# 回归

本节代码来自: 书籍代码 https://github.com/ageron/handson-ml3 推荐自学

## 1.1 线性回归

- np.random.rand()
  - 该函数括号内的参数指定的是返回结果的形状
  - 返回结果中的每一个元素是服从0~1均匀分布的随机样本值

```python
In [1]: import matplotlib.pyplot as plt

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

```python
In [2]: from pathlib import Path

IMAGES_PATH = Path() / "images"
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)
```
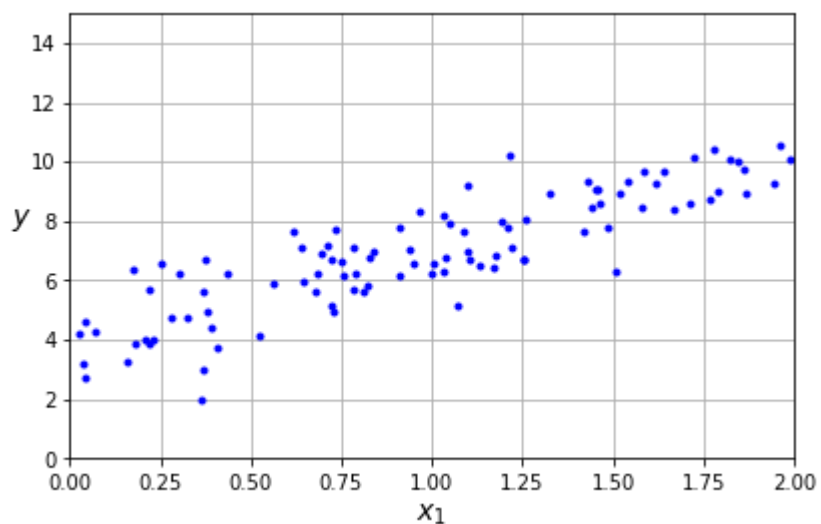
```python
In [3]: import numpy as np

np.random.seed(2023)  # 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
```

In [4]:
```python
# 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()
```



In [5]:
```python
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
```

In [6]:
```python
theta_best
```

Out[6]:
```
array([[3.87411963],
       [3.13856283]])
```

In [7]:
```python
X_new = np.array([[0], [2]])
X_new_b = add_dummy_feature(X_new)  # add x0 = 1 to each instance
```

```python
y_predict = X_new_b @ theta_best
y_predict
```

Out[7]:  array([[ 3.87411963],
                [10.1512453 ]])
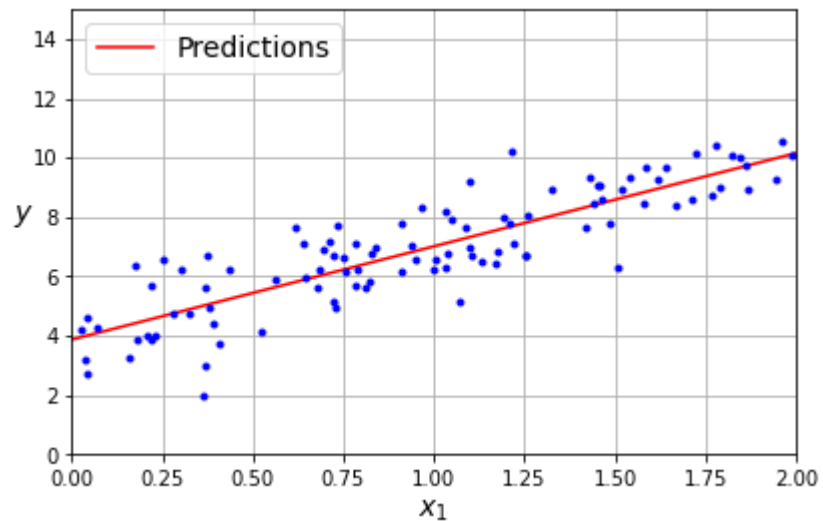
In [8]:
```python
import matplotlib.pyplot as plt

