}(\mathbf{x})$ is the vector of class scores for the input, $\sigma\big(\mathbf{s}(\mathbf{x})\big)_k$ is the probability that $\mathbf{x}$ belongs to class $k$, and $K$ is the total number of classes. In short, equation 4.8 exponentiates each score and then normalises the results by dividing by the sum of all exponentials so the probabilities sum to one.

Like logistic regression, a Softmax Regression model assigns an input to the class whose predicted probability is largest, which coincides with the class that has the highest score:

$$\hat{y} = \underset{k}{\operatorname{argmax}} \, \sigma\big(\mathbf{s}(\mathbf{x})\big)_k = \underset{k}{\operatorname{argmax}} \, s_k(\mathbf{x}) = \underset{k}{\operatorname{argmax}} \big((\boldsymbol{\theta}^{(k)})^{\mathrm{T}} \cdot \mathbf{x}\big). \tag{4.9}$$

---

**NOTE**

The argmax operator returns the argument that maximises a function. In this setting, it yields the index $k$ for which the estimated probability $\sigma\big(s(\mathbf{x})\big)_k$ attains its largest value.

---

With probability estimation and prediction established, we next consider training. The task is to learn parameters that place a large probability on the correct class and correspondingly small probabilities on all others. This objective is met by minimising the cost in equation 4.10, known as the cross entropy, which heavily penalises the model when it assigns a low probability to the true

label. Cross entropy is widely used to measure how well predicted class probabilities agree with the actual classes. The cost function is

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{K} y_k^{(i)} \log(\hat{p}_k^{(i)}). \tag{4.10}$$

---

**NOTE**

In this expression, $y_k^{(i)} = 1$ when the $i^{\text{th}}$ sample's true label is class $k$, and $y_k^{(i)} = 0$ otherwise.

---

When considering only two classes ($K = 2$), it's important to highlight that this cost function aligns with the Logistic Regression's cost function, commonly referred to as log loss (refer to Equation 4.5).

The gradient of the cross-entropy cost with respect to $\theta^{(k)}$ is

$$\nabla_{\theta^{(k)}} J(\Theta) = \frac{1}{m} \sum_{i=1}^{m} (\hat{p}_k^{(i)} - y_k^{(i)} X^{(i)}). \tag{4.11}$$

By computing this vector for each class, you obtain the full gradient, which can then be fed to Gradient Descent or another optimizer to find the parameter matrix $\Theta$ that minimises the cost.

Applying Softmax Regression to the three-class iris problem in Scikit-Learn is straightforward. The `LogisticRegression` estimator normally uses one-versus-all when more than two classes are present, but setting `multi_class=''multinomial''` activates true Softmax learning. Choose a solver that supports this option, such as `lbfgs` (see the library documentation). The model applies $\ell_2$ regularisation by default, governed by the hyperparameter $C$. After training, a flower with 5 cm long and 2 cm wide petals is classified as Iris-Virginica with probability 94.2%, while the probability of Iris versicolor is 5.8%.

---

**NOTE**

The Softmax Regression classifier is multiclass, not multioutput. As such, Softmax Regression can only predict one class at a time, so it works for problems where each input belongs to exactly one category-like classifying an email as spam, promotions, or updates. It can't be used for cases where multiple labels may apply, such as tagging a news article with topics like politics, economics, and technology all at once.

---

Figure 4.7 displays the decision regions, shaded with background colours to mark each class. The borders separating any two classes are straight lines. The plot also includes contour curves for the predicted probability of the Iris-Versicolor class. At the point where all three borders meet, every class receives a probability of 33%, so the selected class can have a confidence below 50%.
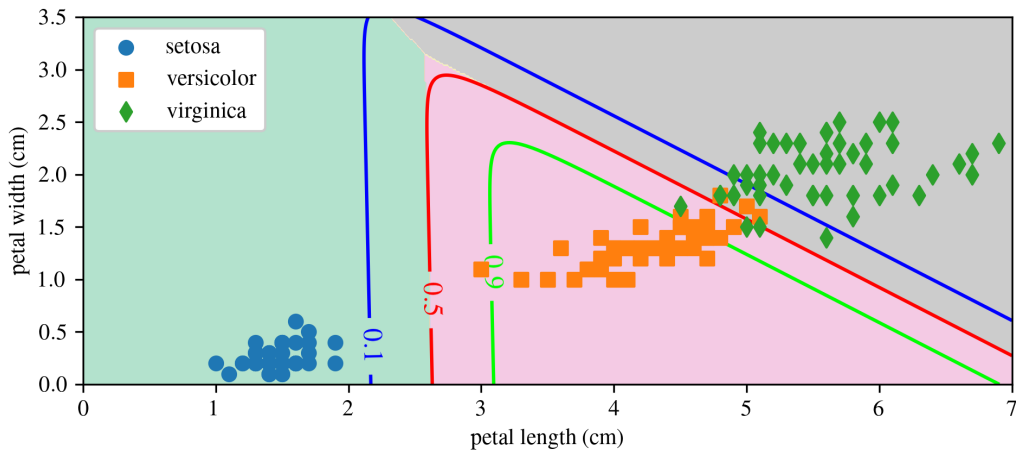
Figure 4.7: Softmax classification for the three iris plant species.

