

# DMML2024-Lecture08-01feb2024-Regression

February 6, 2024

## Chapter 4 – Training Models

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

### 1 Setup

This project requires Python 3.7 or above:

```
[1]: import sys

assert sys.version_info >= (3, 7)
```

It also requires Scikit-Learn 1.0.1:

```
[2]: from packaging import version
import sklearn

assert version.parse(sklearn.__version__) >= version.parse("1.0.1")
```

As we did in previous chapters, let's define the default font sizes to make the figures prettier:

```
[3]: import matplotlib.pyplot as plt

plt.rc('font', size=14)
plt.rc('axes', labelsz=14, titlesz=14)
plt.rc('legend', fontsize=14)
plt.rc('xtick', labelsz=10)
plt.rc('ytick', labelsz=10)
```

And let's create the `images/training_linear_models` folder (if it doesn't already exist), and define the `save_fig()` function which is used through this notebook to save the figures in high-res for the book:

```
[4]: from pathlib import Path

IMAGES_PATH = Path() / "images" / "training_linear_models"
IMAGES_PATH.mkdir(parents=True, exist_ok=True)

def save_fig(fig_id, tight_layout=True, fig_extension="png", resolution=300):
    path = IMAGES_PATH / f"{fig_id}.{fig_extension}"
```

```
if tight_layout:
    plt.tight_layout()
plt.savefig(path, format=fig_extension, dpi=resolution)
```

## 2 Linear Regression

### 2.1 The Normal Equation

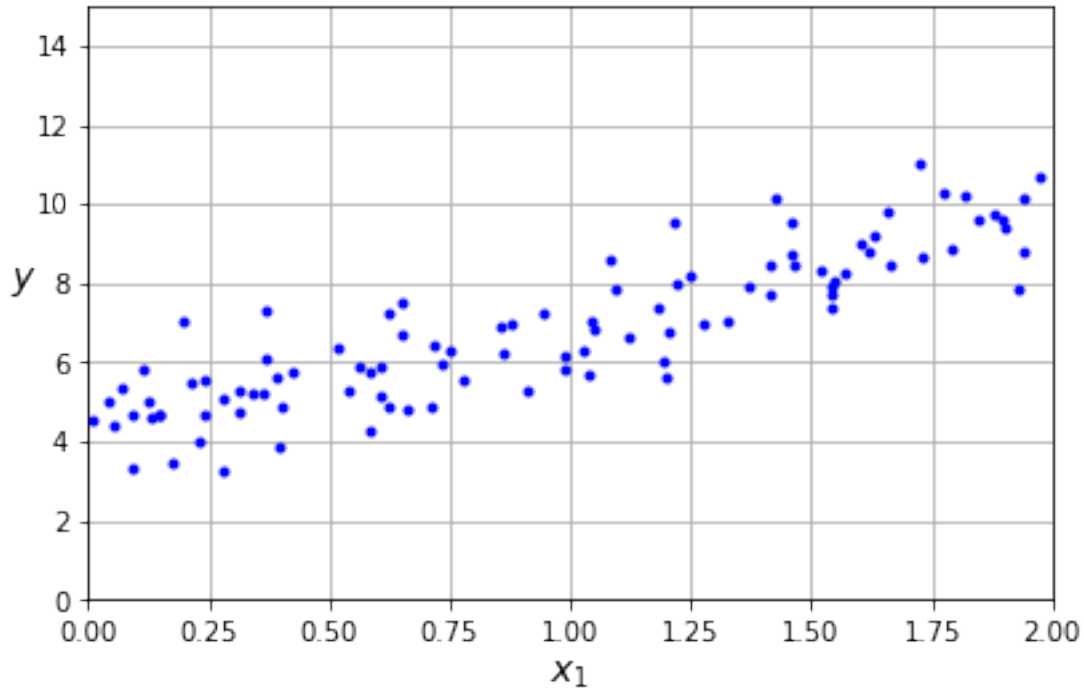
```
[5]: import numpy as np

np.random.seed(42) # to make this code example reproducible
m = 100 # number of instances
X = 2 * np.random.rand(m, 1) # column vector
y = 4 + 3 * X + np.random.randn(m, 1) # column vector
```

```
[6]: # extra code - generates and saves Figure 4-1
```

```
import matplotlib.pyplot as plt

plt.figure(figsize=(6, 4))
plt.plot(X, y, "b.")
plt.xlabel("$x_1$")
plt.ylabel("$y$", rotation=0)
plt.axis([0, 2, 0, 15])
plt.grid()
save_fig("generated_data_plot")
plt.show()
```



```
[7]: from sklearn.preprocessing import add_dummy_feature
```

```
X_b = add_dummy_feature(X) # add x0 = 1 to each instance
theta_best = np.linalg.inv(X_b.T @ X_b) @ X_b.T @ y
```

```
[8]: theta_best
```

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

```
[9]: X_new = np.array([[0], [2]])
X_new_b = add_dummy_feature(X_new) # add x0 = 1 to each instance
y_predict = X_new_b @ theta_best
y_predict
```

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

```
[10]: import matplotlib.pyplot as plt
```

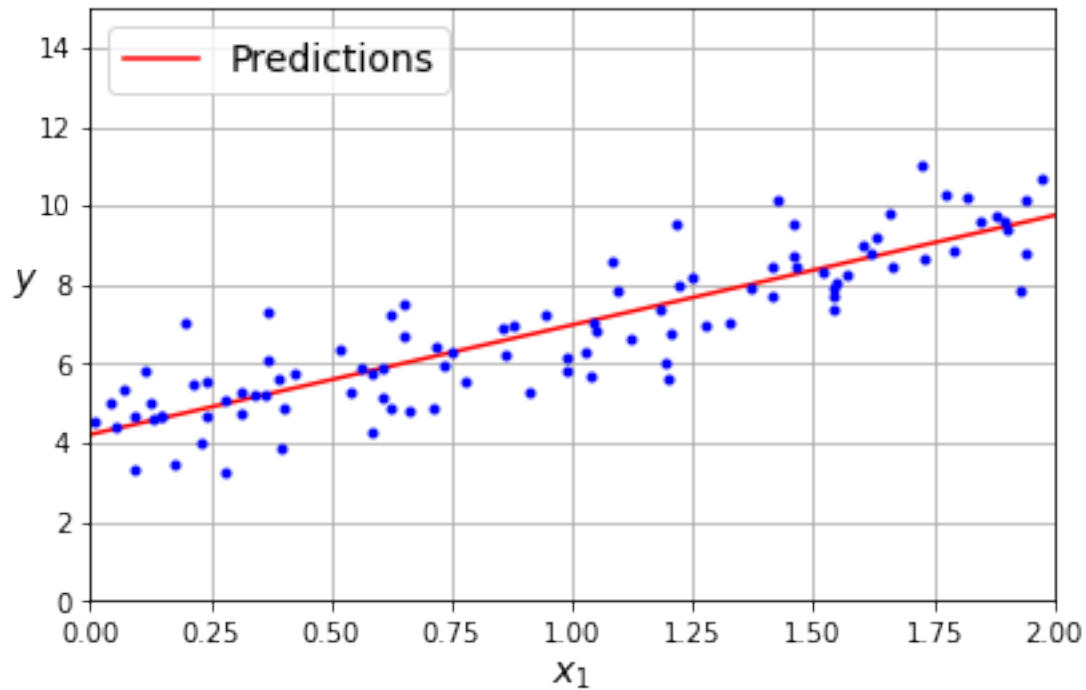
```
plt.figure(figsize=(6, 4)) # extra code - not needed, just formatting
plt.plot(X_new, y_predict, "r-", label="Predictions")
plt.plot(X, y, "b.")
```

```

# extra code - beautifies and saves Figure 4-2
plt.xlabel("$x_1$")
plt.ylabel("$y$", rotation=0)
plt.axis([0, 2, 0, 15])
plt.grid()
plt.legend(loc="upper left")
save_fig("linear_model_predictions_plot")

plt.show()

```



```
[11]: from sklearn.linear_model import LinearRegression
```

```

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

```

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

```
[12]: lin_reg.predict(X_new)
```

```
[12]: 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:

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