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()
```



In [9]:
```python
from sklearn.linear_model import LinearRegression

lin_reg = LinearRegression()
```

```python
lin_reg.fit(X, y)
lin_reg.intercept_, lin_reg.coef_
```

Out[9]: (array([3.87411963]), array([[3.13856283]]))

## 1.2 梯度下降

### 批量梯度下降

In [10]:
```python
eta = 0.1  # learning rate
n_epochs = 1000
m = len(X_b)  # number of instances

np.random.seed(2023)
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
```

In [11]:
```python
theta
```

Out[11]:
```python
array([[3.87411963],
       [3.13856283]])
```

In [12]:
```python
# 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)
```
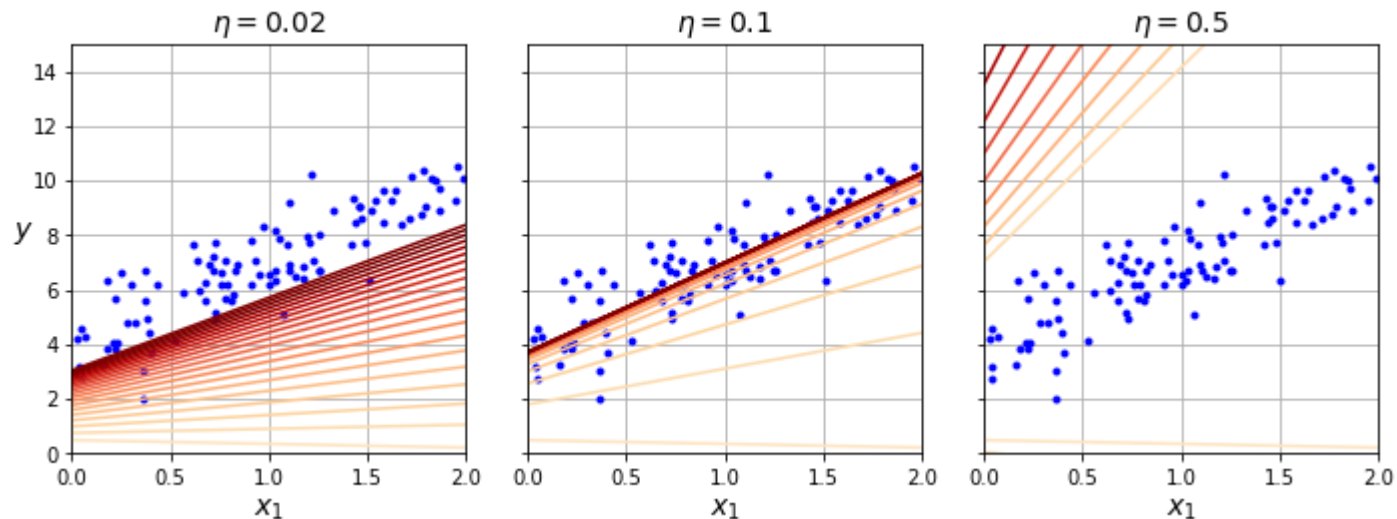
```python
        theta = theta - eta * gradients
        theta_path.append(theta)
    plt.xlabel("$x_1$")
    plt.axis([0, 2, 0, 15])
    plt.grid()
    plt.title(fr"$\eta = {eta}$")
    return theta_path

