

[更多机器学习资源](#)
[目录](#)
[Open in Colab](#)

Appendix: Figure Code

附录：生成图像的代码

Many of the figures used throughout this text are created in-place by code that appears in print. In a few cases, however, the required code is long enough (or not immediately relevant enough) that we instead put it here for reference.

本书中的大多数图表已经使用正文中的代码生成了。然而在一些情况下，需要的代码可能很长（或者其与正文内容直接相关），因此这部分代码就放在附录中供参考。

```
In [1]: %matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns

In [2]: import os
import os.path.exists('figures'):
    os.makedirs('figures')
```

Broadcasting

广播

[查看所在正文](#)

```
In [3]: # Matplotlib中沿用: 参见http://www.astroml.org/book_figures/appendix/fig_broadcast_visual.html
import numpy as np
from matplotlib import pyplot as plt

#-----绘制没有边界的图表和轴
fig = plt.figure(figsize=(6, 4.5), facecolor='w')
ax = plt.axes([0, 0, 1, 1], xticks=[], yticks=[], frameon=False)

def draw_cube(ax, xy, size, depth=0.4,
              edges=None, label=None, label_kwargs=None, **kwargs):
    """draw and label cube. edges is a list of numbers between
    1 and 12, specifying which of the 12 cube edges to draw"""
    if edges is None:
        edges = range(1, 13)
    x, y = xy
    if 1 in edges:
        ax.plot([x, x + size],
                [y + size, y + size], **kwargs)
    if 2 in edges:
        ax.plot([x + size, x + size],
                [y, y + size], **kwargs)
    if 3 in edges:
        ax.plot([x, x + size],
                [y, y], **kwargs)
    if 4 in edges:
        ax.plot([x, x],
                [y, y + size], **kwargs)
    if 5 in edges:
        ax.plot([x, x + depth],
                [y + size, y + depth + size], **kwargs)
    if 6 in edges:
        ax.plot([x + size, x + size + depth],
                [y + size, y + depth + size], **kwargs)
    if 7 in edges:
        ax.plot([x + size, x + size + depth],
                [y, y + depth], **kwargs)
    if 8 in edges:
        ax.plot([x + depth, x + depth],
                [y, y + size], **kwargs)
