**Chapter 4 – Training Models**

*This notebook contains all the sample code and solutions to the exercises in chapter 4.*

[CO Open in Colab] [k Open in Kaggle]

## ▾ Setup

First, let's import a few common modules, ensure MatplotLib plots figures inline and prepare a function to save the figures. We also check that Python 3.5 or later is installed (although Python 2.x may work, it is deprecated so we strongly recommend you use Python 3 instead), as well as Scikit-Learn ≥0.20.

```python
1  # Python ≥3.5 is required
2  import sys
3  assert sys.version_info >= (3, 5)
4
5  # Scikit-Learn ≥0.20 is required
6  import sklearn
7  assert sklearn.__version__ >= "0.20"
8
9  # Common imports
10 import numpy as np
11 import os
12
13 # to make this notebook's output stable across runs
14 np.random.seed(42)
15
16 # To plot pretty figures
17 %matplotlib inline
18 import matplotlib as mpl
19 import matplotlib.pyplot as plt
20 mpl.rc('axes', labelsize=14)
21 mpl.rc('xtick', labelsize=12)
22 mpl.rc('ytick', labelsize=12)
23
24 # Where to save the figures
25 PROJECT_ROOT_DIR = "."
26 CHAPTER_ID = "training_linear_models"
27 IMAGES_PATH = os.path.join(PROJECT_ROOT_DIR, "images", CHAPTER_ID)
28 os.makedirs(IMAGES_PATH, exist_ok=True)
29
30 def save_fig(fig_id, tight_layout=True, fig_extension="png", resolution=300):
31     path = os.path.join(IMAGES_PATH, fig_id + "." + fig_extension)
32     print("Saving figure", fig_id)
33     if tight_layout:
34         plt.tight_layout()
35     plt.savefig(path, format=fig_extension, dpi=resolution)
```

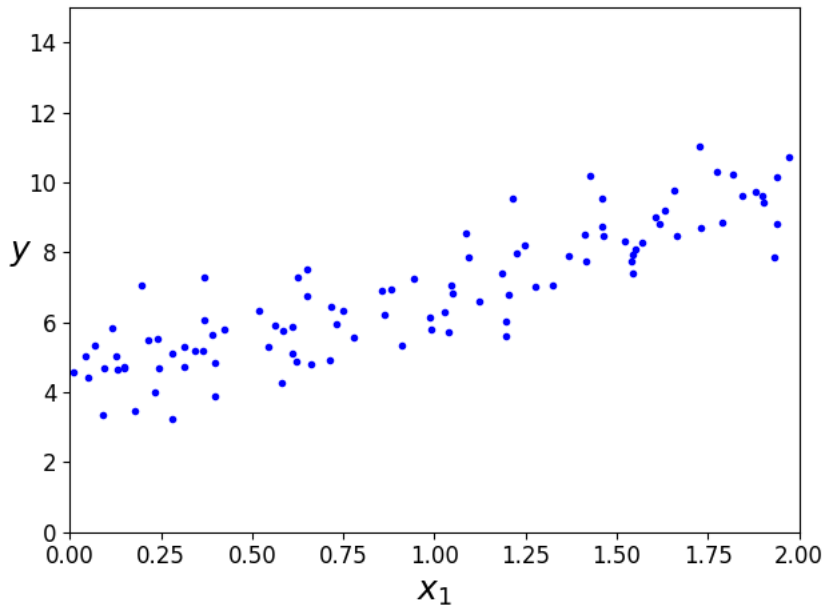## ▾ Linear Regression

## ▾ The Normal Equation

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

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

Saving figure generated_data_plot



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

```
1 theta_best
```

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

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

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

```
1 plt.plot(X_new, y_predict, "r-")
2 plt.plot(X, y, "b.")
3 plt.axis([0, 2, 0, 15])
4 plt.show()
```
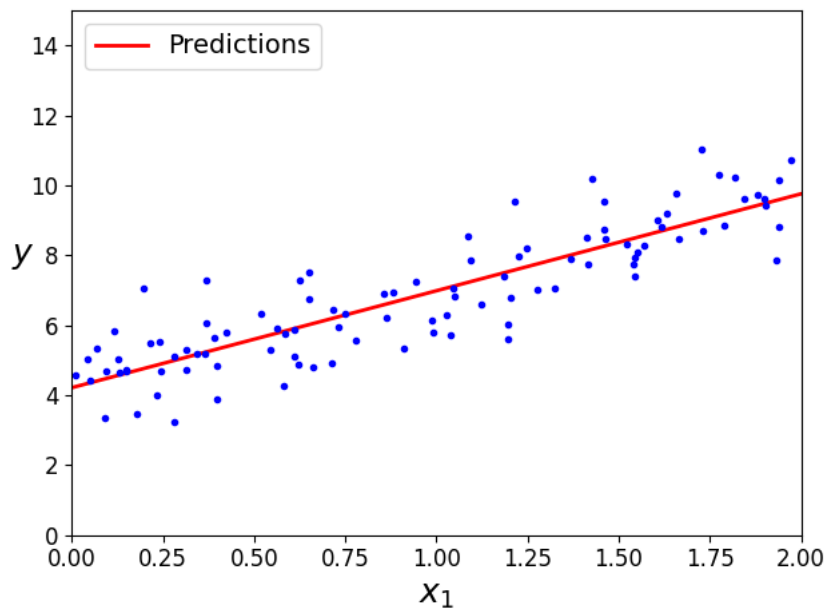
14 ┤

The figure in the book actually corresponds to the following code, with a legend and axis labels:

12 ┤

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

Saving figure linear_model_predictions_plot



```
1 from sklearn.linear_model import LinearRegression
2
3 lin_reg = LinearRegression()
4 lin_reg.fit(X, y)
5 lin_reg.intercept_, lin_reg.coef_
```

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

```
1 lin_reg.predict(X_new)
```

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

The `LinearRegression` class is based on the `scipy.linalg.lstsq()` function (the name stands for "least squares"), which you could call directly:

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

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

This function computes $\mathbf{X}^+\mathbf{y}$, where $\mathbf{X}^+$ is the *pseudoinverse* of $\mathbf{X}$ (specifically the Moore-Penrose inverse). You can use `np.linalg.pinv()` to compute the pseudoinverse directly:

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

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

# ▾ Gradient Descent

## Batch Gradient Descent

```
1 eta = 0.1  # learning rate
2 n_iterations = 1000
3 m = 100
4
5 theta = np.random.randn(2,1)  # random initialization
6
7 for iteration in range(n_iterations):
8     gradients = 2/m * X_b.T.dot(X_b.dot(theta) - y)
9     theta = theta - eta * gradients
```

```
1 theta
```

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

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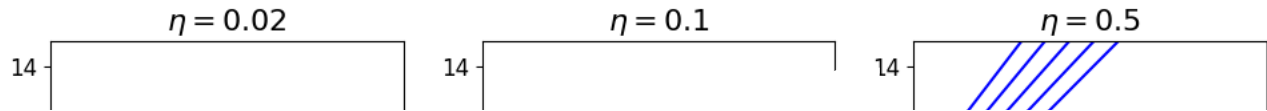
```
array([[4.21509616],
       [9.75532293]])
```

```
1 theta_path_bgd = []
2
3 def plot_gradient_descent(theta, eta, theta_path=None):
4     m = len(X_b)
5     plt.plot(X, y, "b.")
6     n_iterations = 1000
7     for iteration in range(n_iterations):
8         if iteration < 10:
9             y_predict = X_new_b.dot(theta)
10            style = "b-" if iteration > 0 else "r--"
11            plt.plot(X_new, y_predict, style)
12        gradients = 2/m * X_b.T.dot(X_b.dot(theta) - y)
13        theta = theta - eta * gradients
14        if theta_path is not None:
15            theta_path.append(theta)
16    plt.xlabel("$x_1$", fontsize=18)
17    plt.axis([0, 2, 0, 15])
18    plt.title(r"$\eta = {}$".format(eta), fontsize=16)
```

```
1 np.random.seed(42)
2 theta = np.random.randn(2,1)  # random initialization
3
4 plt.figure(figsize=(10,4))
5 plt.subplot(131); plot_gradient_descent(theta, eta=0.02)
6 plt.ylabel("$y$", rotation=0, fontsize=18)
7 plt.subplot(132); plot_gradient_descent(theta, eta=0.1, theta_path=theta_path_bgd)
8 plt.subplot(133); plot_gradient_descent(theta, eta=0.5)
9
10 save_fig("gradient_descent_plot")
11 plt.show()
```
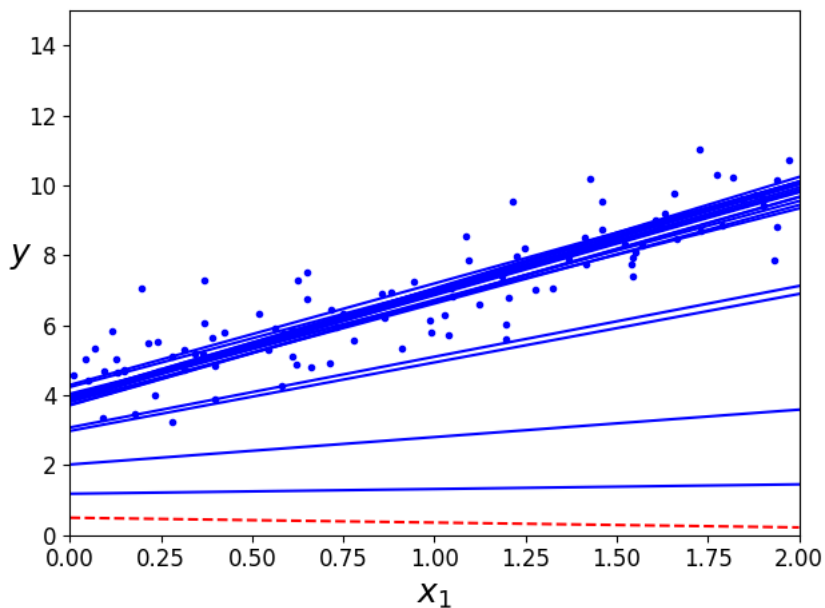
Saving figure gradient_descent_plot

|  $\eta = 0.02$  |  $\eta = 0.1$  |  $\eta = 0.5$  |
|---|---|---|
| 14 | 14 | 14 |