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()
```



```python
In [13]: theta_path_sgd = []  # extra code – we need to store the path of theta in the
                              #              parameter space to plot the next figure
```

In [14]:
```python
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()
```
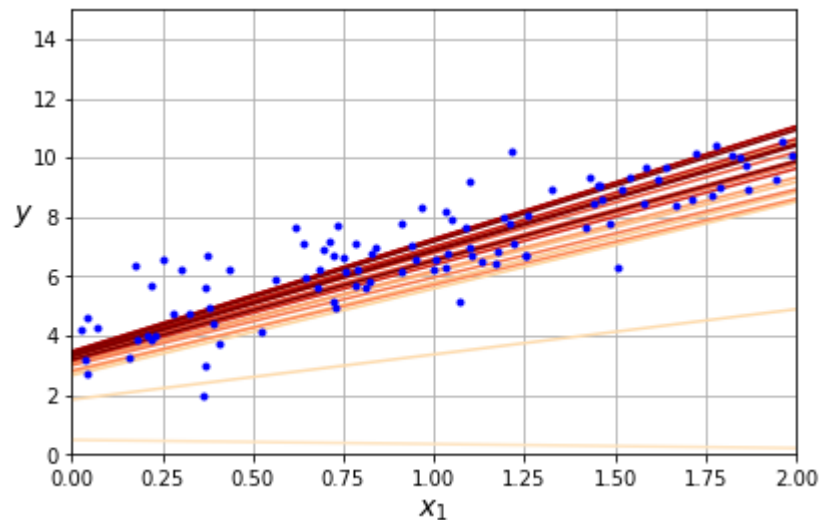
```
In [15]: from sklearn.linear_model import SGDRegressor

         sgd_reg = SGDRegressor(max_iter=1000, tol=1e-5, penalty=None, eta0=0.01,
                                n_iter_no_change=100, random_state=2023)
         sgd_reg.fit(X, y.ravel())  # y.ravel() because fit() expects 1D targets
         sgd_reg.intercept_, sgd_reg.coef_
```
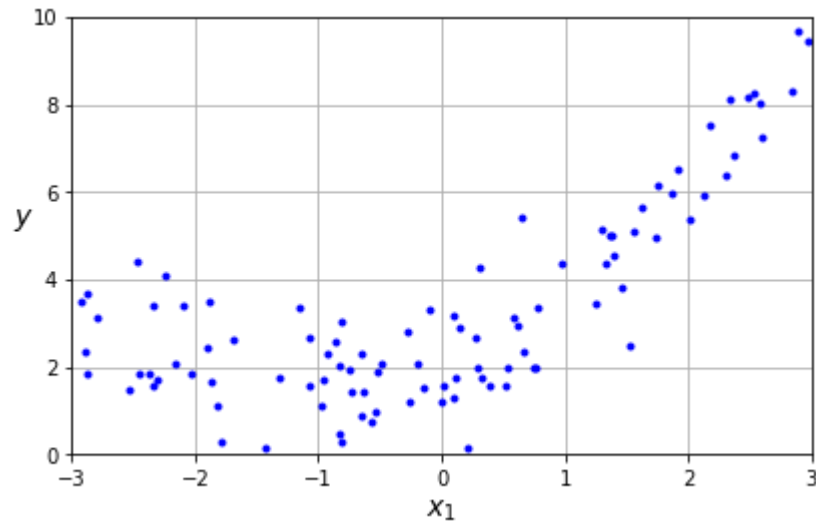
Out[15]: (array([3.87253499]), array([3.13933332]))

## 1.3 多项式回归

```
In [16]: np.random.seed(2023)
         m = 100
         X = 6 * np.random.rand(m, 1) - 3
         y = 0.5 * X ** 2 + X + 2 + np.random.randn(m, 1)
```

```
In [17]: # 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()
```



In [18]:
```
from sklearn.preprocessing import PolynomialFeatures

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

```
[-1.06807018]
[-1.06807018  1.1407739 ]
```

In [19]:
```
lin_reg = LinearRegression()
lin_reg.fit(X_poly, y)
lin_reg.intercept_, lin_reg.coef_
```
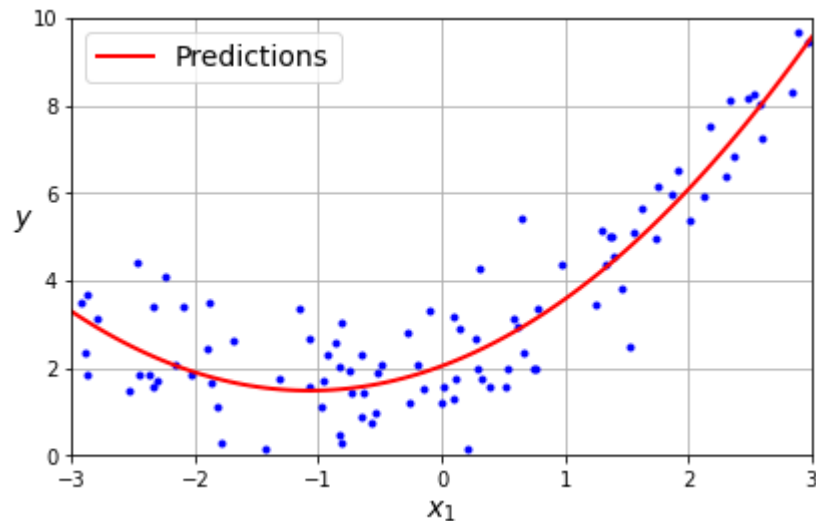
Out[19]:  (array([2.04672232]), array([[1.04476021, 0.48703958]]))