    if 9 in edges:
        ax.plot([x + depth, x + depth + size],
                [y + depth + size, y + depth + size], **kwargs)
    if 10 in edges:
        ax.plot([x + depth + size, x + depth + size],
                [y + depth, y + depth + size], **kwargs)
    if 11 in edges:
        ax.plot([x + depth, x + depth + size],
                [y + depth, y + depth], **kwargs)
    if 12 in edges:
        ax.plot([x + depth + size, x + depth],
                [y + depth, y + depth + size], **kwargs)
    if label:
        if label_kwargs is None:
            label_kwargs = {}
        ax.text(x + 0.5 * size, y + 0.5 * size, label,
                ha='center', va='center', **label_kwargs)
    solid = dict(c='black', ls='-', lw=1,
                 label_kwargs=dict(color='k'))
    dotted = dict(c='black', ls='--', lw=0.5, alpha=0.5,
                  label_kwargs=dict(color='gray'))
    depth = 0.3

#-----绘制上部操作: 向量加标量
draw_cube(ax, (1, 10), 1, depth, [1, 2, 3, 4, 5, 6, 9], '1', **solid)
draw_cube(ax, (2, 10), 1, depth, [1, 2, 3, 6, 9], '2', **solid)
draw_cube(ax, (3, 10), 1, depth, [1, 2, 3, 6, 7, 9, 10], '2', **solid)

draw_cube(ax, (6, 10), 1, depth, [1, 2, 3, 4, 5, 6, 7, 9, 10], '5', **dotted)
draw_cube(ax, (7, 10), 1, depth, [1, 2, 3, 6, 7, 9, 10, 11], '5', **dotted)
draw_cube(ax, (8, 7.5), 1, depth, [1, 2, 3, 6, 7, 9, 10], '2', **solid)

draw_cube(ax, (12, 10), 1, depth, [1, 2, 3, 4, 5, 6, 9], '1', **solid)
draw_cube(ax, (13, 10), 1, depth, [1, 2, 3, 6, 9], '1', **dotted)
draw_cube(ax, (14, 10), 1, depth, [1, 2, 3, 6, 7, 9, 10], '1', **dotted)

ax.text(5, 10.5, '+', size=12, ha='center', va='center')
ax.text(10.5, 10.5, '=', size=12, ha='center', va='center')
ax.text(11, 11.5, r'$\text{np.arange}(3) + 5$',
        size=12, ha='left', va='bottom')

#-----绘制中部操作: 矩阵加向量

# 第一部分
draw_cube(ax, (1, 7.5), 1, depth, [1, 2, 3, 4, 5, 6, 9], '1', **solid)
draw_cube(ax, (2, 7.5), 1, depth, [1, 2, 3, 6, 9], '1', **solid)