## ▾ Stochastic Gradient Descent

```
1 theta_path_sgd = []
2 m = len(X_b)
3 np.random.seed(42)
```

```
1 n_epochs = 50
2 t0, t1 = 5, 50  # learning schedule hyperparameters
3
4 def learning_schedule(t):
5     return t0 / (t + t1)
6
7 theta = np.random.randn(2,1)  # random initialization
8
9 for epoch in range(n_epochs):
10    for i in range(m):
11        if epoch == 0 and i < 20:                # not shown in the book
12            y_predict = X_new_b.dot(theta)       # not shown
13            style = "b-" if i > 0 else "r--"     # not shown
14            plt.plot(X_new, y_predict, style)    # not shown
15        random_index = np.random.randint(m)
16        xi = X_b[random_index:random_index+1]
17        yi = y[random_index:random_index+1]
18        gradients = 2 * xi.T.dot(xi.dot(theta) - yi)
19        eta = learning_schedule(epoch * m + i)
20        theta = theta - eta * gradients
21        theta_path_sgd.append(theta)             # not shown
22
23 plt.plot(X, y, "b.")                            # not shown
24 plt.xlabel("$x_1$", fontsize=18)               # not shown
25 plt.ylabel("$y$", rotation=0, fontsize=18)     # not shown
26 plt.axis([0, 2, 0, 15])                        # not shown
27 save_fig("sgd_plot")                           # not shown
28 plt.show()                                     # not shown
```

Saving figure sgd_plot



```
1 theta
```

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

```
1 from sklearn.linear_model import SGDRegressor
2
3 sgd_reg = SGDRegressor(max_iter=1000, tol=1e-3, penalty=None, eta0=0.1, random_state=42)
4 sgd_reg.fit(X, y.ravel())
```

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

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

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

## ▾ Mini-batch gradient descent

```
 1 theta_path_mgd = []
 2
 3 n_iterations = 50
 4 minibatch_size = 20
 5
 6 np.random.seed(42)
 7 theta = np.random.randn(2,1)  # random initialization
 8
 9 t0, t1 = 200, 1000
10 def learning_schedule(t):
11     return t0 / (t + t1)
12
13 t = 0
14 for epoch in range(n_iterations):
15     shuffled_indices = np.random.permutation(m)
16     X_b_shuffled = X_b[shuffled_indices]
17     y_shuffled = y[shuffled_indices]
18     for i in range(0, m, minibatch_size):
19         t += 1
20         xi = X_b_shuffled[i:i+minibatch_size]
21         yi = y_shuffled[i:i+minibatch_size]
22         gradients = 2/minibatch_size * xi.T.dot(xi.dot(theta) - yi)
23         eta = learning_schedule(t)
24         theta = theta - eta * gradients
25         theta_path_mgd.append(theta)
```

```
1 theta
```

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

```
1 theta_path_bgd = np.array(theta_path_bgd)
2 theta_path_sgd = np.array(theta_path_sgd)
3 theta_path_mgd = np.array(theta_path_mgd)
```
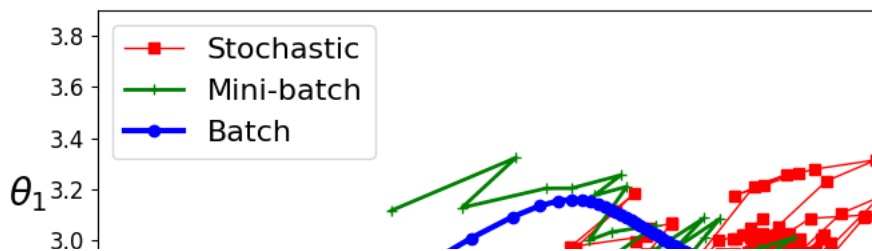
```
 1 plt.figure(figsize=(7,4))
 2 plt.plot(theta_path_sgd[:, 0], theta_path_sgd[:, 1], "r-s", linewidth=1, label="Stochastic")
 3 plt.plot(theta_path_mgd[:, 0], theta_path_mgd[:, 1], "g-+", linewidth=2, label="Mini-batch")
 4 plt.plot(theta_path_bgd[:, 0], theta_path_bgd[:, 1], "b-o", linewidth=3, label="Batch")
 5 plt.legend(loc="upper left", fontsize=16)
 6 plt.xlabel(r"$\theta_0$", fontsize=20)
 7 plt.ylabel(r"$\theta_1$    ", fontsize=20, rotation=0)
 8 plt.axis([2.5, 4.5, 2.3, 3.9])
 9 save_fig("gradient_descent_paths_plot")
10 plt.show()
```
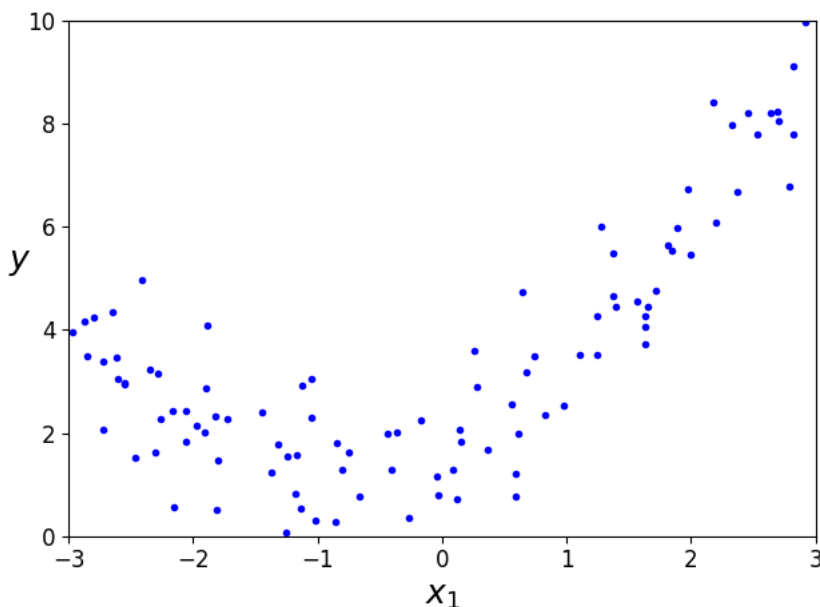
Saving figure gradient_descent_paths_plot



## Polynomial Regression

```
1 import numpy as np
2 import numpy.random as rnd
3
4 np.random.seed(42)
```

```
1 m = 100
2 X = 6 * np.random.rand(m, 1) - 3
3 y = 0.5 * X**2 + X + 2 + np.random.randn(m, 1)
```

```
1 plt.plot(X, y, "b.")
2 plt.xlabel("$x_1$", fontsize=18)
3 plt.ylabel("$y$", rotation=0, fontsize=18)
4 plt.axis([-3, 3, 0, 10])
5 save_fig("quadratic_data_plot")
6 plt.show()
```

Saving figure quadratic_data_plot



```
1 from sklearn.preprocessing import PolynomialFeatures
2 poly_features = PolynomialFeatures(degree=2, include_bias=False)
3 X_poly = poly_features.fit_transform(X)
4 X[0]
```

array([-0.75275929])

```
1 X_poly[0]
```

array([-0.75275929,  0.56664654])

```
1 lin_reg = LinearRegression()
2 lin_reg.fit(X_poly, y)
3 lin_reg.intercept_, lin_reg.coef_
```

```
(array([1.78134581]), array([[0.93366893, 0.56456263]]))
```
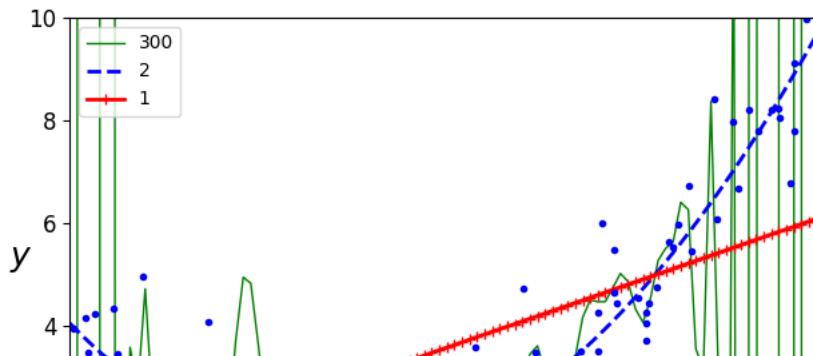
```
 1 X_new=np.linspace(-3, 3, 100).reshape(100, 1)
 2 X_new_poly = poly_features.transform(X_new)
 3 y_new = lin_reg.predict(X_new_poly)
 4 plt.plot(X, y, "b.")
 5 plt.plot(X_new, y_new, "r-", linewidth=2, label="Predictions")
 6 plt.xlabel("$x_1$", fontsize=18)
 7 plt.ylabel("$y$", rotation=0, fontsize=18)
 8 plt.legend(loc="upper left", fontsize=14)
 9 plt.axis([-3, 3, 0, 10])
10 save_fig("quadratic_predictions_plot")
11 plt.show()
```