```
[13]: 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:

```
[14]: np.linalg.pinv(X_b) @ y
```

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

## 3 Gradient Descent

### 3.1 Batch Gradient Descent

```
[15]: eta = 0.1 # learning rate
      n_epochs = 1000
      m = len(X_b) # number of instances

      np.random.seed(42)
      theta = np.random.randn(2, 1) # randomly initialized model parameters

      for epoch in range(n_epochs):
          gradients = 2 / m * X_b.T @ (X_b @ theta - y)
          theta = theta - eta * gradients
```

The trained model parameters:

```
[16]: theta
```

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

```
[17]: # extra code - generates and saves Figure 4-8
```

```
import matplotlib as mpl

def plot_gradient_descent(theta, eta):
    m = len(X_b)
    plt.plot(X, y, "b.")
    n_epochs = 1000
    n_shown = 20
    theta_path = []
    for epoch in range(n_epochs):
        if epoch < n_shown:
```

```

        y_predict = X_new_b @ theta
        color = mpl.colors.rgb2hex(plt.cm.OrRd(epoch / n_shown + 0.15))
        plt.plot(X_new, y_predict, linestyle="solid", color=color)
        gradients = 2 / m * X_b.T @ (X_b @ theta - y)
        theta = theta - eta * gradients
        theta_path.append(theta)
plt.xlabel("$x_1$")
plt.axis([0, 2, 0, 15])
plt.grid()
plt.title(fr"$\eta = {eta}$")
return theta_path

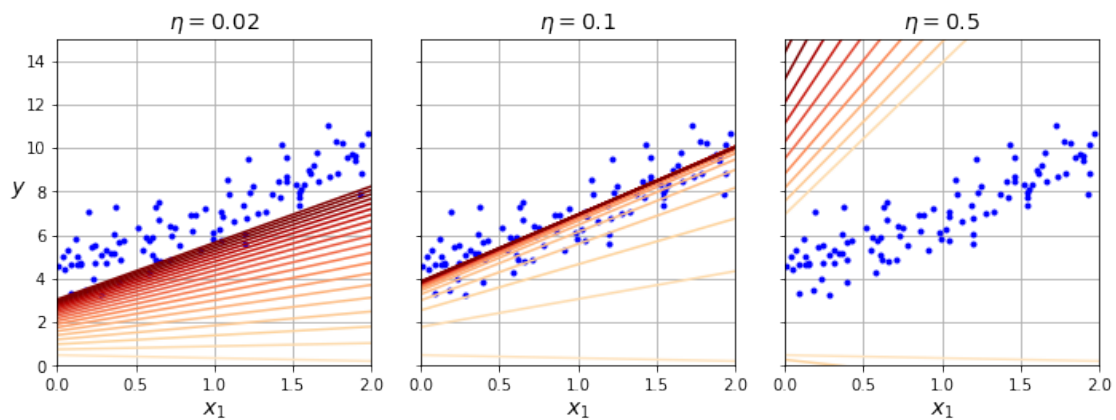
```

```

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)
plt.subplot(132)
theta_path_bgd = plot_gradient_descent(theta, eta=0.1)
plt.gca().axes.yaxis.set_ticklabels([])
plt.subplot(133)
plt.gca().axes.yaxis.set_ticklabels([])
plot_gradient_descent(theta, eta=0.5)
save_fig("gradient_descent_plot")
plt.show()

```



## 3.2 Stochastic Gradient Descent

```
[18]: theta_path_sgd = [] # extra code - we need to store the path of theta in the
      #                  #                  parameter space to plot the next figure

[19]: n_epochs = 50
      t0, t1 = 5, 50 # learning schedule hyperparameters

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

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

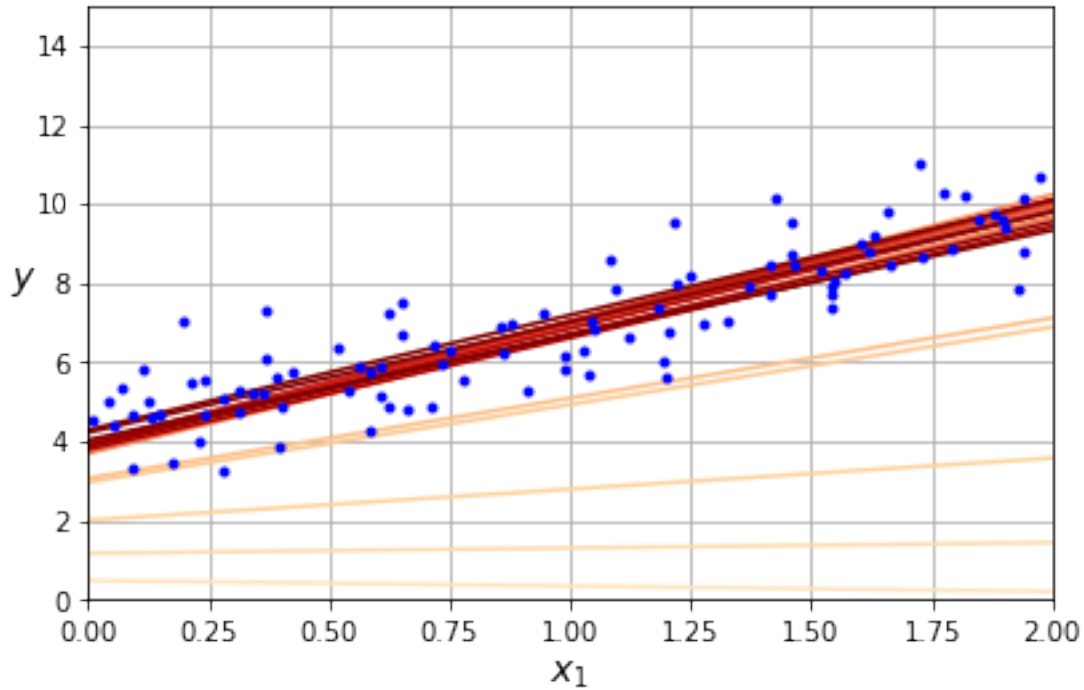
      n_shown = 20 # extra code - just needed to generate the figure below
      plt.figure(figsize=(6, 4)) # extra code - not needed, just formatting

      for epoch in range(n_epochs):
          for iteration in range(m):

              # extra code - these 4 lines are used to generate the figure
              if epoch == 0 and iteration < n_shown:
                  y_predict = X_new_b @ theta
                  color = mpl.colors.rgb2hex(plt.cm.OrRd(iteration / n_shown + 0.15))
                  plt.plot(X_new, y_predict, color=color)

              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 @ (xi @ theta - yi) # for SGD, do not divide by m
              eta = learning_schedule(epoch * m + iteration)
              theta = theta - eta * gradients
              theta_path_sgd.append(theta) # extra code - to generate the figure

      # extra code - this section beautifies and saves Figure 4-10
      plt.plot(X, y, "b.")
      plt.xlabel("$x_1$")
      plt.ylabel("$y$", rotation=0)
      plt.axis([0, 2, 0, 15])
      plt.grid()
      save_fig("sgd_plot")
      plt.show()
```



```
[20]: theta
```

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

```
[21]: from sklearn.linear_model import SGDRegressor
```

```
sgd_reg = SGDRegressor(max_iter=1000, tol=1e-5, penalty=None, eta0=0.01,
                       n_iter_no_change=100, random_state=42)
sgd_reg.fit(X, y.ravel()) # y.ravel() because fit() expects 1D targets
```

```
[21]: SGDRegressor(n_iter_no_change=100, penalty=None, random_state=42, tol=1e-05)
```

```
[22]: sgd_reg.intercept_, sgd_reg.coef_
```

```
[22]: (array([4.21278812]), array([2.77270267]))
```