**Example 4.4  Softmax Classification**
This example applies softmax regression to classify all three iris species using petal length and width. It builds a multinomial logistic regression model and visualizes the decision boundaries along with confidence contours over the 2D petal feature space.

## 4.3   Examples

### Example 4.1

```python
1   """
2   Example 4.1 Introduction to the IRIS data set
3   @author: austin_downey
4   """
5
6   import IPython as IP
7   IP.get_ipython().run_line_magic('reset', '-sf')
8
9   import matplotlib.pyplot as plt
10  import sklearn as sk
11
12  cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
13  plt.close('all')
14
15
16  #%% Load your data
17
18  # We will use the Iris data set.  This dataset was created by biologist Ronald
19  # Fisher in his 1936 paper "The use of multiple measurements in taxonomic
20  # problems" as an example of linear discriminant analysis
21
22  iris = sk.datasets.load_iris()
23
24  # for simplicity, extract some of the data sets
25  X = iris['data'] # this contains the length of the pedals and sepals
26  Y = iris['target'] # contains what type of flower it is
27  Y_names = iris['target_names'] # contains the name that aligns with the type of the
    flower
28  feature_names = iris['feature_names'] # the names of the features
29
30  # plot the Sepal data
31  plt.figure(figsize=(6.5,3))
32  plt.subplot(121)
33  plt.grid(True)
34  plt.scatter(X[Y==0,0],X[Y==0,1],marker='o')
35  plt.scatter(X[Y==1,0],X[Y==1,1],marker='s')
36  plt.scatter(X[Y==2,0],X[Y==2,1],marker='d')
37  plt.xlabel(feature_names[0])
38  plt.ylabel(feature_names[1])
39
40
41  plt.subplot(122)
42  plt.grid(True)
43  plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0])
44  plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1])
45  plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2])
46  plt.xlabel(feature_names[2])
47  plt.ylabel(feature_names[3])
48  plt.legend(framealpha=1)
49  plt.tight_layout()
50
```

## Example 4.2

```python
1  """
2  Example 4.2 1D Decision boundary for the Iris dataset
3  @author: austin_downey
4  """
5
6  import IPython as IP
7  IP.get_ipython().run_line_magic('reset', '-sf')
8
9  import numpy as np
10 import matplotlib.pyplot as plt
11 import sklearn as sk
12
13
14 cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
15 plt.close('all')
16
17
18 #%% Load your data
19
20 # We will use the Iris data set.  This dataset was created by biologist Ronald
21 # Fisher in his 1936 paper "The use of multiple measurements in taxonomic
22 # problems" as an example of linear discriminant analysis
23
24 iris = sk.datasets.load_iris()
25
26 # for simplicity, extract some of the data sets
27 X = iris['data'] # this contains the length of the pedals and sepals
28 Y = iris['target'] # contains what type of flower it is
29 Y_names = iris['target_names'] # contains the name that aligns with the type of the
   flower
30 feature_names = iris['feature_names'] # the names of the features
31
32 # plot the Sepal data
33 plt.figure(figsize=(6.5,3))
34 plt.subplot(121)
35 plt.grid(True)
36 plt.scatter(X[Y==0,0],X[Y==0,1],marker='o')
37 plt.scatter(X[Y==1,0],X[Y==1,1],marker='s')
38 plt.scatter(X[Y==2,0],X[Y==2,1],marker='d')
39 plt.xlabel(feature_names[0])
40 plt.ylabel(feature_names[1])
41
42
43 plt.subplot(122)
44 plt.grid(True)
45 plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0])
46 plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1])
47 plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2])
48 plt.xlabel(feature_names[2])
49 plt.ylabel(feature_names[3])
50 plt.legend(framealpha=1)
51 plt.tight_layout()
52
53
54 #%% Train a Logistic Regression model
55
56 # define the features (X) and the output (Y)
57 X_pedal = iris["data"][:, 3:] # consider just the petal width
58 y_pedal = iris["target"] == 2 # 1 if Iris-Virginica, else 0
59
60 # Build the logistic Regression model and train it.
61 log_reg = sk.linear_model.LogisticRegression( C=1)
62 log_reg.fit(X_pedal, y_pedal)
63 # Note: The hyper-parameter controlling the regularization strength of a
64 # Scikit-Learn LogisticRegression model is not alpha (as in other linear models),
65 # but its  inverse: C. The higher the value of C, the less the model is regularized.
66
```

```
67   # Build a range of the feature (X) to predict over. Here we just consider pedal width.
68   X_new = np.linspace(0, 3, 1000)
69   X_new = np.expand_dims(X_new, axis=1)
70
71   # Use the Logistic Regression Model to predict the pedal type based on pedal width
72   y_proba = log_reg.predict_proba(X_new)
73
74   # plot the probability plots
75   plt.figure(figsize=(6.5,3))
76   plt.grid(True)
77
78   # plot the data used for training, set at 0 and 1
79   plt.scatter(X[Y==0,3],np.zeros(50),marker='o',label=Y_names[0],zorder=10)
80   plt.scatter(X[Y==1,3],np.zeros(50),marker='s',label=Y_names[1],zorder=10)
81   plt.scatter(X[Y==2,3],np.ones(50),marker='d',label=Y_names[2],zorder=10)
82
83   # plot the probability
84   plt.plot(X_new,y_proba[:,0], '--',color=cc[3],label='Not Iris - Virgincia')
85   plt.plot(X_new,y_proba[:,1],color=cc[2], label='Iris - Virgincia')
86
87   # find and plot the 50% decision boundary
88   x_at_50 = X_new[np.argmin(np.abs(y_proba[:,0] - 0.5))]
89   plt.vlines(x_at_50,0,1,color='k',linestyles=':',label='Decision Boundary')
90
91   plt.xlabel('pedal width (cm)')
92   plt.ylabel('probability')
93   plt.legend(framealpha=1)
94   plt.tight_layout()
95
96   # make predictions on the model trained on the data.
97   log_reg.predict([[1.7]])
```