In [20]:
```
# 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))
```

```python
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()
```
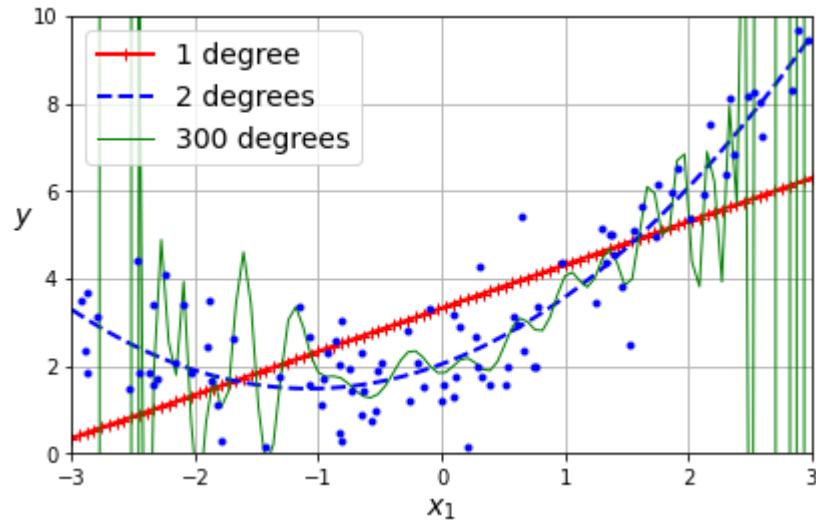