### 3.3 Mini-batch gradient descent

The code in this section is used to generate the next figure, it is not in the book.

```
[23]: # extra code - this cell generates and saves Figure 4-11
```

```
from math import ceil
```



```

n_epochs = 50
minibatch_size = 20
n_batches_per_epoch = ceil(m / minibatch_size)

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

t0, t1 = 200, 1000 # learning schedule hyperparameters

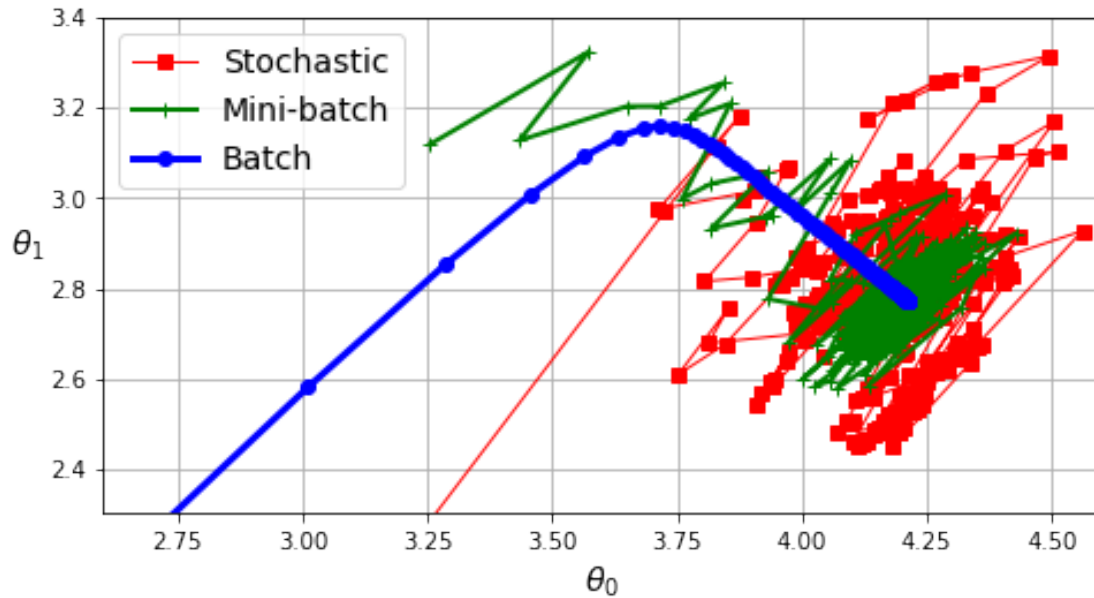
def learning_schedule(t):
    return t0 / (t + t1)

theta_path_mgd = []
for epoch in range(n_epochs):
    shuffled_indices = np.random.permutation(m)
    X_b_shuffled = X_b[shuffled_indices]
    y_shuffled = y[shuffled_indices]
    for iteration in range(0, n_batches_per_epoch):
        idx = iteration * minibatch_size
        xi = X_b_shuffled[idx : idx + minibatch_size]
        yi = y_shuffled[idx : idx + minibatch_size]
        gradients = 2 / minibatch_size * xi.T @ (xi @ theta - yi)
        eta = learning_schedule(iteration)
        theta = theta - eta * gradients
        theta_path_mgd.append(theta)

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")
plt.xlabel(r"$\theta_0$")
plt.ylabel(r"$\theta_1$ ", rotation=0)
plt.axis([2.6, 4.6, 2.3, 3.4])
plt.grid()
save_fig("gradient_descent_paths_plot")
plt.show()

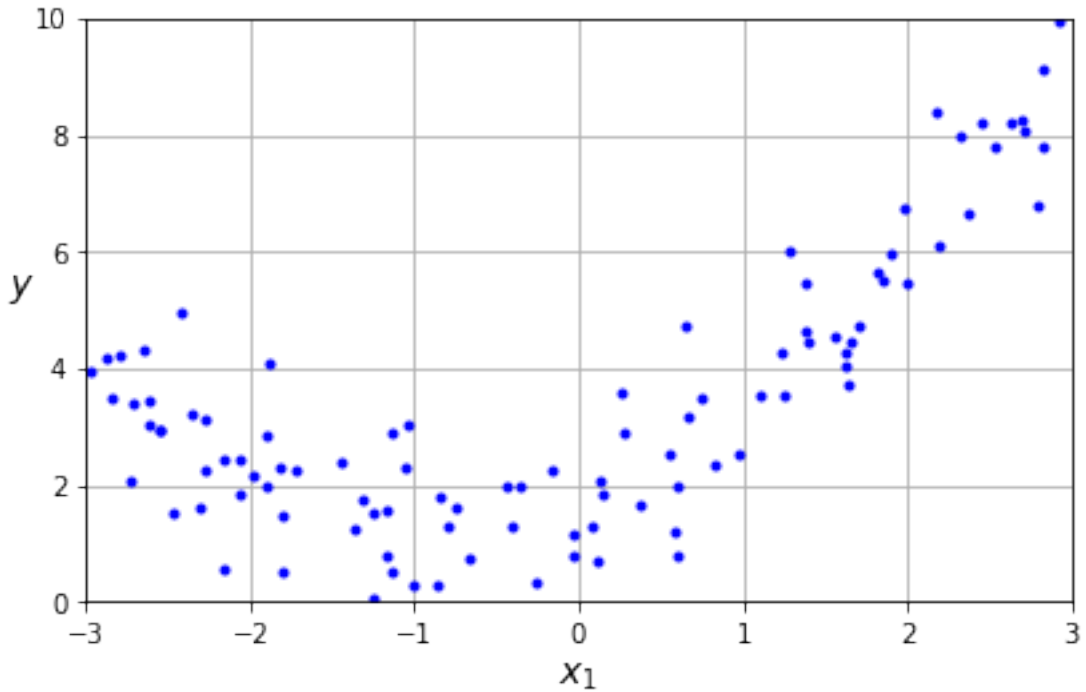
```



## 4 Polynomial Regression

```
[24]: np.random.seed(42)
m = 100
X = 6 * np.random.rand(m, 1) - 3
y = 0.5 * X ** 2 + X + 2 + np.random.randn(m, 1)
```

```
[25]: # extra code - this cell generates and saves Figure 4-12
plt.figure(figsize=(6, 4))
plt.plot(X, y, "b.")
plt.xlabel("$x_1$")
plt.ylabel("$y$", rotation=0)
plt.axis([-3, 3, 0, 10])
plt.grid()
save_fig("quadratic_data_plot")
plt.show()
```



```
[26]: from sklearn.preprocessing import PolynomialFeatures
```

```
poly_features = PolynomialFeatures(degree=2, include_bias=False)
X_poly = poly_features.fit_transform(X)
X[0]
```

```
[26]: array([-0.75275929])
```

```
[27]: X_poly[0]
```

```
[27]: array([-0.75275929,  0.56664654])
```

```
[28]: lin_reg = LinearRegression()
lin_reg.fit(X_poly, y)
lin_reg.intercept_, lin_reg.coef_
```

```
[28]: (array([1.78134581]), array([[0.93366893, 0.56456263]]))
```