    Saving figure quadratic_predictions_plot



```
 1 from sklearn.preprocessing import StandardScaler
 2 from sklearn.pipeline import Pipeline
 3
 4 for style, width, degree in (("g-", 1, 300), ("b--", 2, 2), ("r-+", 2, 1)):
 5     polybig_features = PolynomialFeatures(degree=degree, include_bias=False)
 6     std_scaler = StandardScaler()
 7     lin_reg = LinearRegression()
 8     polynomial_regression = Pipeline([
 9             ("poly_features", polybig_features),
10             ("std_scaler", std_scaler),
11             ("lin_reg", lin_reg),
12         ])
13     polynomial_regression.fit(X, y)
14     y_newbig = polynomial_regression.predict(X_new)
15     plt.plot(X_new, y_newbig, style, label=str(degree), linewidth=width)
16
17 plt.plot(X, y, "b.", linewidth=3)
18 plt.legend(loc="upper left")
19 plt.xlabel("$x_1$", fontsize=18)
20 plt.ylabel("$y$", rotation=0, fontsize=18)
21 plt.axis([-3, 3, 0, 10])
22 save_fig("high_degree_polynomials_plot")
23 plt.show()
```

Saving figure high_degree_polynomials_plot
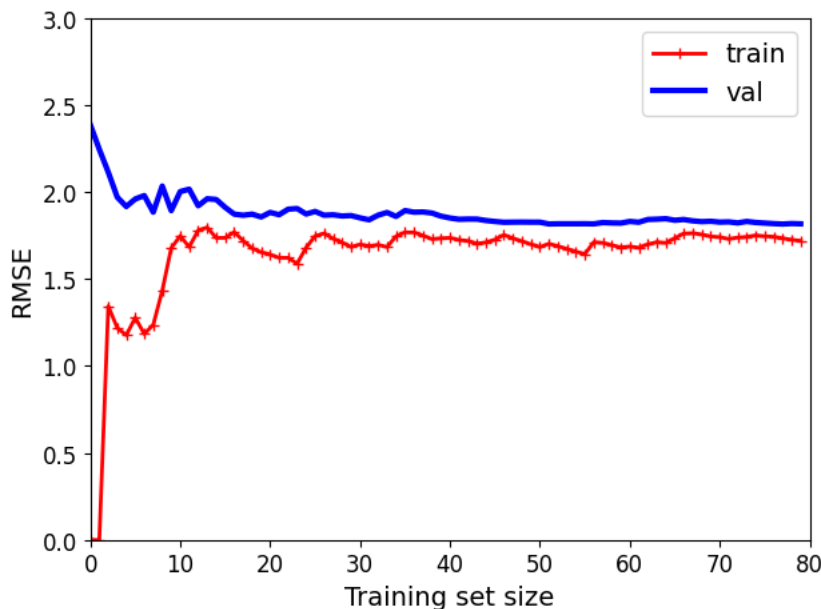


## Learning Curves

```
1 from sklearn.metrics import mean_squared_error
2 from sklearn.model_selection import train_test_split
3
4 def plot_learning_curves(model, X, y):
5     X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=10)
6     train_errors, val_errors = [], []
7     for m in range(1, len(X_train) + 1):
8         model.fit(X_train[:m], y_train[:m])
9         y_train_predict = model.predict(X_train[:m])
10        y_val_predict = model.predict(X_val)
11        train_errors.append(mean_squared_error(y_train[:m], y_train_predict))
12        val_errors.append(mean_squared_error(y_val, y_val_predict))
13
14    plt.plot(np.sqrt(train_errors), "r-+", linewidth=2, label="train")
15    plt.plot(np.sqrt(val_errors), "b-", linewidth=3, label="val")
16    plt.legend(loc="upper right", fontsize=14)    # not shown in the book
17    plt.xlabel("Training set size", fontsize=14) # not shown
18    plt.ylabel("RMSE", fontsize=14)              # not shown
```

```
1 lin_reg = LinearRegression()
2 plot_learning_curves(lin_reg, X, y)
3 plt.axis([0, 80, 0, 3])                    # not shown in the book
4 save_fig("underfitting_learning_curves_plot")  # not shown
5 plt.show()                                 # not shown
```

Saving figure underfitting_learning_curves_plot



```
1 from sklearn.pipeline import Pipeline
2
3 polynomial_regression = Pipeline([
```
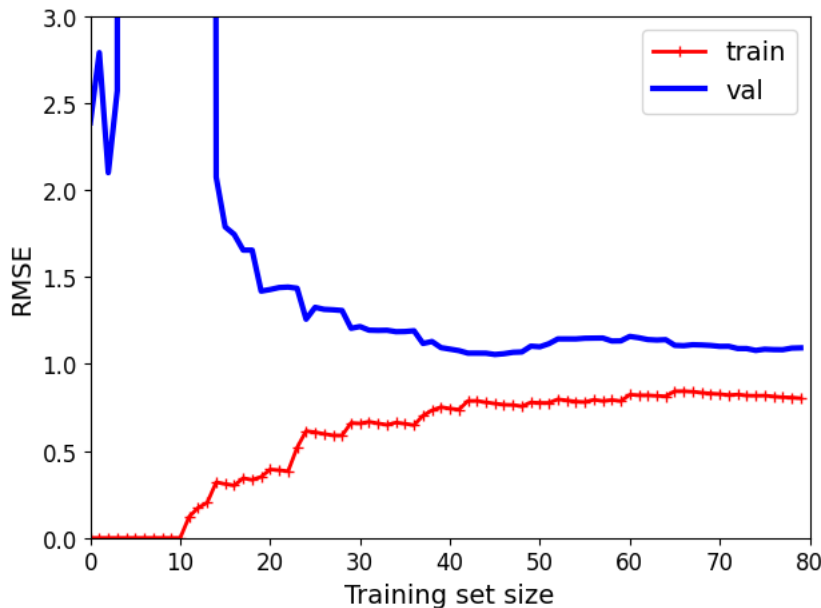
```
 4          ("poly_features", PolynomialFeatures(degree=10, include_bias=False)),
 5          ("lin_reg", LinearRegression()),
 6      ])
 7
 8 plot_learning_curves(polynomial_regression, X, y)
 9 plt.axis([0, 80, 0, 3])          # not shown
10 save_fig("learning_curves_plot")  # not shown
11 plt.show()                        # not shown
```

```
Saving figure learning_curves_plot
```