In [21]:
```python
# extra code – this cell generates and saves Figure 4–14

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)
```

```python
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()
```



```python
In [22]:   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()
```
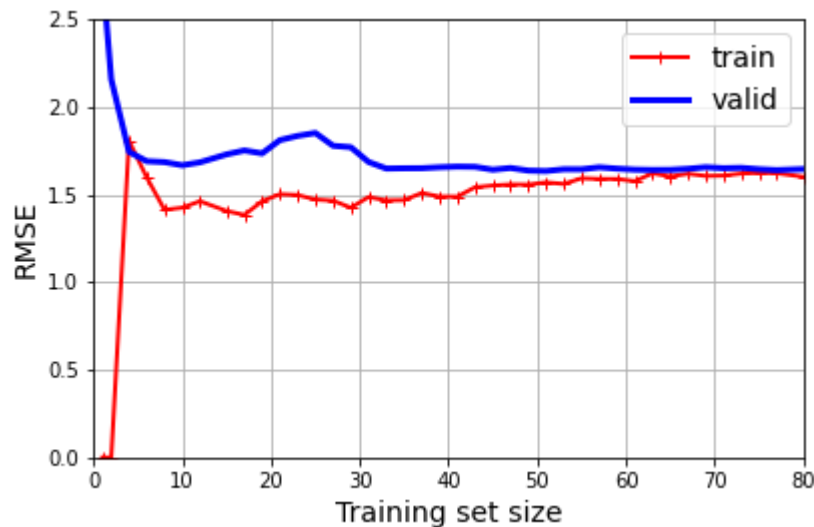
```python
plt.legend(loc="upper right")
plt.axis([0, 80, 0, 2.5])
save_fig("underfitting_learning_curves_plot")

plt.show()
```



```python
In [23]:  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")
```
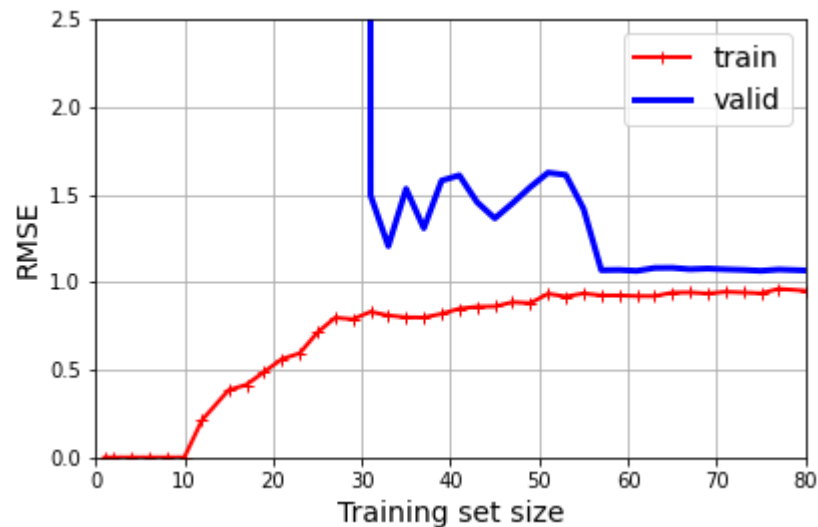
```python
In [24]:  # extra code – generates and saves Figure 4–16

          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")
```
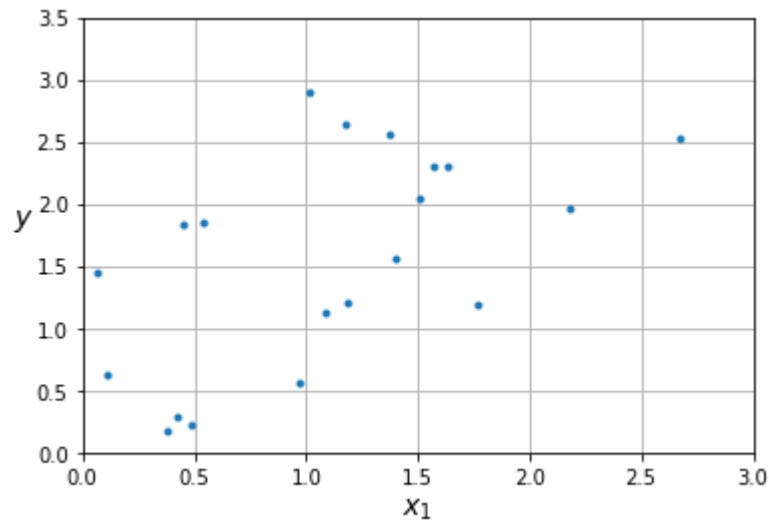
```
plt.ylabel("RMSE")
plt.grid()
plt.axis([0, 80, 0, 2.5])
save_fig("learning_curves_plot")
plt.show()
```



In [25]:
```
# extra code – we've done this type of generation several times before
np.random.seed(2023)
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)
```

In [26]:
```
# 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()
```

In [27]:
```python
sgd_reg = SGDRegressor(penalty="l2", alpha=0.1 / m, tol=None,
                       max_iter=1000, eta0=0.01, random_state=2023)
sgd_reg.fit(X, y.ravel())  # y.ravel() because fit() expects 1D targets
sgd_reg.predict([[1.5]])
```

Out[27]: array([1.85152736])

In [28]:
```python
from sklearn.linear_model import Ridge

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

Out[28]: array([[1.85177963]])

In [29]:
```python
ridge_reg.intercept_, ridge_reg.coef_  # extra code
```

Out[29]: (array([0.81292784]), array([[0.69256786]]))

In [30]:
```python
from sklearn.pipeline import Pipeline
model_poly = Pipeline([("poly",PolynomialFeatures(degree=10, include_bias=False)),
("scaler", StandardScaler()),
("ridge", Ridge(alpha=0.1, solver="cholesky"))
])
```