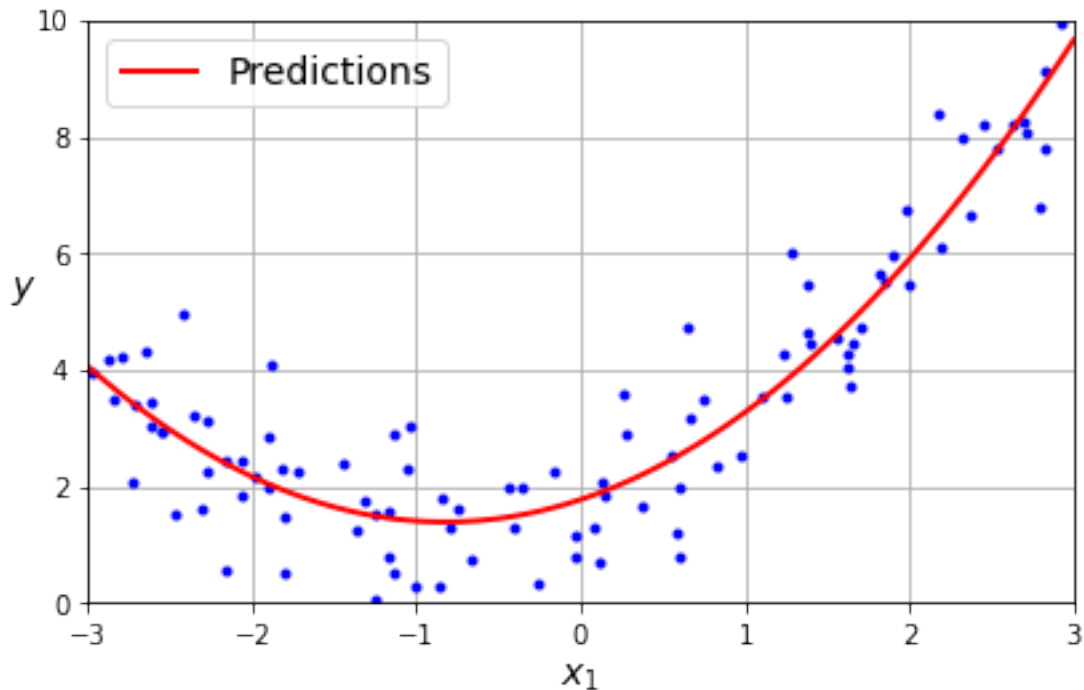
```
[29]: # extra code - this cell generates and saves Figure 4-13
```

```
X_new = np.linspace(-3, 3, 100).reshape(100, 1)
X_new_poly = poly_features.transform(X_new)
y_new = lin_reg.predict(X_new_poly)
```

```

plt.figure(figsize=(6, 4))
plt.plot(X, y, "b.")
plt.plot(X_new, y_new, "r-", linewidth=2, label="Predictions")
plt.xlabel("$x_1$")
plt.ylabel("$y$", rotation=0)
plt.legend(loc="upper left")
plt.axis([-3, 3, 0, 10])
plt.grid()
save_fig("quadratic_predictions_plot")
plt.show()

```



[30]: *# extra code - this cell generates and saves Figure 4-14*

```

from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import make_pipeline

plt.figure(figsize=(6, 4))

for style, width, degree in (("r+", 2, 1), ("b--", 2, 2), ("g-", 1, 300)):
    polybig_features = PolynomialFeatures(degree=degree, include_bias=False)
    std_scaler = StandardScaler()
    lin_reg = LinearRegression()
    polynomial_regression = make_pipeline(polybig_features, std_scaler, lin_reg)
    polynomial_regression.fit(X, y)

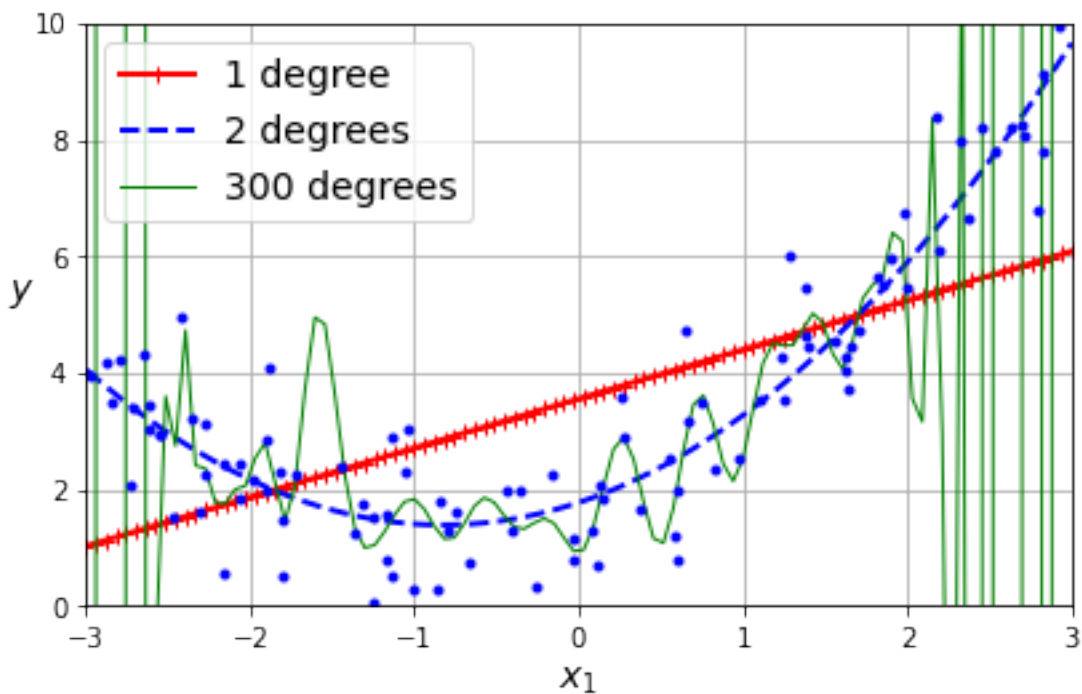
```

```

y_newbig = polynomial_regression.predict(X_new)
label = f"{degree} degree{'s' if degree > 1 else ''}"
plt.plot(X_new, y_newbig, style, label=label, linewidth=width)

plt.plot(X, y, "b.", linewidth=3)
plt.legend(loc="upper left")
plt.xlabel("$x_1$")
plt.ylabel("$y$", rotation=0)
plt.axis([-3, 3, 0, 10])
plt.grid()
save_fig("high_degree_polynomials_plot")
plt.show()

```



## 5 Learning Curves

```

[31]: from sklearn.model_selection import learning_curve

train_sizes, train_scores, valid_scores = learning_curve(
    LinearRegression(), X, y, train_sizes=np.linspace(0.01, 1.0, 40), cv=5,
    scoring="neg_root_mean_squared_error")
train_errors = -train_scores.mean(axis=1)
valid_errors = -valid_scores.mean(axis=1)

```

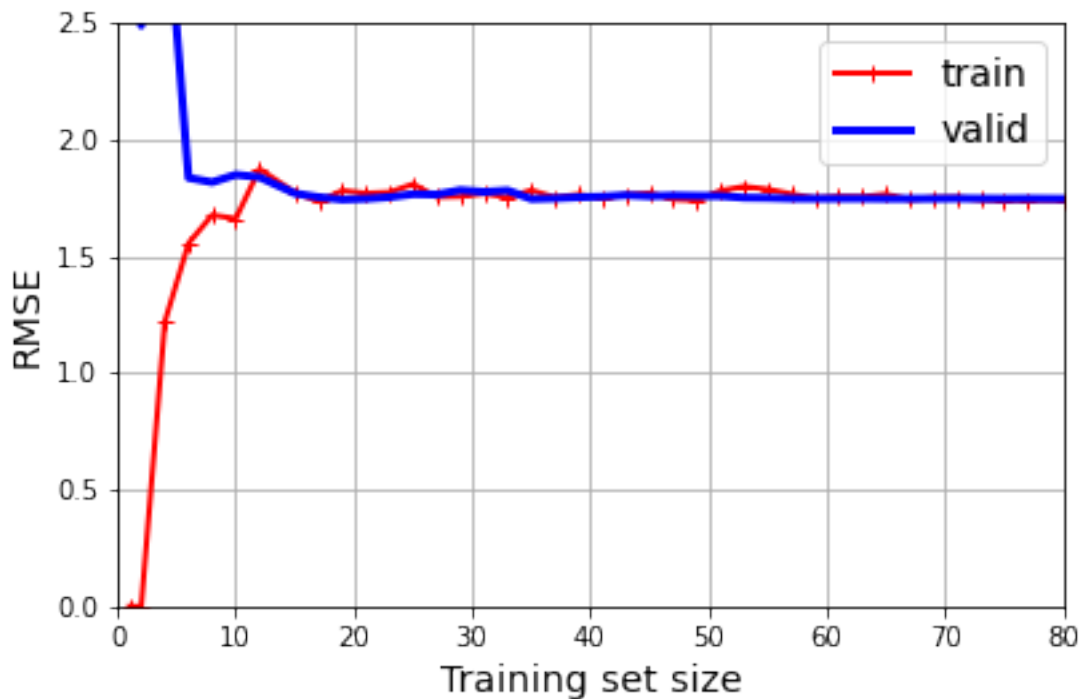
```

plt.figure(figsize=(6, 4)) # extra code - not needed, just formatting
plt.plot(train_sizes, train_errors, "r-+", linewidth=2, label="train")
plt.plot(train_sizes, valid_errors, "b-", linewidth=3, label="valid")

# extra code - beautifies and saves Figure 4-15
plt.xlabel("Training set size")
plt.ylabel("RMSE")
plt.grid()
plt.legend(loc="upper right")
plt.axis([0, 80, 0, 2.5])
save_fig("underfitting_learning_curves_plot")

plt.show()