**Example 4.3**

```python
1   #!/usr/bin/env python3
2   # -*- coding: utf-8 -*-
3   """
4   Example 4.3 2D Decision boundary for the Iris dataset
5
6   Developed for Machine Learning for Mechanical Engineers at the University of
7   South Carolina
8
9   @author: austin_downey
10  """
11
12  import IPython as IP
13  IP.get_ipython().run_line_magic('reset', '-sf')
14
15
16  import numpy as np
17  import matplotlib.pyplot as plt
18  import sklearn as sk
19
20
21  cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
22  plt.close('all')
23
24
25  #%% Load your data
26
27  # We will use the Iris data set.  This dataset was created by biologist Ronald
28  # Fisher in his 1936 paper "The use of multiple measurements in taxonomic
29  # problems" as an example of linear discriminant analysis
30
31  iris = sk.datasets.load_iris()
32
33  # for simplicity, extract some of the data sets
34  X = iris['data'] # this contains the length of the petals and sepals
35  Y = iris['target'] # contains what type of flower it is
36  Y_names = iris['target_names'] # contains the name that aligns with the type of the
        flower
37  feature_names = iris['feature_names'] # the names of the features
38
39  # plot the Sepal data
40  plt.figure(figsize=(6.5,3))
41  plt.subplot(121)
42  plt.grid(True)
43  plt.scatter(X[Y==0,0],X[Y==0,1],marker='o',zorder=10)
44  plt.scatter(X[Y==1,0],X[Y==1,1],marker='s',zorder=10)
45  plt.scatter(X[Y==2,0],X[Y==2,1],marker='d',zorder=10)
46  plt.xlabel(feature_names[0])
47  plt.ylabel(feature_names[1])
48
49
50  plt.subplot(122)
51  plt.grid(True)
52  plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0],zorder=10)
53  plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1],zorder=10)
54  plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2],zorder=10)
55  plt.xlabel(feature_names[2])
56  plt.ylabel(feature_names[3])
57  plt.legend(framealpha=1)
58  plt.tight_layout()
59
60
61  #%% plot the Linear decision boundary in 2D "Petal" space
62
63  # build the training and target set.
64  X_train = X[:, (2, 3)]  # petal length, petal width
65  y_train = Y == 2
66
```

```
67    # build the Logistic Regression model
68    log_reg = sk.linear_model.LogisticRegression(C=10**10)
69    # Note: The hyper-parameter controlling the regularization strength of a Scikit-Learn
70    # LogisticRegression model is not alpha (as in other linear models), but its
71    # inverse: C. The higher the value of C, the less the model is regularized.
72
73    # train the Logistic Regression model
74    log_reg.fit(X_train, y_train)
75
76    # build the x values for the predictions over the entire "petal space"
77    x_grid, y_grid = np.meshgrid(
78          np.linspace(2.8, 7, 500),
79          np.linspace(0.3, 3, 200),
80       )
81    X_new = np.vstack((x_grid.reshape(-1), y_grid.reshape(-1))).T # build a vector format of
      the mesh grid
82
83    # predict on the vectorized format
84    y_predict = log_reg.predict(X_new)
85    y_proba = log_reg.predict_proba(X_new)
86
87    # convert back to meshgrid shape for plotting
88    zz_predict = y_predict.reshape(x_grid.shape)
89    zz_proba = y_proba[:, 1].reshape(x_grid.shape)
90
91    # plot the 2D "petal space"
92    plt.figure(figsize=(6.5,3))
93    plt.grid(True)
94    plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',color=cc[1],label=Y_names[1],zorder=10)
95    plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',color=cc[2],label=Y_names[2],zorder=10)
96    plt.contourf(x_grid, y_grid, zz_predict, cmap='Pastel2')
97    contour = plt.contour(x_grid, y_grid, zz_proba, [0.100,0.5,0.900],cmap=plt.cm.brg)
98    plt.clabel(contour, inline=1) # add the labels to the plot
99    plt.xlabel(feature_names[2])
100   plt.ylabel(feature_names[3])
101   plt.legend()
102   plt.tight_layout()
103
104
105
```