draw_cube(ax, (3, 7.5), 1, depth, [1, 2, 3, 6, 7, 9, 10], '1', **solid)

draw_cube(ax, (1, 6.5), 1, depth, [2, 3, 4], '1', **solid)
draw_cube(ax, (2, 6.5), 1, depth, [2, 3], '1', **solid)
draw_cube(ax, (3, 6.5), 1, depth, [2, 3, 6, 9], '1', **solid)

draw_cube(ax, (1, 5.5), 1, depth, [2, 3, 4], '1', **solid)
draw_cube(ax, (2, 5.5), 1, depth, [2, 3], '1', **dotted)
draw_cube(ax, (3, 5.5), 1, depth, [2, 3, 7, 10], '1', **solid)

# 第二部分
draw_cube(ax, (6, 7.5), 1, depth, [1, 2, 3, 4, 5, 6, 9], '0', **solid)
draw_cube(ax, (7, 7.5), 1, depth, [1, 2, 3, 6, 9], '1', **solid)
draw_cube(ax, (8, 7.5), 1, depth, [1, 2, 3, 7, 10], '2', **solid)

draw_cube(ax, (6, 6.5), 1, depth, range(2, 13), '0', **dotted)
draw_cube(ax, (7, 6.5), 1, depth, [2, 3, 6, 7, 9, 10, 11], '1', **dotted)
draw_cube(ax, (8, 6.5), 1, depth, [2, 3, 6, 7, 9, 10, 11], '2', **dotted)

draw_cube(ax, (6, 5.5), 1, depth, [2, 3, 4, 7, 8, 10, 11], '0', **dotted)
draw_cube(ax, (7, 5.5), 1, depth, [2, 3, 7, 10, 11], '1', **dotted)
draw_cube(ax, (8, 5.5), 1, depth, [2, 3, 7, 10, 11], '2', **dotted)

# 第三部分
draw_cube(ax, (12, 7.5), 1, depth, [1, 2, 3, 4, 5, 6, 9], '1', **solid)
draw_cube(ax, (14, 7.5), 1, depth, [1, 2, 3, 6, 7, 9, 10], '3', **solid)

draw_cube(ax, (12, 6.5), 1, depth, [2, 3, 4], '1', **solid)
draw_cube(ax, (13, 6.5), 1, depth, [2, 3], '1', **solid)
draw_cube(ax, (14, 6.5), 1, depth, [2, 3, 7, 10], '3', **solid)

draw_cube(ax, (12, 5.5), 1, depth, [2, 3, 4], '1', **solid)
draw_cube(ax, (13, 5.5), 1, depth, [2, 3], '3', **solid)
draw_cube(ax, (14, 5.5), 1, depth, [2, 3, 7, 10], '3', **solid)

ax.text(5, 7.0, '+', size=12, ha='center', va='center')
ax.text(10.5, 7.0, '=', size=12, ha='center', va='center')
ax.text(11, 9.0, r'$\text{np.ones}(3, 3) + \text{np.arange}(3)$',
        size=12, ha='left', va='bottom')

#-----绘制底部操作: 向量加向量, 双广播

# 第一部分
draw_cube(ax, (1, 3), 1, depth, [1, 2, 3, 4, 5, 6, 7, 9, 10], '0', **solid)
draw_cube(ax, (1, 2), 1, depth, [2, 3, 4, 7, 10], '2', **solid)