```



```

[32]: from sklearn.pipeline import make_pipeline

polynomial_regression = make_pipeline(
    PolynomialFeatures(degree=10, include_bias=False),
    LinearRegression())

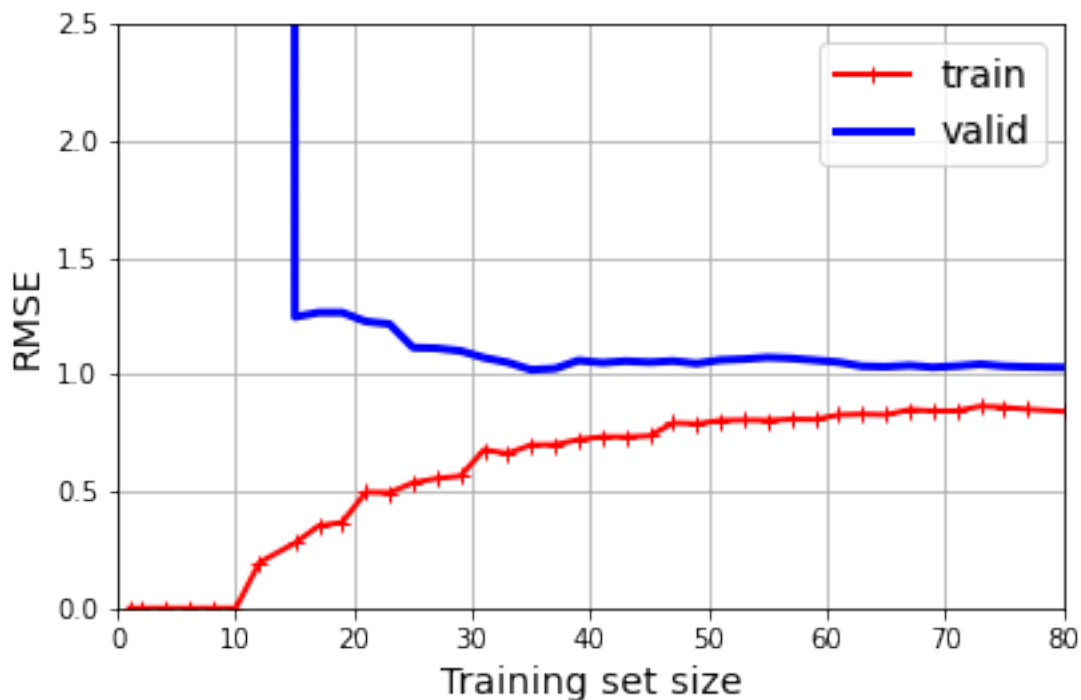
train_sizes, train_scores, valid_scores = learning_curve(
    polynomial_regression, X, y, train_sizes=np.linspace(0.01, 1.0, 40), cv=5,
    scoring="neg_root_mean_squared_error")

```

```
[33]: # extra code - generates and saves Figure 4-16
```

```
train_errors = -train_scores.mean(axis=1)
valid_errors = -valid_scores.mean(axis=1)

plt.figure(figsize=(6, 4))
plt.plot(train_sizes, train_errors, "r+", linewidth=2, label="train")
plt.plot(train_sizes, valid_errors, "b-", linewidth=3, label="valid")
plt.legend(loc="upper right")
plt.xlabel("Training set size")
plt.ylabel("RMSE")
plt.grid()
plt.axis([0, 80, 0, 2.5])
save_fig("learning_curves_plot")
plt.show()
```



## 6 Regularized Linear Models

### 6.1 Ridge Regression

Let's generate a very small and noisy linear dataset:

```
[34]: # extra code - we've done this type of generation several times before
np.random.seed(42)
```

```

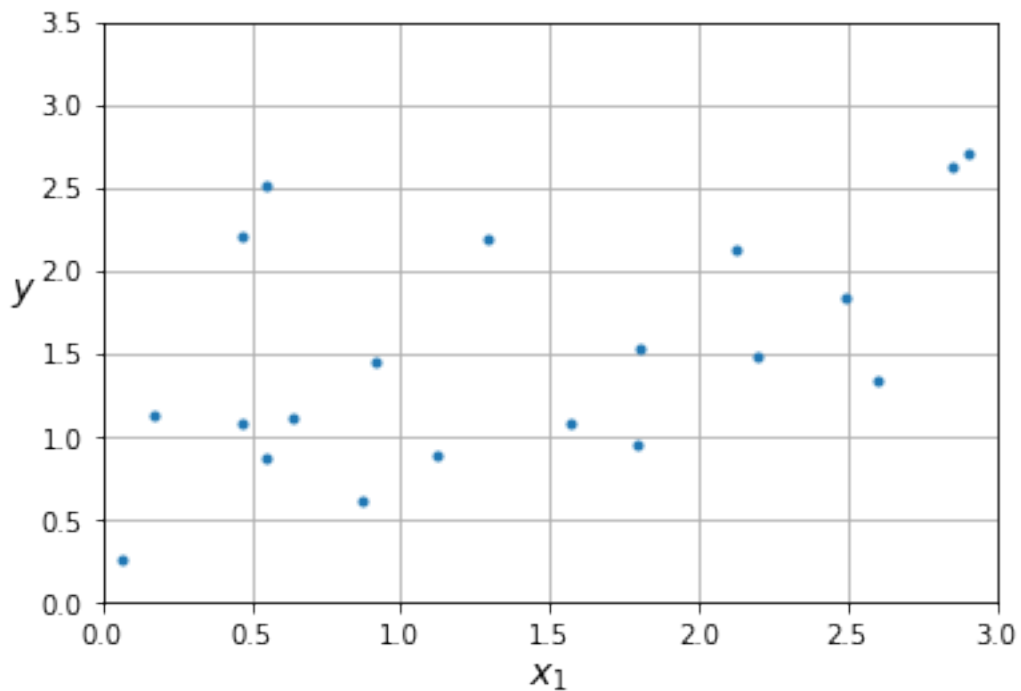
m = 20
X = 3 * np.random.rand(m, 1)
y = 1 + 0.5 * X + np.random.randn(m, 1) / 1.5
X_new = np.linspace(0, 3, 100).reshape(100, 1)

```

```

[35]: # extra code - a quick peek at the dataset we just generated
plt.figure(figsize=(6, 4))
plt.plot(X, y, ".")
plt.xlabel("$x_1$")
plt.ylabel("$y$ ", rotation=0)
plt.axis([0, 3, 0, 3.5])
plt.grid()
plt.show()

```



```

[36]: from sklearn.linear_model import Ridge

ridge_reg = Ridge(alpha=0.1, solver="cholesky")
ridge_reg.fit(X, y)
ridge_reg.predict([[1.5]])

```

```

[36]: array([[1.55325833]])