## Regularized Linear Models

## Ridge Regression

```
1 np.random.seed(42)
2 m = 20
3 X = 3 * np.random.rand(m, 1)
4 y = 1 + 0.5 * X + np.random.randn(m, 1) / 1.5
5 X_new = np.linspace(0, 3, 100).reshape(100, 1)
```

```
1 from sklearn.linear_model import Ridge
2 ridge_reg = Ridge(alpha=1, solver="cholesky", random_state=42)
3 ridge_reg.fit(X, y)
4 ridge_reg.predict([[1.5]])
```

```
    array([[1.55071465]])
```

```
1 ridge_reg = Ridge(alpha=1, solver="sag", random_state=42)
2 ridge_reg.fit(X, y)
3 ridge_reg.predict([[1.5]])
```

```
    array([[1.55072189]])
```

```
 1 from sklearn.linear_model import Ridge
 2
 3 def plot_model(model_class, polynomial, alphas, **model_kargs):
 4     for alpha, style in zip(alphas, ("b-", "g--", "r:")):
 5         model = model_class(alpha, **model_kargs) if alpha > 0 else LinearRegression()
 6         if polynomial:
 7             model = Pipeline([
 8                     ("poly_features", PolynomialFeatures(degree=10, include_bias=False)),
 9                     ("std_scaler", StandardScaler()),
10                     ("regul_reg", model),
11                 ])
```
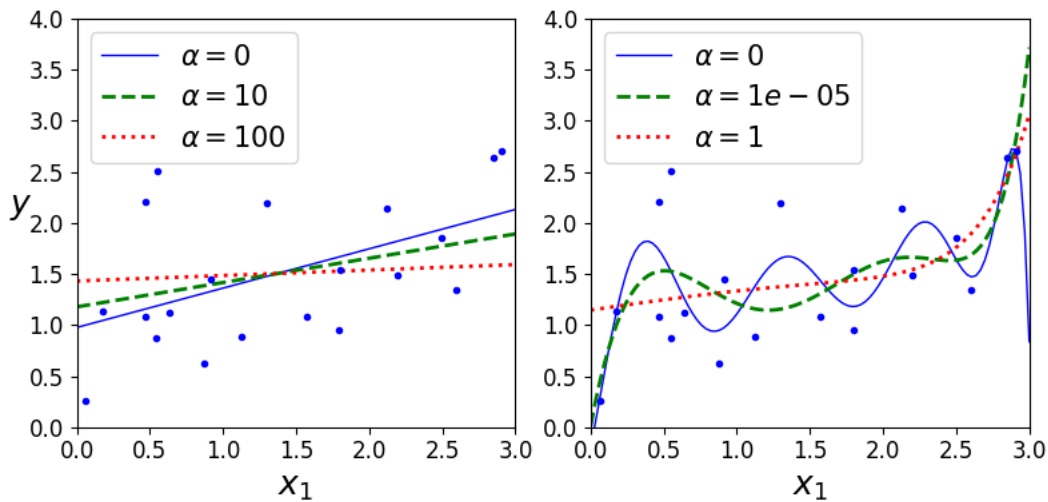
```
12        model.fit(X, y)
13        y_new_regul = model.predict(X_new)
14        lw = 2 if alpha > 0 else 1
15        plt.plot(X_new, y_new_regul, style, linewidth=lw, label=r"$\alpha = {}$".format(alpha))
16    plt.plot(X, y, "b.", linewidth=3)
17    plt.legend(loc="upper left", fontsize=15)
18    plt.xlabel("$x_1$", fontsize=18)
19    plt.axis([0, 3, 0, 4])
20
21 plt.figure(figsize=(8,4))
22 plt.subplot(121)
23 plot_model(Ridge, polynomial=False, alphas=(0, 10, 100), random_state=42)
24 plt.ylabel("$y$", rotation=0, fontsize=18)
25 plt.subplot(122)
26 plot_model(Ridge, polynomial=True, alphas=(0, 10**-5, 1), random_state=42)
27
28 save_fig("ridge_regression_plot")
29 plt.show()
```

    Saving figure ridge_regression_plot



**Note**: to be future-proof, we set `max_iter=1000` and `tol=1e-3` because these will be the default values in Scikit-Learn 0.21.

```
1 sgd_reg = SGDRegressor(penalty="l2", max_iter=1000, tol=1e-3, random_state=42)
2 sgd_reg.fit(X, y.ravel())
3 sgd_reg.predict([[1.5]])
```
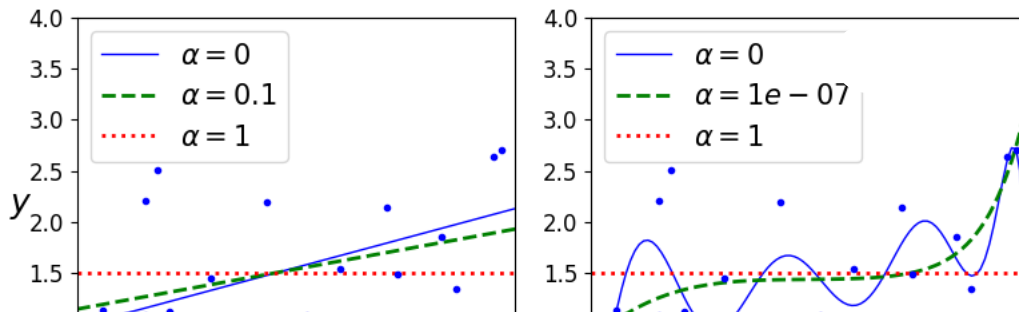
    array([1.47012588])

## ▾ Lasso Regression

```
1 from sklearn.linear_model import Lasso
2
3 plt.figure(figsize=(8,4))
4 plt.subplot(121)
5 plot_model(Lasso, polynomial=False, alphas=(0, 0.1, 1), random_state=42)
6 plt.ylabel("$y$", rotation=0, fontsize=18)
7 plt.subplot(122)
8 plot_model(Lasso, polynomial=True, alphas=(0, 10**-7, 1), random_state=42)
9
10 save_fig("lasso_regression_plot")
11 plt.show()
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_coordinate_descent.py:631: ConvergenceWarning: Objective did not converge.
  model = cd_fast.enet_coordinate_descent(
Saving figure lasso_regression_plot
```



```
1 from sklearn.linear_model import Lasso
2 lasso_reg = Lasso(alpha=0.1)
3 lasso_reg.fit(X, y)
4 lasso_reg.predict([[1.5]])
```

```
array([1.53788174])
```

## Elastic Net

```
1 from sklearn.linear_model import ElasticNet
2 elastic_net = ElasticNet(alpha=0.1, l1_ratio=0.5, random_state=42)
3 elastic_net.fit(X, y)
4 elastic_net.predict([[1.5]])
```

```
array([1.54333232])
```

## Early Stopping

```
1 np.random.seed(42)
2 m = 100
3 X = 6 * np.random.rand(m, 1) - 3
4 y = 2 + X + 0.5 * X**2 + np.random.randn(m, 1)
5
6 X_train, X_val, y_train, y_val = train_test_split(X[:50], y[:50].ravel(), test_size=0.5, random_state=10)
```

```
1 from copy import deepcopy
2
3 poly_scaler = Pipeline([
4         ("poly_features", PolynomialFeatures(degree=90, include_bias=False)),
5         ("std_scaler", StandardScaler())
6     ])
7
8 X_train_poly_scaled = poly_scaler.fit_transform(X_train)
9 X_val_poly_scaled = poly_scaler.transform(X_val)
10
11 sgd_reg = SGDRegressor(max_iter=1, tol=None, warm_start=True,
12                        penalty=None, learning_rate="constant", eta0=0.0005, random_state=42)
13
14 minimum_val_error = float("inf")
15 best_epoch = None
16 best_model = None
17 for epoch in range(1000):
18     # continues where it left off
19     sgd_reg.fit(X_train_poly_scaled, y_train)
20     y_val_predict = sgd_reg.predict(X_val_poly_scaled)
21     val_error = mean_squared_error(y_val, y_val_predict)
22     if val_error < minimum_val_error:
23         minimum_val_error = val_error
24         best_epoch = epoch
25         best_model = deepcopy(sgd_reg)
```
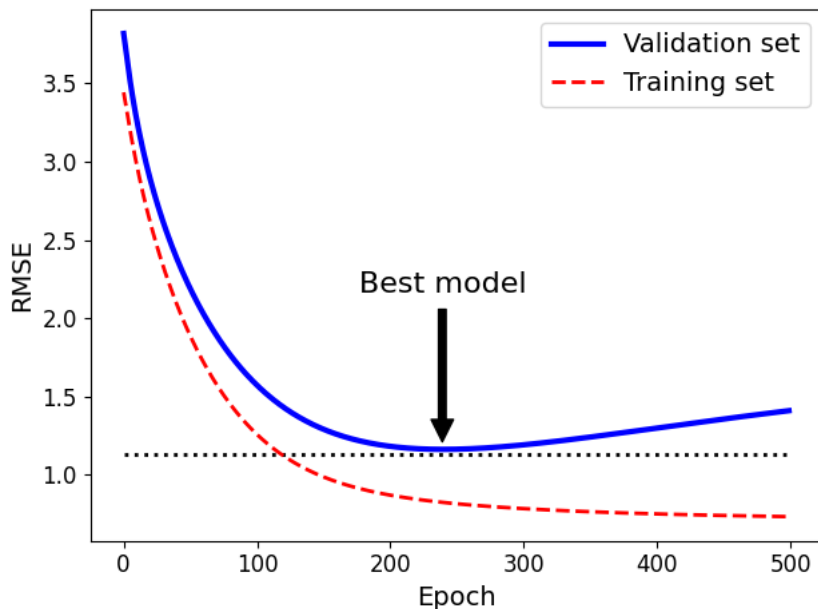
Create the graph:

```
1 sgd_reg = SGDRegressor(max_iter=1, tol=None, warm_start=True,
2                        penalty=None, learning_rate="constant", eta0=0.0005, random_state=42)
3
4 n_epochs = 500
5 train_errors, val_errors = [], []
6 for epoch in range(n_epochs):
7     sgd_reg.fit(X_train_poly_scaled, y_train)
8     y_train_predict = sgd_reg.predict(X_train_poly_scaled)
9     y_val_predict = sgd_reg.predict(X_val_poly_scaled)
10    train_errors.append(mean_squared_error(y_train, y_train_predict))
11    val_errors.append(mean_squared_error(y_val, y_val_predict))
12
13 best_epoch = np.argmin(val_errors)
14 best_val_rmse = np.sqrt(val_errors[best_epoch])
15
16 plt.annotate('Best model',
17              xy=(best_epoch, best_val_rmse),
18              xytext=(best_epoch, best_val_rmse + 1),
19              ha="center",
20              arrowprops=dict(facecolor='black', shrink=0.05),
21              fontsize=16,
22              )
23
24 best_val_rmse -= 0.03  # just to make the graph look better
25 plt.plot([0, n_epochs], [best_val_rmse, best_val_rmse], "k:", linewidth=2)
26 plt.plot(np.sqrt(val_errors), "b-", linewidth=3, label="Validation set")
27 plt.plot(np.sqrt(train_errors), "r--", linewidth=2, label="Training set")
28 plt.legend(loc="upper right", fontsize=14)
29 plt.xlabel("Epoch", fontsize=14)
30 plt.ylabel("RMSE", fontsize=14)
31 save_fig("early_stopping_plot")
32 plt.show()
```