```
model_poly.fit(X, y)
model_poly.predict([[1.5]])
```

Out[30]: array([[1.97939809]])

In [31]:
```
from sklearn.linear_model import Lasso

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

Out[31]: array([1.76662782])

In [32]:
```
lasso_reg.intercept_, lasso_reg.coef_   # extra code
```

Out[32]: (array([1.04624973]), array([0.48025206]))

In [33]:
```
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]])
```

Out[33]: array([1.78727936])

In [34]:
```
elastic_net.intercept_, elastic_net.coef_   # extra code
```

Out[34]: (array([0.98966306]), array([0.5317442]))

## 1.4 逻辑回归

In [35]:
```
from sklearn.datasets import load_iris

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

Out[35]: ['data',
         'target',
         'frame',
         'target_names',
         'DESCR',
         'feature_names',
         'filename',
         'data_module']

In [36]: `iris.data.head()`

Out[36]:

|   | 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 |
| **3** | 4.6 | 3.1 | 1.5 | 0.2 |
| **4** | 5.0 | 3.6 | 1.4 | 0.2 |

In [37]: `iris.target_names`

Out[37]: array(['setosa', 'versicolor', 'virginica'], dtype='<U10')

In [38]:
```python
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=2023)

log_reg = LogisticRegression(random_state=2023)
log_reg.fit(X_train, y_train)
```

Out[38]:
```
▼        LogisticRegression

LogisticRegression(random_state=2023)
```

In [39]: `X_new = np.linspace(0, 3, 1000).reshape(-1, 1)`  # *reshape to get a column vector*
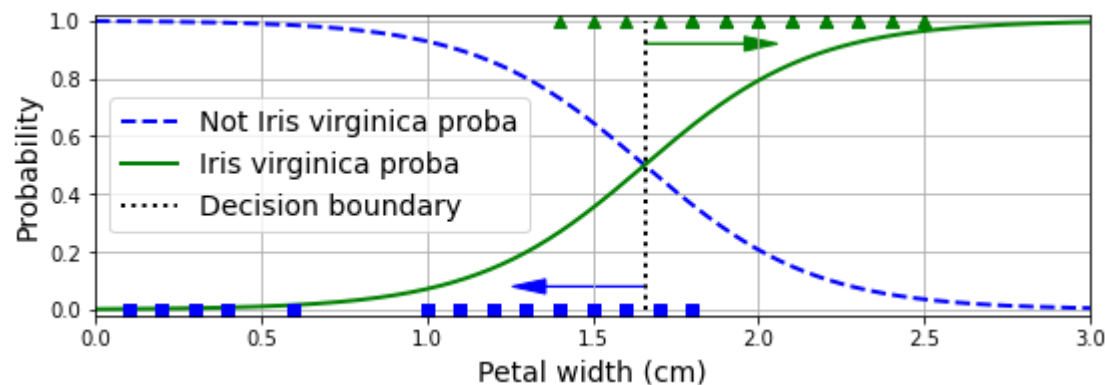
```python
y_proba = log_reg.predict_proba(X_new)
decision_boundary = X_new[y_proba[:, 1] >= 0.5][0, 0]

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()
```



```
In [40]:  decision_boundary
```

```
Out[40]:  1.6576576576576576
```

In [41]:
```python
log_reg.predict([[1.7], [1.5]])
```

Out[41]: array([ True, False])

In [42]:
```python
X = iris.data[["petal length (cm)", "petal width (cm)"]].values
y = iris["target"]
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=2023)

softmax_reg = LogisticRegression(C=10,random_state=2023)
softmax_reg.fit(X_train, y_train)
```

Out[42]:
▼            LogisticRegression

LogisticRegression(C=10, random_state=2023)

In [43]:
```python
softmax_reg.predict([[5, 2]])
```

Out[43]: array([2])

In [44]:
```python
softmax_reg.predict_proba([[5, 2]]).round(2)
```

Out[44]: array([[0.  , 0.06, 0.94]])