```

```

[37]: # extra code - this cell generates and saves Figure 4-17

```

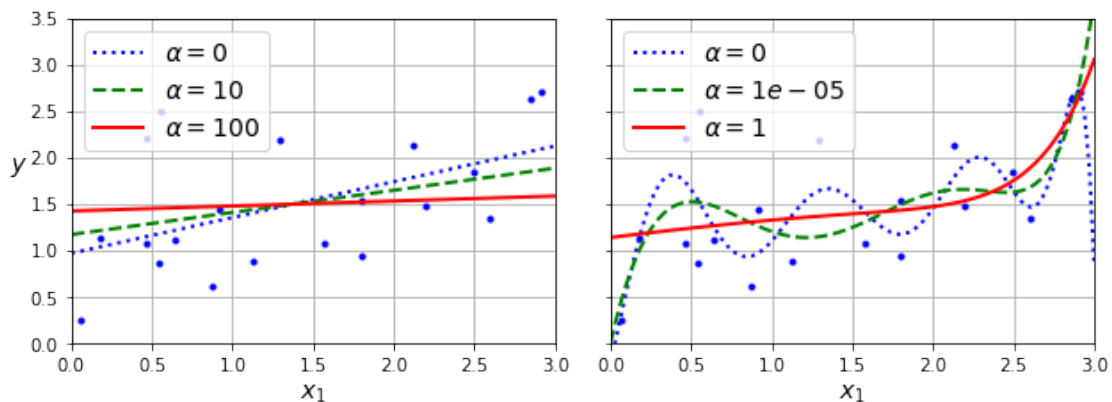


```

def plot_model(model_class, polynomial, alphas, **model_kwargs):
    plt.plot(X, y, "b.", linewidth=3)
    for alpha, style in zip(alphas, ("b:", "g--", "r-")):
        if alpha > 0:
            model = model_class(alpha, **model_kwargs)
        else:
            model = LinearRegression()
        if polynomial:
            model = make_pipeline(
                PolynomialFeatures(degree=10, include_bias=False),
                StandardScaler(),
                model)
        model.fit(X, y)
        y_new_regul = model.predict(X_new)
        plt.plot(X_new, y_new_regul, style, linewidth=2,
                 label=fr"\alpha = {alpha}")
    plt.legend(loc="upper left")
    plt.xlabel("$x_1$")
    plt.axis([0, 3, 0, 3.5])
    plt.grid()

plt.figure(figsize=(9, 3.5))
plt.subplot(121)
plot_model(Ridge, polynomial=False, alphas=(0, 10, 100), random_state=42)
plt.ylabel("$y$ ", rotation=0)
plt.subplot(122)
plot_model(Ridge, polynomial=True, alphas=(0, 10**-5, 1), random_state=42)
plt.gca().axes.yaxis.set_ticklabels([])
save_fig("ridge_regression_plot")
plt.show()

```



```
[38]: sgd_reg = SGDRegressor(penalty="l2", alpha=0.1 / m, tol=None,
                             max_iter=1000, eta0=0.01, random_state=42)
sgd_reg.fit(X, y.ravel()) # y.ravel() because fit() expects 1D targets
sgd_reg.predict([[1.5]])
```

```
[38]: array([1.55302613])
```

```
[39]: # extra code - show that we get roughly the same solution as earlier when
#       we use Stochastic Average GD (solver="sag")
ridge_reg = Ridge(alpha=0.1, solver="sag", random_state=42)
ridge_reg.fit(X, y)
ridge_reg.predict([[1.5]])
```

```
[39]: array([[1.55321535]])
```

```
[40]: # extra code - shows the closed form solution of Ridge regression,
#       compare with the next Ridge model's learned parameters below
alpha = 0.1
A = np.array([[0., 0.], [0., 1.]])
X_b = np.c_[np.ones(m), X]
np.linalg.inv(X_b.T @ X_b + alpha * A) @ X_b.T @ y
```

```
[40]: array([[0.97898394],
            [0.3828496 ]])
```

```
[41]: ridge_reg.intercept_, ridge_reg.coef_ # extra code
```

```
[41]: (array([0.97944909]), array([[0.38251084]]))
```

## 6.2 Lasso Regression

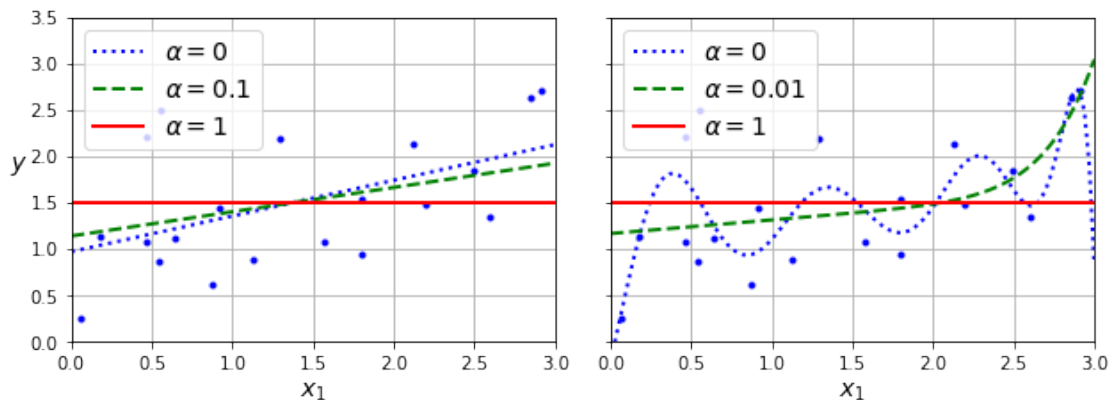
```
[42]: from sklearn.linear_model import Lasso

lasso_reg = Lasso(alpha=0.1)
lasso_reg.fit(X, y)
lasso_reg.predict([[1.5]])
```

```
[42]: array([1.53788174])
```

```
[43]: # extra code - this cell generates and saves Figure 4-18
plt.figure(figsize=(9, 3.5))
plt.subplot(121)
plot_model(Lasso, polynomial=False, alphas=(0, 0.1, 1), random_state=42)
plt.ylabel("$y$ ", rotation=0)
plt.subplot(122)
plot_model(Lasso, polynomial=True, alphas=(0, 1e-2, 1), random_state=42)
plt.gca().axes.yaxis.set_ticklabels([])
save_fig("lasso_regression_plot")
```

```
plt.show()
```



```
[44]: # extra code - this BIG cell generates and saves Figure 4-19
```

```
t1a, t1b, t2a, t2b = -1, 3, -1.5, 1.5

t1s = np.linspace(t1a, t1b, 500)
t2s = np.linspace(t2a, t2b, 500)
t1, t2 = np.meshgrid(t1s, t2s)
T = np.c_[t1.ravel(), t2.ravel()]
Xr = np.array([[1, 1], [1, -1], [1, 0.5]])
yr = 2 * Xr[:, :1] + 0.5 * Xr[:, 1:]

J = (1 / len(Xr) * ((T @ Xr.T - yr.T) ** 2).sum(axis=1)).reshape(t1.shape)

N1 = np.linalg.norm(T, ord=1, axis=1).reshape(t1.shape)
N2 = np.linalg.norm(T, ord=2, axis=1).reshape(t1.shape)

t_min_idx = np.unravel_index(J.argmin(), J.shape)
t1_min, t2_min = t1[t_min_idx], t2[t_min_idx]

t_init = np.array([[0.25], [-1]])

def bgd_path(theta, X, y, l1, l2, core=1, eta=0.05, n_iterations=200):
    path = [theta]
    for iteration in range(n_iterations):
        gradients = (core * 2 / len(X) * X.T @ (X @ theta - y)
                    + l1 * np.sign(theta) + l2 * theta)
        theta = theta - eta * gradients
        path.append(theta)
    return np.array(path)
```

```

fig, axes = plt.subplots(2, 2, sharex=True, sharey=True, figsize=(10.1, 8))

for i, N, l1, l2, title in ((0, N1, 2.0, 0, "Lasso"), (1, N2, 0, 2.0, "Ridge")):
    JR = J + l1 * N1 + l2 * 0.5 * N2 ** 2

    tr_min_idx = np.unravel_index(JR.argmin(), JR.shape)
    t1r_min, t2r_min = t1[tr_min_idx], t2[tr_min_idx]

    levels = np.exp(np.linspace(0, 1, 20)) - 1
    levelsJ = levels * (J.max() - J.min()) + J.min()
    levelsJR = levels * (JR.max() - JR.min()) + JR.min()
    levelsN = np.linspace(0, N.max(), 10)

    path_J = bgd_path(t_init, Xr, yr, l1=0, l2=0)
    path_JR = bgd_path(t_init, Xr, yr, l1, l2)
    path_N = bgd_path(theta=np.array([[2.0], [0.5]]), X=Xr, y=yr,
                        l1=np.sign(l1) / 3, l2=np.sign(l2), core=0)

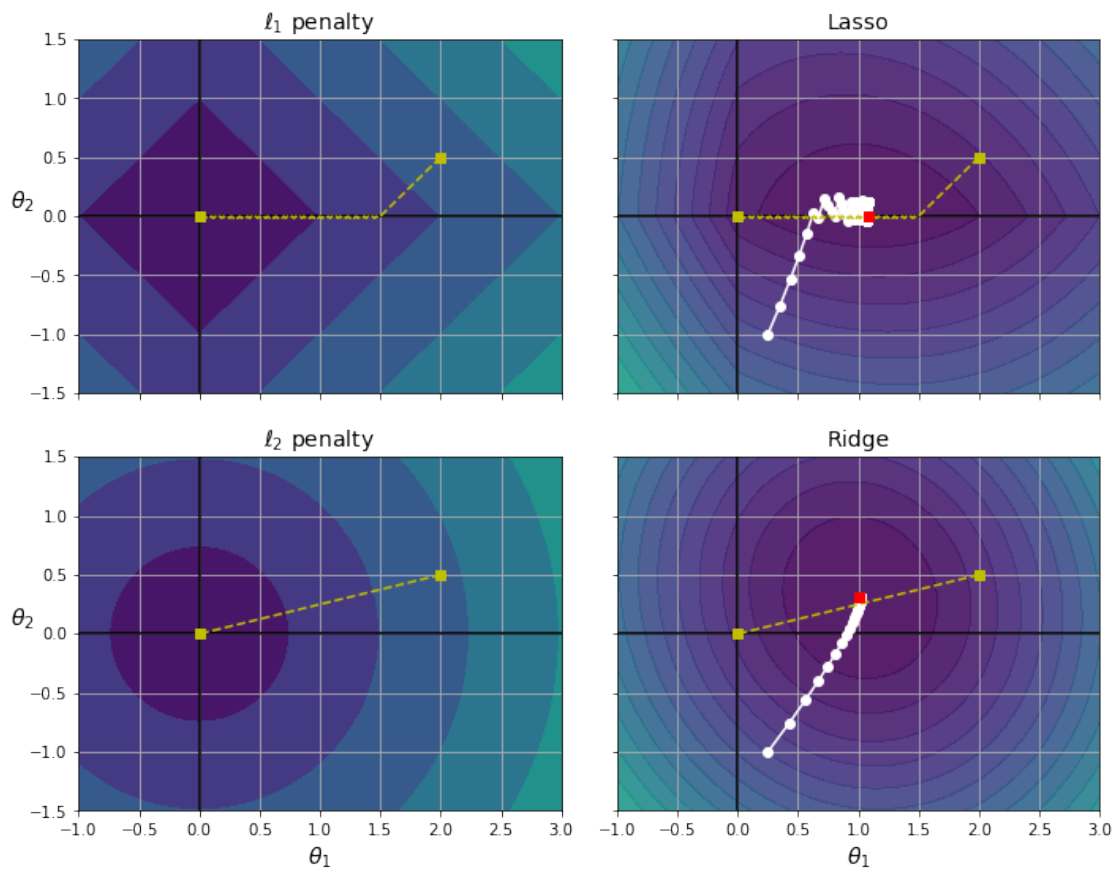
    ax = axes[i, 0]
    ax.grid()
    ax.axhline(y=0, color="k")
    ax.axvline(x=0, color="k")
    ax.contourf(t1, t2, N / 2.0, levels=levelsN)
    ax.plot(path_N[:, 0], path_N[:, 1], "y--")
    ax.plot(0, 0, "ys")
    ax.plot(t1_min, t2_min, "ys")
    ax.set_title(fr"\ell_{i + 1} penalty")
    ax.axis([t1a, t1b, t2a, t2b])
    if i == 1:
        ax.set_xlabel(r"\theta_1")
    ax.set_ylabel(r"\theta_2", rotation=0)

    ax = axes[i, 1]
    ax.grid()
    ax.axhline(y=0, color="k")
    ax.axvline(x=0, color="k")
    ax.contourf(t1, t2, JR, levels=levelsJR, alpha=0.9)
    ax.plot(path_JR[:, 0], path_JR[:, 1], "w-o")
    ax.plot(path_N[:, 0], path_N[:, 1], "y--")
    ax.plot(0, 0, "ys")
    ax.plot(t1_min, t2_min, "ys")
    ax.plot(t1r_min, t2r_min, "rs")
    ax.set_title(title)
    ax.axis([t1a, t1b, t2a, t2b])
    if i == 1:
        ax.set_xlabel(r"\theta_1")

save_fig("lasso_vs_ridge_plot")