Saving figure early_stopping_plot



```
1 best_epoch, best_model
```

```
(239,
 SGDRegressor(eta0=0.0005, learning_rate='constant', max_iter=1, penalty=None,
              random_state=42, tol=None, warm_start=True))
```

```
1 %matplotlib inline
2 import matplotlib.pyplot as plt
3 import numpy as np
```

```
1 t1a, t1b, t2a, t2b = -1, 3, -1.5, 1.5
2
3 t1s = np.linspace(t1a, t1b, 500)
4 t2s = np.linspace(t2a, t2b, 500)
5 t1, t2 = np.meshgrid(t1s, t2s)
6 T = np.c_[t1.ravel(), t2.ravel()]
```
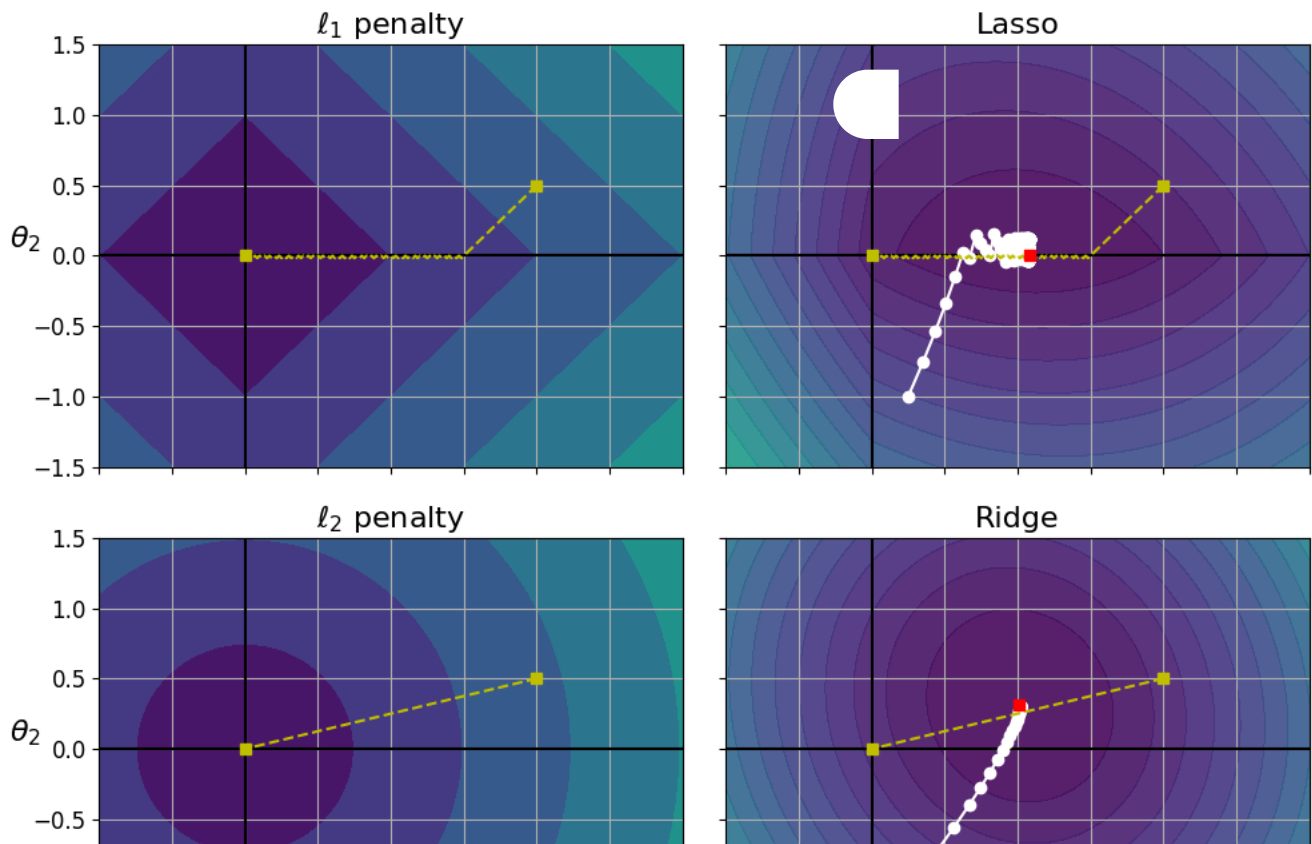
```
 7 Xr = np.array([[1, 1], [1, -1], [1, 0.5]])
 8 yr = 2 * Xr[:, :1] + 0.5 * Xr[:, 1:]
 9
10 J = (1/len(Xr) * np.sum((T.dot(Xr.T) - yr.T)**2, axis=1)).reshape(t1.shape)
11
12 N1 = np.linalg.norm(T, ord=1, axis=1).reshape(t1.shape)
13 N2 = np.linalg.norm(T, ord=2, axis=1).reshape(t1.shape)
14
15 t_min_idx = np.unravel_index(np.argmin(J), J.shape)
16 t1_min, t2_min = t1[t_min_idx], t2[t_min_idx]
17
18 t_init = np.array([[0.25], [-1]])
```

```
 1 def bgd_path(theta, X, y, l1, l2, core = 1, eta = 0.05, n_iterations = 200):
 2     path = [theta]
 3     for iteration in range(n_iterations):
 4         gradients = core * 2/len(X) * X.T.dot(X.dot(theta) - y) + l1 * np.sign(theta) + l2 * theta
 5         theta = theta - eta * gradients
 6         path.append(theta)
 7     return np.array(path)
 8
 9 fig, axes = plt.subplots(2, 2, sharex=True, sharey=True, figsize=(10.1, 8))
10 for i, N, l1, l2, title in ((0, N1, 2., 0, "Lasso"), (1, N2, 0,  2., "Ridge")):
11     JR = J + l1 * N1 + l2 * 0.5 * N2**2
12
13     tr_min_idx = np.unravel_index(np.argmin(JR), JR.shape)
14     t1r_min, t2r_min = t1[tr_min_idx], t2[tr_min_idx]
15
16     levelsJ=(np.exp(np.linspace(0, 1, 20)) - 1) * (np.max(J) - np.min(J)) + np.min(J)
17     levelsJR=(np.exp(np.linspace(0, 1, 20)) - 1) * (np.max(JR) - np.min(JR)) + np.min(JR)
18     levelsN=np.linspace(0, np.max(N), 10)
19
20     path_J = bgd_path(t_init, Xr, yr, l1=0, l2=0)
21     path_JR = bgd_path(t_init, Xr, yr, l1, l2)
22     path_N = bgd_path(np.array([[2.0], [0.5]]), Xr, yr, np.sign(l1)/3, np.sign(l2), core=0)
23
24     ax = axes[i, 0]
25     ax.grid(True)
26     ax.axhline(y=0, color='k')
27     ax.axvline(x=0, color='k')
28     ax.contourf(t1, t2, N / 2., levels=levelsN)
29     ax.plot(path_N[:, 0], path_N[:, 1], "y--")
30     ax.plot(0, 0, "ys")
31     ax.plot(t1_min, t2_min, "ys")
32     ax.set_title(r"$\ell_{}$ penalty".format(i + 1), fontsize=16)
33     ax.axis([t1a, t1b, t2a, t2b])
34     if i == 1:
35         ax.set_xlabel(r"$\theta_1$", fontsize=16)
36     ax.set_ylabel(r"$\theta_2$", fontsize=16, rotation=0)
37
38     ax = axes[i, 1]
39     ax.grid(True)
40     ax.axhline(y=0, color='k')
41     ax.axvline(x=0, color='k')
42     ax.contourf(t1, t2, JR, levels=levelsJR, alpha=0.9)
43     ax.plot(path_JR[:, 0], path_JR[:, 1], "w-o")
44     ax.plot(path_N[:, 0], path_N[:, 1], "y--")
45     ax.plot(0, 0, "ys")
46     ax.plot(t1_min, t2_min, "ys")
47     ax.plot(t1r_min, t2r_min, "rs")
48     ax.set_title(title, fontsize=16)
49     ax.axis([t1a, t1b, t2a, t2b])
50     if i == 1:
51         ax.set_xlabel(r"$\theta_1$", fontsize=16)
52
53 save_fig("lasso_vs_ridge_plot")
54 plt.show()
```