draw_cube(ax, (2, 2), 1, depth, [1, 2, 3, 6, 7, 9, 10, 11], '0', **dotted)
draw_cube(ax, (2, 3), 1, depth, [2, 3, 7, 10, 11], '1', **dotted)
draw_cube(ax, (3, 1), 1, depth, [2, 3, 6, 7, 10, 11], '2', **dotted)

draw_cube(ax, (3, 3), 1, depth, [1, 2, 3, 6, 7, 9, 10, 11], '0', **dotted)
draw_cube(ax, (3, 2), 1, depth, [2, 3, 7, 10, 11], '1', **dotted)
draw_cube(ax, (3, 1), 1, depth, [2, 3, 7, 10, 11], '2', **dotted)

# 第二部分
draw_cube(ax, (6, 3), 1, depth, [1, 2, 3, 4, 5, 6, 9], '0', **solid)
draw_cube(ax, (7, 3), 1, depth, [1, 2, 3, 6, 9], '1', **solid)
draw_cube(ax, (8, 3), 1, depth, [1, 2, 3, 6, 7, 9, 10], '2', **solid)

draw_cube(ax, (6, 2), 1, depth, range(2, 13), '0', **dotted)
draw_cube(ax, (7, 2), 1, depth, [2, 3, 6, 7, 9, 10, 11], '1', **dotted)
draw_cube(ax, (8, 2), 1, depth, [2, 3, 6, 7, 9, 10, 11], '2', **dotted)

draw_cube(ax, (6, 1), 1, depth, [2, 3, 4, 7, 8, 10, 11, 12], '0', **dotted)
draw_cube(ax, (7, 1), 1, depth, [2, 3, 7, 10, 11], '1', **dotted)
draw_cube(ax, (8, 1), 1, depth, [2, 3, 7, 10, 11], '2', **dotted)

# 第三部分
draw_cube(ax, (12, 3), 1, depth, [1, 2, 3, 4, 5, 6, 9], '0', **solid)
draw_cube(ax, (13, 3), 1, depth, [1, 2, 3, 6, 9], '1', **solid)
draw_cube(ax, (14, 3), 1, depth, [2, 3, 7, 10], '3', **solid)

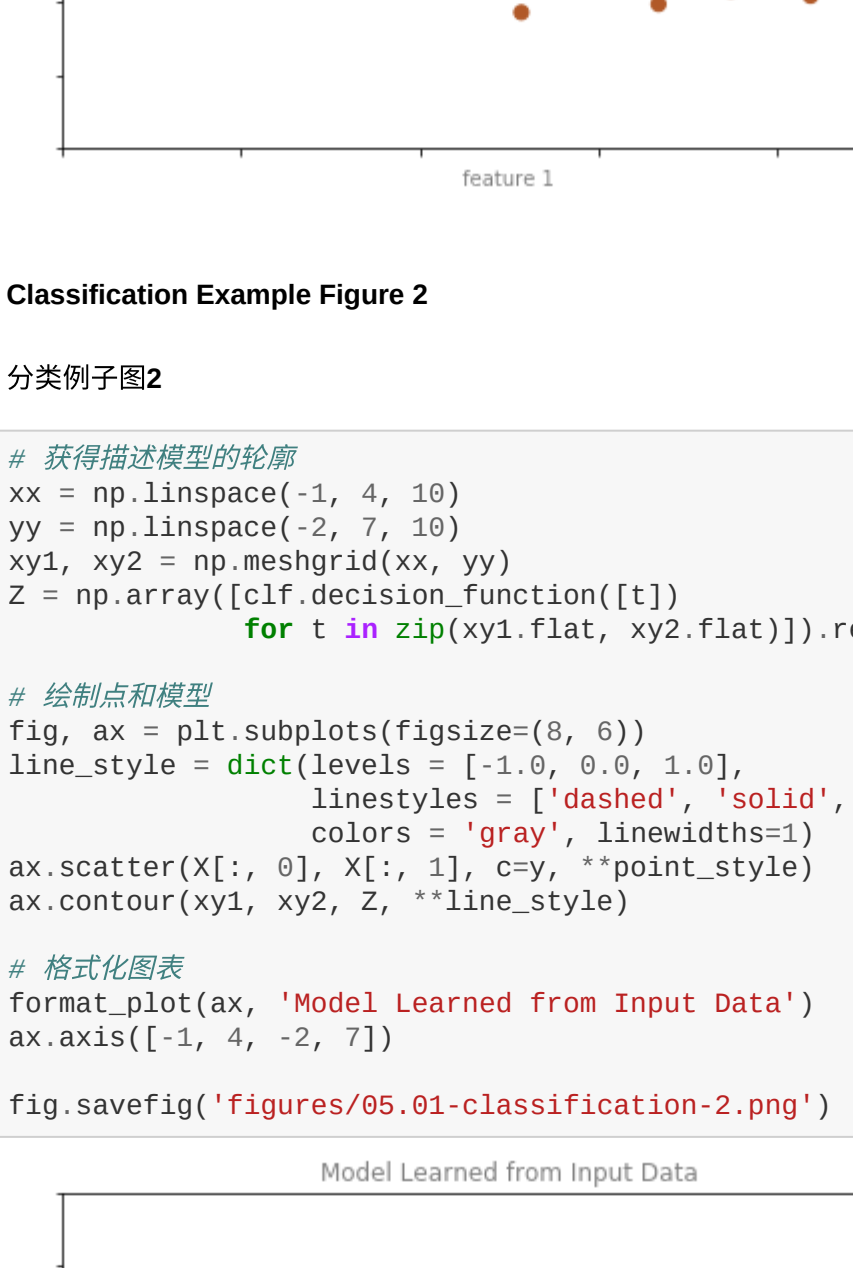
draw_cube(ax, (12, 2), 1, depth, [2, 3, 4], '1', **solid)
draw_cube(ax, (13, 2), 1, depth, [2, 3], '3', **solid)
draw_cube(ax, (14, 2), 1, depth, [2, 3, 7, 10], '3', **solid)

draw_cube(ax, (12, 1), 1, depth, [2, 3, 4], '1', **solid)
draw_cube(ax, (13, 1), 1, depth, [2, 3], '3', **solid)
draw_cube(ax, (14, 1), 1, depth, [2, 3, 7, 10], '4', **solid)

ax.text(5, 2.5, '+', size=12, ha='center', va='center')
ax.text(10.5, 2.5, '=', size=12, ha='center', va='center')
ax.text(11, 4.5, r'$\text{np.arange}(3).reshape(3, 1) + \text{np.arange}(3)$',
        ha='left', size=12, va='bottom')

ax.set_xlim(0, 16)
ax.set_ylim(0.5, 12.5)

fig.savefig('figures/02.05-broadcasting.png')
```