```

```
plt.show()
```



### 6.3 Elastic Net

```
[45]: from sklearn.linear_model import ElasticNet

elastic_net = ElasticNet(alpha=0.1, l1_ratio=0.5)
elastic_net.fit(X, y)
elastic_net.predict([[1.5]])
```

```
[45]: array([1.54333232])
```

### 6.4 Early Stopping

Let's go back to the quadratic dataset we used earlier:

```
[46]: from copy import deepcopy
from sklearn.metrics import mean_squared_error
from sklearn.preprocessing import StandardScaler
```

```

# extra code - creates the same quadratic dataset as earlier and splits it
np.random.seed(42)
m = 100
X = 6 * np.random.rand(m, 1) - 3
y = 0.5 * X ** 2 + X + 2 + np.random.randn(m, 1)
X_train, y_train = X[: m // 2], y[: m // 2, 0]
X_valid, y_valid = X[m // 2 :], y[m // 2 :, 0]

preprocessing = make_pipeline(PolynomialFeatures(degree=90, include_bias=False),
                              StandardScaler())
X_train_prep = preprocessing.fit_transform(X_train)
X_valid_prep = preprocessing.transform(X_valid)
sgd_reg = SGDRegressor(penalty=None, eta0=0.002, random_state=42)
n_epochs = 500
best_valid_rmse = float('inf')
train_errors, val_errors = [], [] # extra code - it's for the figure below

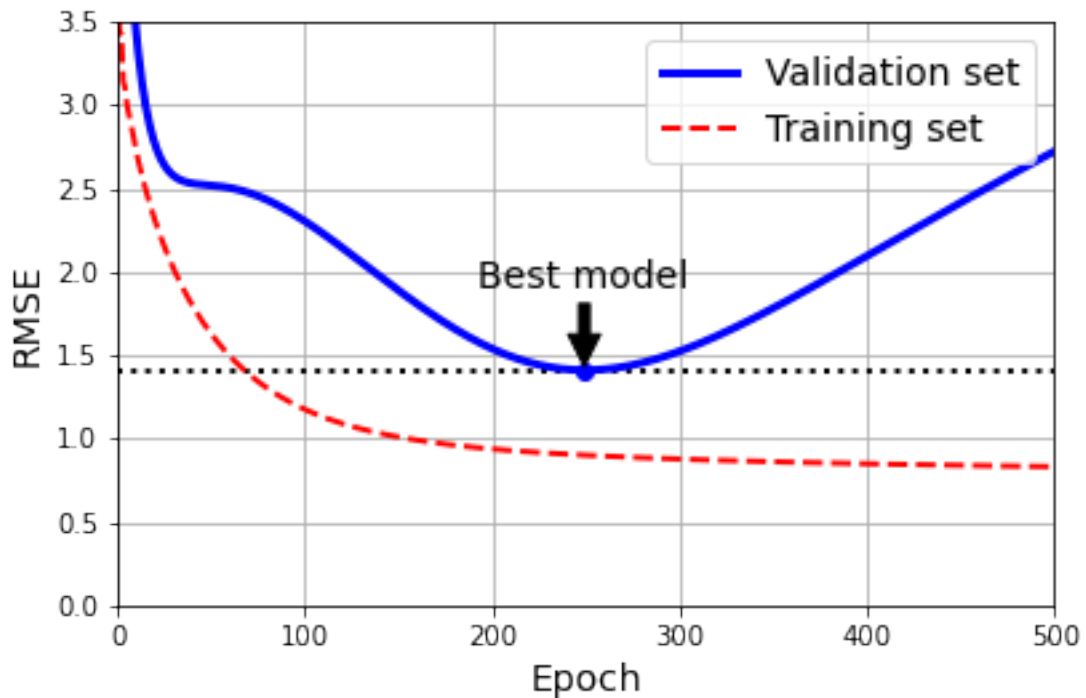
for epoch in range(n_epochs):
    sgd_reg.partial_fit(X_train_prep, y_train)
    y_valid_predict = sgd_reg.predict(X_valid_prep)
    val_error = mean_squared_error(y_valid, y_valid_predict, squared=False)
    if val_error < best_valid_rmse:
        best_valid_rmse = val_error
        best_model = deepcopy(sgd_reg)

    # extra code - we evaluate the train error and save it for the figure
    y_train_predict = sgd_reg.predict(X_train_prep)
    train_error = mean_squared_error(y_train, y_train_predict, squared=False)
    val_errors.append(val_error)
    train_errors.append(train_error)

# extra code - this section generates and saves Figure 4-20
best_epoch = np.argmin(val_errors)
plt.figure(figsize=(6, 4))
plt.annotate('Best model',
            xy=(best_epoch, best_valid_rmse),
            xytext=(best_epoch, best_valid_rmse + 0.5),
            ha="center",
            arrowprops=dict(facecolor='black', shrink=0.05))
plt.plot([0, n_epochs], [best_valid_rmse, best_valid_rmse], "k:", linewidth=2)
plt.plot(val_errors, "b-", linewidth=3, label="Validation set")
plt.plot(best_epoch, best_valid_rmse, "bo")
plt.plot(train_errors, "r--", linewidth=2, label="Training set")
plt.legend(loc="upper right")
plt.xlabel("Epoch")
plt.ylabel("RMSE")
plt.axis([0, n_epochs, 0, 3.5])