Saving figure lasso_vs_ridge_plot
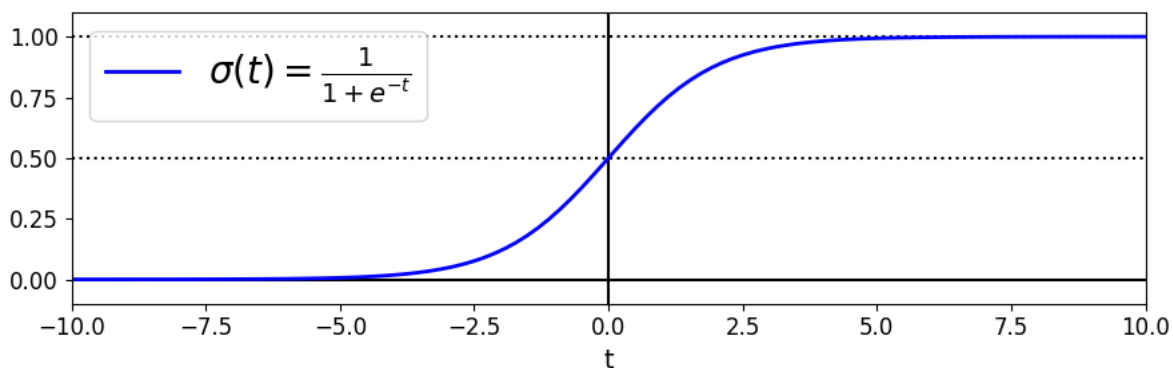


## Logistic Regression

## Decision Boundaries

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

Saving figure logistic_function_plot

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

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

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

```
**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.
    - Gates, G.W. (1972) "The Reduced Nearest Neighbor Rule".  IEEE Transactions
      on Information Theory, May 1972, 431-433.
    - See also: 1988 MLC Proceedings, 54-64.  Cheeseman et al"s AUTOCLASS II
      conceptual clustering system finds 3 classes in the data.
    - Many, many more ...
```

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

```
<ipython-input-56-c3494bf9af66>:2: DeprecationWarning: `np.int` is a deprecated alias for the builtin `int`. To silence this warning, us
Deprecated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
  y = (iris["target"] == 2).astype(np.int)  # 1 if Iris virginica, else 0
```
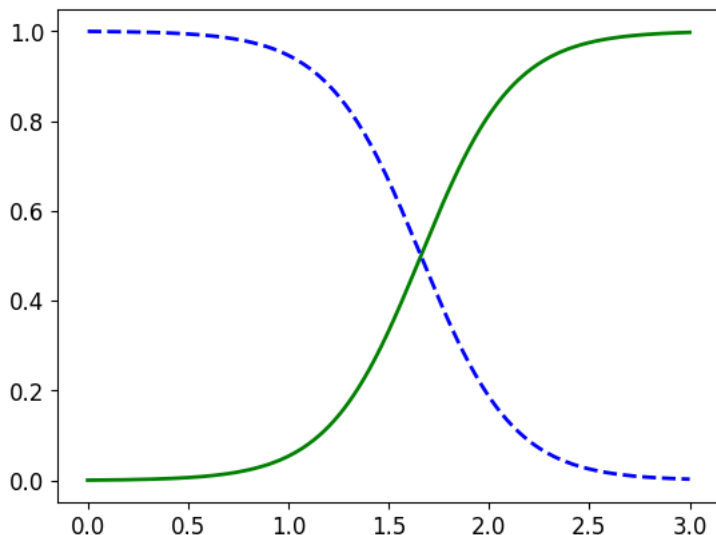
**Note**: To be future-proof we set `solver="lbfgs"` since this will be the default value in Scikit-Learn 0.22.

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

```
▾          LogisticRegression
LogisticRegression(random_state=42)
```

```
1 X_new = np.linspace(0, 3, 1000).reshape(-1, 1)
2 y_proba = log_reg.predict_proba(X_new)
3
4 plt.plot(X_new, y_proba[:, 1], "g-", linewidth=2, label="Iris virginica")
5 plt.plot(X_new, y_proba[:, 0], "b--", linewidth=2, label="Not Iris virginica")
```

```
[<matplotlib.lines.Line2D at 0x7bb350bdf730>]
```



The figure in the book actually is actually a bit fancier:

```
 1 X_new = np.linspace(0, 3, 1000).reshape(-1, 1)
 2 y_proba = log_reg.predict_proba(X_new)
 3 decision_boundary = X_new[y_proba[:, 1] >= 0.5][0]
 4
 5 plt.figure(figsize=(8, 3))
 6 plt.plot(X[y==0], y[y==0], "bs")
 7 plt.plot(X[y==1], y[y==1], "g^")
 8 plt.plot([decision_boundary, decision_boundary], [-1, 2], "k:", linewidth=2)
 9 plt.plot(X_new, y_proba[:, 1], "g-", linewidth=2, label="Iris virginica")
10 plt.plot(X_new, y_proba[:, 0], "b--", linewidth=2, label="Not Iris virginica")
11 plt.text(decision_boundary+0.02, 0.15, "Decision  boundary", fontsize=14, color="k", ha="center")
12 plt.arrow(decision_boundary, 0.08, -0.3, 0, head_width=0.05, head_length=0.1, fc='b', ec='b')
13 plt.arrow(decision_boundary, 0.92, 0.3, 0, head_width=0.05, head_length=0.1, fc='g', ec='g')
14 plt.xlabel("Petal width (cm)", fontsize=14)
15 plt.ylabel("Probability", fontsize=14)
16 plt.legend(loc="center left", fontsize=14)
17 plt.axis([0, 3, -0.02, 1.02])
18 save_fig("logistic_regression_plot")
19 plt.show()
```

```
/usr/local/lib/python3.10/dist-packages/matplotlib/patches.py:1475: VisibleDeprecationWarning: Creating an ndarray from ragged nested se
    self.verts = np.dot(coords, M) + [
Saving figure logistic_regression_plot
```



```
1 decision_boundary
```

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

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

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

## ▾ Softmax Regression

```
 1 from sklearn.linear_model import LogisticRegression
 2
 3 X = iris["data"][:, (2, 3)]  # petal length, petal width
 4 y = (iris["target"] == 2).astype(np.int)
 5
 6 log_reg = LogisticRegression(solver="lbfgs", C=10**10, random_state=42)
 7 log_reg.fit(X, y)
 8
 9 x0, x1 = np.meshgrid(
10         np.linspace(2.9, 7, 500).reshape(-1, 1),
11         np.linspace(0.8, 2.7, 200).reshape(-1, 1),
12     )
13 X_new = np.c_[x0.ravel(), x1.ravel()]
14
15 y_proba = log_reg.predict_proba(X_new)
16
17 plt.figure(figsize=(10, 4))
18 plt.plot(X[y==0, 0], X[y==0, 1], "bs")
19 plt.plot(X[y==1, 0], X[y==1, 1], "g^")
20
21 zz = y_proba[:, 1].reshape(x0.shape)
22 contour = plt.contour(x0, x1, zz, cmap=plt.cm.brg)
23
24
25 left_right = np.array([2.9, 7])
26 boundary = -(log_reg.coef_[0][0] * left_right + log_reg.intercept_[0]) / log_reg.coef_[0][1]
27
28 plt.clabel(contour, inline=1, fontsize=12)
29 plt.plot(left_right, boundary, "k--", linewidth=3)
30 plt.text(3.5, 1.5, "Not Iris virginica", fontsize=14, color="b", ha="center")
31 plt.text(6.5, 2.3, "Iris virginica", fontsize=14, color="g", ha="center")
32 plt.xlabel("Petal length", fontsize=14)
33 plt.ylabel("Petal width", fontsize=14)
34 plt.axis([2.9, 7, 0.8, 2.7])
35 save_fig("logistic_regression_contour_plot")
36 plt.show()
```

```
<ipython-input-62-1a12a15f0956>:4: DeprecationWarning: `np.int` is a deprecated alias for the builtin `int`. To silence this warning, us
Deprecated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
  y = (iris["target"] == 2).astype(np.int)
Saving figure logistic_regression_contour_plot
```



```
1 X = iris["data"][:, (2, 3)]  # petal length, petal width
2 y = iris["target"]
3
4 softmax_reg = LogisticRegression(multi_class="multinomial",solver="lbfgs", C=10, random_state=42)
5 softmax_reg.fit(X, y)
```

```
                    LogisticRegression
LogisticRegression(C=10, multi_class='multinomial', random_state=42)
```



```
 1 x0, x1 = np.meshgrid(
 2         np.linspace(0, 8, 500).reshape(-1, 1),
 3         np.linspace(0, 3.5, 200).reshape(-1, 1),
 4     )
 5 X_new = np.c_[x0.ravel(), x1.ravel()]
 6
 7
 8 y_proba = softmax_reg.predict_proba(X_new)
 9 y_predict = softmax_reg.predict(X_new)
10
11 zz1 = y_proba[:, 1].reshape(x0.shape)
12 zz = y_predict.reshape(x0.shape)
13
14 plt.figure(figsize=(10, 4))
15 plt.plot(X[y==2, 0], X[y==2, 1], "g^", label="Iris virginica")
16 plt.plot(X[y==1, 0], X[y==1, 1], "bs", label="Iris versicolor")
17 plt.plot(X[y==0, 0], X[y==0, 1], "yo", label="Iris setosa")
18
19 from matplotlib.colors import ListedColormap
20 custom_cmap = ListedColormap(['#fafab0','#9898ff','#a0faa0'])
21
22 plt.contourf(x0, x1, zz, cmap=custom_cmap)
23 contour = plt.contour(x0, x1, zz1, cmap=plt.cm.brg)
24 plt.clabel(contour, inline=1, fontsize=12)
25 plt.xlabel("Petal length", fontsize=14)
26 plt.ylabel("Petal width", fontsize=14)
27 plt.legend(loc="center left", fontsize=14)
28 plt.axis([0, 7, 0, 3.5])
29 save_fig("softmax_regression_contour_plot")
30 plt.show()
```

```
Saving figure softmax_regression_contour_plot
```



```
1 softmax_reg.predict([[5, 2]])
```

```
array([2])
```

```
1 softmax_reg.predict_proba([[5, 2]])
```

```
    array([[6.38014896e-07, 5.74929995e-02, 9.42506362e-01]])
```

## ▾ Exercise solutions

## ▾ 1. to 11.

See appendix A.