## Example 4.4

```
1    """
2    Example 4.4 Softmax decision boundary for the Iris dataset
3    @author: austin_downey
4    """
5
6    import IPython as IP
7    IP.get_ipython().magic('reset -sf')
8
9    import numpy as np
10   import matplotlib.pyplot as plt
11   import sklearn as sk
12
13
14   cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
15   plt.close('all')
16
17
18   #%% Load your data
19
20   # We will use the Iris data set.  This dataset was created by biologist Ronald
21   # Fisher in his 1936 paper "The use of multiple measurements in taxonomic
22   # problems" as an example of linear discriminant analysis
23   iris = sk.datasets.load_iris()
24
25   # for simplicity, extract some of the data sets
26   X = iris['data'] # this contains the length of the petals and sepals
27   Y = iris['target'] # contains what type of flower it is
28   Y_names = iris['target_names'] # contains the name that aligns with the type of the
     flower
29   feature_names = iris['feature_names'] # the names of the features
30
31   # plot the Sepal data
32   plt.figure(figsize=(6.5,3))
33   plt.subplot(121)
34   plt.grid(True)
35   plt.scatter(X[Y==0,0],X[Y==0,1],marker='o',zorder=10)
36   plt.scatter(X[Y==1,0],X[Y==1,1],marker='s',zorder=10)
37   plt.scatter(X[Y==2,0],X[Y==2,1],marker='d',zorder=10)
38   plt.xlabel(feature_names[0])
39   plt.ylabel(feature_names[1])
40
41   plt.subplot(122)
42   plt.grid(True)
43   plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0],zorder=10)
44   plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1],zorder=10)
45   plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2],zorder=10)
46   plt.xlabel(feature_names[2])
47   plt.ylabel(feature_names[3])
48   plt.legend(framealpha=1)
49   plt.tight_layout()
50
51
52   #%% Softmax Regression
53
54   # build the training and target set.
55   X_train = X[:, (2, 3)]  # petal length, petal width
56   y_train = Y
57
58   # build and train the softmax model
59   softmax_reg = sk.linear_model.LogisticRegression(multi_class="multinomial",
60           solver="lbfgs", C=10)
61   softmax_reg.fit(X_train, y_train)
62
63   # build the x values for the predictions over the entire "petal space"
64   x_grid, y_grid = np.meshgrid(
65           np.linspace(0, 7, 500),
66           np.linspace(0, 4, 200),
```

```
67        )
68   X_new = np.vstack((x_grid.reshape(-1), y_grid.reshape(-1))).T # build a vector format of
     the mesh grid
69
70   # predict on the vectorized format
71   y_predict = softmax_reg.predict(X_new)
72   y_proba = softmax_reg.predict_proba(X_new)
73
74
75   # convert back to meshgrid shape for plotting
76   zz_predict = y_predict.reshape(x_grid.shape)
77   zz_proba = y_proba[:, 1].reshape(x_grid.shape) # the selected column selects the
     probability that the data falls within this class.
78
79   # plot the 2D "petal space"
80   plt.figure(figsize=(6.5, 4))
81   plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0],zorder=10)
82   plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1],zorder=10)
83   plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2],zorder=10)
84   plt.contourf(x_grid, y_grid, zz_predict, cmap='Pastel2')
85   contour = plt.contour(x_grid, y_grid, zz_proba, [0.100,0.5,0.900], cmap=plt.cm.brg)
86   plt.clabel(contour, inline=1)
87   plt.xlabel(feature_names[2])
88   plt.ylabel(feature_names[3])
89   plt.legend()
90   plt.tight_layout()
```