```

```
plt.grid()
save_fig("early_stopping_plot")
plt.show()
```



## 7 Logistic Regression

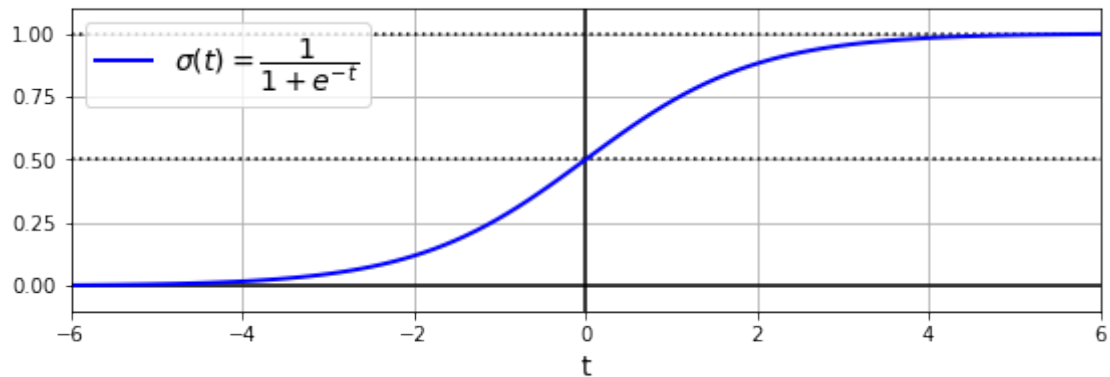
### 7.1 Estimating Probabilities

[47]: # extra code - generates and saves Figure 4-21

```
lim = 6
t = np.linspace(-lim, lim, 100)
sig = 1 / (1 + np.exp(-t))

plt.figure(figsize=(8, 3))
plt.plot([-lim, lim], [0, 0], "k-")
plt.plot([-lim, lim], [0.5, 0.5], "k:")
plt.plot([-lim, lim], [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")
```

```
plt.axis([-lim, lim, -0.1, 1.1])
plt.gca().set_yticks([0, 0.25, 0.5, 0.75, 1])
plt.grid()
save_fig("logistic_function_plot")
plt.show()
```



## 7.2 Decision Boundaries

```
[48]: from sklearn.datasets import load_iris
```

```
iris = load_iris(as_frame=True)
list(iris)
```

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

```
[49]: print(iris.DESCR) # extra code - it's a bit too long
```

```
.. _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.

- 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 ...

```
[50]: iris.data.head(3)
```

```
[50]:   sepal length (cm)  sepal width (cm)  petal length (cm)  petal width (cm)
0          5.1         3.5         1.4         0.2
1          4.9         3.0         1.4         0.2
2          4.7         3.2         1.3         0.2
```

```
[51]: iris.target.head(3) # note that the instances are not shuffled
```

```
[51]: 0    0
      1    0
      2    0
      Name: target, dtype: int64
```

```
[52]: iris.target_names
```

```
[52]: array(['setosa', 'versicolor', 'virginica'], dtype='<U10')
```

```
[53]: from sklearn.linear_model import LogisticRegression
      from sklearn.model_selection import train_test_split

      X = iris.data[["petal width (cm)"].values
      y = iris.target_names[iris.target] == 'virginica'
      X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)

      log_reg = LogisticRegression(random_state=42)
      log_reg.fit(X_train, y_train)
```

```
[53]: LogisticRegression(random_state=42)
```

```
[54]: X_new = np.linspace(0, 3, 1000).reshape(-1, 1) # reshape to get a column vector
      y_proba = log_reg.predict_proba(X_new)
      decision_boundary = X_new[y_proba[:, 1] >= 0.5][0, 0]

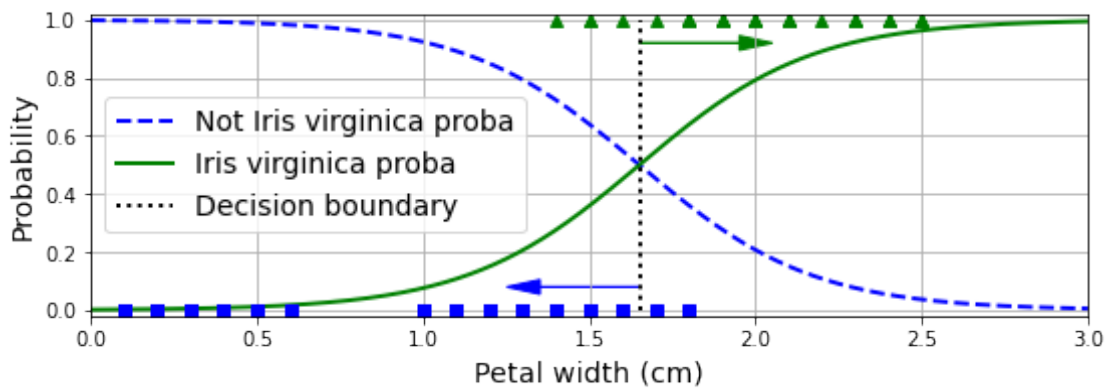
      plt.figure(figsize=(8, 3)) # extra code - not needed, just formatting
      plt.plot(X_new, y_proba[:, 0], "b--", linewidth=2,
              label="Not Iris virginica proba")
      plt.plot(X_new, y_proba[:, 1], "g-", linewidth=2, label="Iris virginica proba")
      plt.plot([decision_boundary, decision_boundary], [0, 1], "k:", linewidth=2,
              label="Decision boundary")
```

```

# extra code - this section beautifies and saves Figure 4-23
plt.arrow(x=decision_boundary, y=0.08, dx=-0.3, dy=0,
          head_width=0.05, head_length=0.1, fc="b", ec="b")
plt.arrow(x=decision_boundary, y=0.92, dx=0.3, dy=0,
          head_width=0.05, head_length=0.1, fc="g", ec="g")
plt.plot(X_train[y_train == 0], y_train[y_train == 0], "bs")
plt.plot(X_train[y_train == 1], y_train[y_train == 1], "g^")
plt.xlabel("Petal width (cm)")
plt.ylabel("Probability")
plt.legend(loc="center left")
plt.axis([0, 3, -0.02, 1.02])
plt.grid()
save_fig("logistic_regression_plot")

plt.show()

```



```
[55]: decision_boundary
```

```
[55]: 1.6516516516516517
```

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

```
[56]: array([ True, False])
```