## ▾ 12. Batch Gradient Descent with early stopping for Softmax Regression

(without using Scikit-Learn)

Let's start by loading the data. We will just reuse the Iris dataset we loaded earlier.

```
1 X = iris["data"][:, (2, 3)]  # petal length, petal width
2 y = iris["target"]
```

We need to add the bias term for every instance ($x_0 = 1$):

```
1 X_with_bias = np.c_[np.ones([len(X), 1]), X]
```

And let's set the random seed so the output of this exercise solution is reproducible:

```
1 np.random.seed(2042)
```

The easiest option to split the dataset into a training set, a validation set and a test set would be to use Scikit-Learn's `train_test_split()` function, but the point of this exercise is to try understand the algorithms by implementing them manually. So here is one possible implementation:

```
 1 test_ratio = 0.2
 2 validation_ratio = 0.2
 3 total_size = len(X_with_bias)
 4
 5 test_size = int(total_size * test_ratio)
 6 validation_size = int(total_size * validation_ratio)
 7 train_size = total_size - test_size - validation_size
 8
 9 rnd_indices = np.random.permutation(total_size)
10
11 X_train = X_with_bias[rnd_indices[:train_size]]
12 y_train = y[rnd_indices[:train_size]]
13 X_valid = X_with_bias[rnd_indices[train_size:-test_size]]
14 y_valid = y[rnd_indices[train_size:-test_size]]
15 X_test = X_with_bias[rnd_indices[-test_size:]]
16 y_test = y[rnd_indices[-test_size:]]
```

The targets are currently class indices (0, 1 or 2), but we need target class probabilities to train the Softmax Regression model. Each instance will have target class probabilities equal to 0.0 for all classes except for the target class which will have a probability of 1.0 (in other words, the vector of class probabilities for ay given instance is a one-hot vector). Let's write a small function to convert the vector of class indices into a matrix containing a one-hot vector for each instance:

```
1 def to_one_hot(y):
2     n_classes = y.max() + 1
3     m = len(y)
4     Y_one_hot = np.zeros((m, n_classes))
5     Y_one_hot[np.arange(m), y] = 1
6     return Y_one_hot
```

Let's test this function on the first 10 instances:

```
1 y_train[:10]
```

```
array([0, 1, 2, 1, 1, 0, 1, 1, 1, 0])
```

```
1 to_one_hot(y_train[:10])
```

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

Looks good, so let's create the target class probabilities matrix for the training set and the test set:

```
1 Y_train_one_hot = to_one_hot(y_train)
2 Y_valid_one_hot = to_one_hot(y_valid)
3 Y_test_one_hot = to_one_hot(y_test)
```

Now let's implement the Softmax function. Recall that it is defined by the following equation:

$$\sigma(\mathbf{s}(\mathbf{x}))_k = \frac{\exp(s_k(\mathbf{x}))}{\sum_{j=1}^{K} \exp(s_j(\mathbf{x}))}$$

```
1 def softmax(logits):
2     exps = np.exp(logits)
3     exp_sums = np.sum(exps, axis=1, keepdims=True)
4     return exps / exp_sums
```

We are almost ready to start training. Let's define the number of inputs and outputs:

```
1 n_inputs = X_train.shape[1] # == 3 (2 features plus the bias term)
2 n_outputs = len(np.unique(y_train))   # == 3 (3 iris classes)
```

Now here comes the hardest part: training! Theoretically, it's simple: it's just a matter of translating the math equations into Python code. But in practice, it can be quite tricky: in particular, it's easy to mix up the order of the terms, or the indices. You can even end up with code that looks like it's working but is actually not computing exactly the right thing. When unsure, you should write down the shape of each term in the equation and make sure the corresponding terms in your code match closely. It can also help to evaluate each term independently and print them out. The good news it that you won't have to do this everyday, since all this is well implemented by Scikit-Learn, but it will help you understand what's going on under the hood.

So the equations we will need are the cost function:

$$J(\mathbf{\Theta}) = -\frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{K} y_k^{(i)} \log\left(\hat{p}_k^{(i)}\right)$$

And the equation for the gradients:

$$\nabla_{\theta^{(k)}} J(\mathbf{\Theta}) = \frac{1}{m} \sum_{i=1}^{m} \left(\hat{p}_k^{(i)} - y_k^{(i)}\right) \mathbf{x}^{(i)}$$

Note that $\log\left(\hat{p}_k^{(i)}\right)$ may not be computable if $\hat{p}_k^{(i)} = 0$. So we will add a tiny value $\epsilon$ to $\log\left(\hat{p}_k^{(i)}\right)$ to avoid getting `nan` values.

```
1 eta = 0.01
2 n_iterations = 5001
3 m = len(X_train)
4 epsilon = 1e-7
5
6 Theta = np.random.randn(n_inputs, n_outputs)
7
```

```
 8 for iteration in range(n_iterations):
 9     logits = X_train.dot(Theta)
10     Y_proba = softmax(logits)
11     if iteration % 500 == 0:
12         loss = -np.mean(np.sum(Y_train_one_hot * np.log(Y_proba + epsilon), axis=1))
13         print(iteration, loss)
14     error = Y_proba - Y_train_one_hot
15     gradients = 1/m * X_train.T.dot(error)
16     Theta = Theta - eta * gradients
```

```
0 5.446205811872683
500 0.8350062641405651
1000 0.6878801447192402
1500 0.6012379137693314
2000 0.5444496861981872
2500 0.5038530181431525
3000 0.4729228972192248
3500 0.44824244188957774
4000 0.4278651093928793
4500 0.41060071429187134
5000 0.3956780375390374
```

And that's it! The Softmax model is trained. Let's look at the model parameters:

```
1 Theta
```

```
array([[ 3.32094157, -0.6501102 , -2.99979416],
       [-1.1718465 ,  0.11706172,  0.10507543],
       [-0.70224261, -0.09527802,  1.4786383 ]])
```

Let's make predictions for the validation set and check the accuracy score:

```
1 logits = X_valid.dot(Theta)
2 Y_proba = softmax(logits)
3 y_predict = np.argmax(Y_proba, axis=1)
4
5 accuracy_score = np.mean(y_predict == y_valid)
6 accuracy_score
```
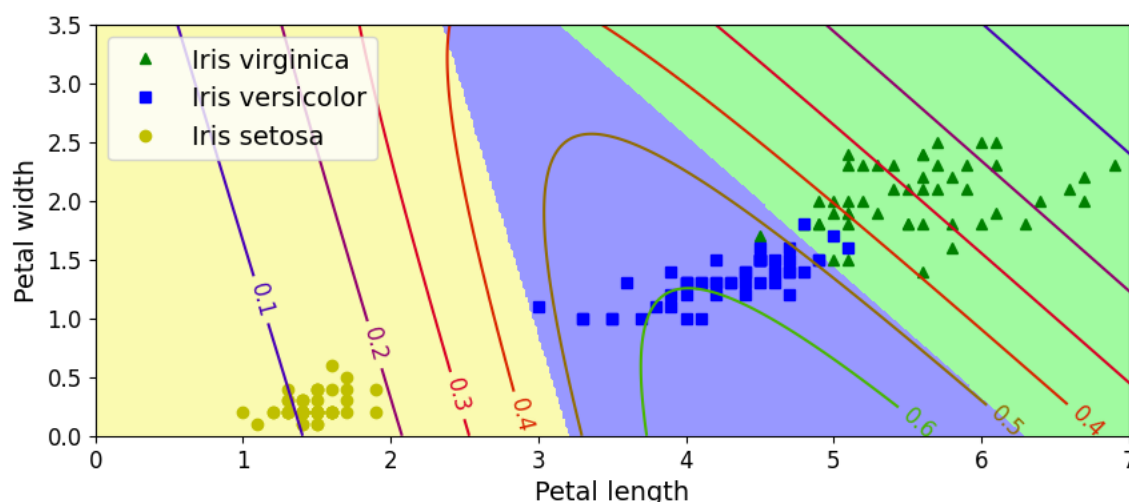
```
0.9666666666666667
```

Well, this model looks pretty good. For the sake of the exercise, let's add a bit of $\ell_2$ regularization. The following training code is similar to the one above, but the loss now has an additional $\ell_2$ penalty, and the gradients have the proper additional term (note that we don't regularize the first element of `Theta` since this corresponds to the bias term). Also, let's try increasing the learning rate `eta`.

```
 1 eta = 0.1
 2 n_iterations = 5001
 3 m = len(X_train)
 4 epsilon = 1e-7
 5 alpha = 0.1  # regularization hyperparameter
 6
 7 Theta = np.random.randn(n_inputs, n_outputs)
 8
 9 for iteration in range(n_iterations):
10     logits = X_train.dot(Theta)
11     Y_proba = softmax(logits)
12     if iteration % 500 == 0:
13         xentropy_loss = -np.mean(np.sum(Y_train_one_hot * np.log(Y_proba + epsilon), axis=1))
14         l2_loss = 1/2 * np.sum(np.square(Theta[1:]))
15         loss = xentropy_loss + alpha * l2_loss
16         print(iteration, loss)
17     error = Y_proba - Y_train_one_hot
18     gradients = 1/m * X_train.T.dot(error) + np.r_[np.zeros([1, n_outputs]), alpha * Theta[1:]]
19     Theta = Theta - eta * gradients
```

```
0 6.629842469083912
500 0.5339667976629505
1000 0.503640075014894
1500 0.4946891059460322
2000 0.4912968418075477
2500 0.4898924700933296
3000 0.4892990598451198
3500 0.48903512443978603
4000 0.4889173621830818
```

```
4500 0.4888643337449303
5000 0.4888403120738818
```

Because of the additional $\ell_2$ penalty, the loss seems greater than earlier, but perhaps this model will perform better? Let's find out:

```
1 logits = X_valid.dot(Theta)
2 Y_proba = softmax(logits)
3 y_predict = np.argmax(Y_proba, axis=1)
4
5 accuracy_score = np.mean(y_predict == y_valid)
6 accuracy_score
```

```
1.0
```

Cool, perfect accuracy! We probably just got lucky with this validation set, but still, it's pleasant.

Now let's add early stopping. For this we just need to measure the loss on the validation set at every iteration and stop when the error starts growing.

```
 1 eta = 0.1
 2 n_iterations = 5001
 3 m = len(X_train)
 4 epsilon = 1e-7
 5 alpha = 0.1  # regularization hyperparameter
 6 best_loss = np.infty
 7
 8 Theta = np.random.randn(n_inputs, n_outputs)
 9
10 for iteration in range(n_iterations):
11     logits = X_train.dot(Theta)
12     Y_proba = softmax(logits)
13     error = Y_proba - Y_train_one_hot
14     gradients = 1/m * X_train.T.dot(error) + np.r_[np.zeros([1, n_outputs]), alpha * Theta[1:]]
15     Theta = Theta - eta * gradients
16
17     logits = X_valid.dot(Theta)
18     Y_proba = softmax(logits)
19     xentropy_loss = -np.mean(np.sum(Y_valid_one_hot * np.log(Y_proba + epsilon), axis=1))
20     l2_loss = 1/2 * np.sum(np.square(Theta[1:]))
21     loss = xentropy_loss + alpha * l2_loss
22     if iteration % 500 == 0:
23         print(iteration, loss)
24     if loss < best_loss:
25         best_loss = loss
26     else:
27         print(iteration - 1, best_loss)
28         print(iteration, loss, "early stopping!")
29         break
```

```
0 4.7096017363419875
500 0.5739711987633519
1000 0.5435638529109128
1500 0.5355752782580261
2000 0.5331959249285544
2500 0.5325946767399383
2765 0.5325460966791898
2766 0.5325460971327977 early stopping!
```

```
1 logits = X_valid.dot(Theta)
2 Y_proba = softmax(logits)
3 y_predict = np.argmax(Y_proba, axis=1)
4
5 accuracy_score = np.mean(y_predict == y_valid)
6 accuracy_score
```

```
1.0
```

Still perfect, but faster.

Now let's plot the model's predictions on the whole dataset:

```
 1 x0, x1 = np.meshgrid(
 2         np.linspace(0, 8, 500).reshape(-1, 1),
 3         np.linspace(0, 3.5, 200).reshape(-1, 1),
 4     )
 5 X_new = np.c_[x0.ravel(), x1.ravel()]
 6 X_new_with_bias = np.c_[np.ones([len(X_new), 1]), X_new]
 7
 8 logits = X_new_with_bias.dot(Theta)
 9 Y_proba = softmax(logits)
10 y_predict = np.argmax(Y_proba, axis=1)
11
12 zz1 = Y_proba[:, 1].reshape(x0.shape)
13 zz = y_predict.reshape(x0.shape)
14
15 plt.figure(figsize=(10, 4))
16 plt.plot(X[y==2, 0], X[y==2, 1], "g^", label="Iris virginica")
17 plt.plot(X[y==1, 0], X[y==1, 1], "bs", label="Iris versicolor")
18 plt.plot(X[y==0, 0], X[y==0, 1], "yo", label="Iris setosa")
19
20 from matplotlib.colors import ListedColormap
21 custom_cmap = ListedColormap(['#fafab0','#9898ff','#a0faa0'])
22
23 plt.contourf(x0, x1, zz, cmap=custom_cmap)
24 contour = plt.contour(x0, x1, zz1, cmap=plt.cm.brg)
25 plt.clabel(contour, inline=1, fontsize=12)
26 plt.xlabel("Petal length", fontsize=14)
27 plt.ylabel("Petal width", fontsize=14)
28 plt.legend(loc="upper left", fontsize=14)
29 plt.axis([0, 7, 0, 3.5])
30 plt.show()
```



And now let's measure the final model's accuracy on the test set:

```
1 logits = X_test.dot(Theta)
2 Y_proba = softmax(logits)
3 y_predict = np.argmax(Y_proba, axis=1)
4
5 accuracy_score = np.mean(y_predict == y_test)
6 accuracy_score
```

```
0.9333333333333333
```

Our perfect model turns out to have slight imperfections. This variability is likely due to the very small size of the dataset: depending on how you sample the training set, validation set and the test set, you can get quite different results. Try changing the random seed and running the code again a few times, you will see that the results will vary.

T͟T   B   I   <>   ⌐⊃   🖾   ⫧≣   ≣   ≣   •••   Ψ   ☺   ▭

ANSWERS TO THE QUESTIONS:

5. A. What is the name of the Scikit Learn estimator class that allows create polynomial features? What is defined by the parameter degree=2?

ANSWERS TO THE QUESTIONS:

The Scikit Learn estimator class that allows you to create polynomial
is called PolynomialFeatures. The degree parameter defines the highest
extend of the polynomial features.

B. On the output on line 36 or of Geron's Ch. 4 notebook, why is the t
set RMSE is equal to zero at the training set size equal to 1 or 2, an
the same training set size the validation set RMSE is very high? (2 po
This happens due to the nature of the learning process in machine lear
models and the concept of overfitting. The training set size is 1 or 2
the model is very simple and can perfectly fit the training data. This
why the training RMSE is zero. However, this model is likely overfitti
data, meaning it has learned the noise in the training data along with
underlying pattern. As a result, when predicting the validation set, t
makes large errors because it has learned patterns that do not general
to new data.

C. After you fit a model using polynomial feature with the degree=20 y
that your model is overfitting. How should you change the degree hyper
and why?
If a model is overfitting when using a polynomial feature with a degre
it means that the model is too complex and is fitting the training dat
well, including its noise and outliers. This is why the model performs
the training data but poorly on the validation data. To address this i
should decrease the degree of the polynomial feature. The degree of th
polynomial feature determines the complexity of the model. A higher de
means a more complex model, which can lead to overfitting.
This can be done by decreasing the degree from 20 to a lower value, su
or 5. You should then evaluate the performance of the model on the val
data to see if the overfitting issue is resolved. If the model still o
you can try decreasing the degree further. The optimal degree is the o
gives the best performance on the validation data without overfitting.
be determined through techniques such as cross-validation or by visual
inspecting the learning curves.

D. Why would you use ridge regression, lasso regression or elastic net
point)
Ridge Regression, Lasso Regression, and Elastic Net are all types of
regularization techniques used in machine learning to prevent overfitt
Ridge regression is a type of linear regression that includes a regula
term which is the sum of the squares of the coefficients, which adds a
for large coefficients. It helps prevent overfitting by discouraging t
from fitting the training data too closely.
Lasso regression also includes a regularization term. However, instead
sum of the squares of the coefficients, the Lasso penalty term is the
value of the coefficients. This can lead to some coefficients being ex
zero, effectively excluding those features from the model. This can be
when you want to perform feature selection, as it can help identify ir
features.
Elastic Net is a hybrid of Ridge and Lasso regression.  It includes bo
sum of the squares of the coefficients and the absolute value of the
coefficients. This can help prevent overfitting while also performing
selection.

E. After you fit a ridge regression with a hyperparameter alpha equal
you see that your model is overfitting. What should you do?
To deal with this you should increase the value of the alpha hyperpara
this scenario you increase the alpha value from 0.001 to a higher valu
as 0.01 or 0.1. After that evaluate the performance of the model on th
validation data to see if the overfitting issue is resolved.

0.01 or 0.1. After that evaluate the performance of the model on the validation data to see if the overfitting issue is resolved.