Aggregation and Grouping

聚合与分组

Figures from the chapter on aggregation and grouping

聚合与分组合并的图表

Split-Apply-Combine

分组-应用-合并

```
In [4]: def draw_dataframe(df, loc=None, width=None, ax=None, linestyle=None,
                        textstyle=None):
    loc = loc or [0, 0]
    width = width or 1

    x, y = loc

    ax = plt.gca()

    ncols = len(df.columns) + 1
    nrows = len(df.index) + 1
    dx = dy = width / ncols

    if linestyle is None:
        linestyle = {'color': 'black'}

    if textstyle is None:
        textstyle = {'size': 12}

    textstyle.update({'ha': 'center', 'va': 'center'})

    # 绘制垂直线
    for i in range(ncols + 1):
        plt.plot(2 * [x + i * dx], [y, y + dy * nrows], **linestyle)

    # 绘制水平线
    for i in range(nrows + 1):
        plt.plot([x, x + dx * ncols], 2 * [y + i * dy], **linestyle)

    # 创建索引标签
    for i in range(nrows - 1):
        plt.text(x + 0.5 * dx, y + (i + 0.5) * dy,
                 str(df.index[i+1:][1]), **textstyle)

    # 创建列标签
    for i in range(ncols - 1):
        plt.text(x + (i + 1.5) * dx, y + (nrows - 0.5) * dy,
                 str(df.columns[i+1]), style='italic', **textstyle)

    # 添加索引标签
    if df.index.name:
        plt.text(x + 0.5 * dx, y + (nrows - 0.5) * dy,
                 str(df.index.name), style='italic', **textstyle)

    # 插入数据
    for j in range(nrows - 1):
        for i in range(ncols - 1):
            plt.text(x + (j + 1.5) * dx,
                    y + (i + 0.5) * dy,
                    str(df.values[j+1:][i+1, j]), **textstyle)

#-----绘制图表
import pandas as pd
df = pd.DataFrame({'data': [1, 2, 3, 4, 5, 6],
                   'index': ['A', 'B', 'C', 'A', 'B', 'C']})
df.index.name = 'key'

fig = plt.figure(figsize=(8, 6), facecolor='white')
ax = plt.axes([0, 0, 1, 1])
ax.axis('off')

draw_dataframe(df, [0, 0])

for y, ind in zip([3, 1, -1], 'ABC'):
    split = df[df.index == ind]
    draw_dataframe(split, [2, y])

    sum = pd.DataFrame(split.sum()).T
    sum.index = [ind]
    sum.index.name = 'key'
    sum.columns = ['data']
    draw_dataframe(sum, [4, y + 0.25])

result = df.groupby(df.index).sum()
draw_dataframe(result, [6, 0.75])

style = dict(fontsize=14, ha='center', weight='bold')
plt.text(0.5, 3.6, "Input", **style)
plt.text(2.5, 4.6, "Split", **style)
plt.text(4.5, 4.35, "Apply (sum)", **style)
plt.text(6.5, 2.85, "Combine", **style)

arrowprops = dict(facecolor='black', width=1, headwidth=6)
plt.annotate('', (1.8, 3.6), (1.2, 2.8), arrowprops=arrowprops)
plt.annotate('', (1.8, 3.6), (1.2, 1.75), arrowprops=arrowprops)
plt.annotate('', (1.8, -0.1), (1.2, 0.7), arrowprops=arrowprops)

plt.annotate('', (3.8, 3.6), (3.2, 3.8), arrowprops=arrowprops)
plt.annotate('', (3.8, 3.6), (3.2, 1.75), arrowprops=arrowprops)
plt.annotate('', (3.8, -0.3), (3.2, -0.3), arrowprops=arrowprops)

plt.annotate('', (5.8, 2.8), (5.2, 3.6), arrowprops=arrowprops)
plt.annotate('', (5.8, 1.75), (5.2, 1.75), arrowprops=arrowprops)
plt.annotate('', (5.8, 0.7), (5.2, -0.1), arrowprops=arrowprops)

plt.axis('equal')
plt.ylim(-1.5, 5);

fig.savefig('figures/03.08-split-apply-combine.png')
```



What Is Machine Learning?

下面文章需要用的格式化工具函数

```
In [5]: def format_plot(ax, title):
    ax.xaxis.set_major_formatter(plt.NullFormatter())
    ax.yaxis.set_major_formatter(plt.NullFormatter())
    ax.set_xlabel('feature 1', color='gray')
    ax.set_ylabel('feature 2', color='gray')
    ax.set_title(title, color='gray')
```

Classification Example Figures

分类例子图表

[图表所在正文](#)

The following code generates the figures from the Classification section.

下面的代码生成了分类小节的图表。

```
In [6]: from sklearn.datasets.samples_generator import make_blobs
from sklearn.svm import SVC

# 创建50个独立的点
X, y = make_blobs(n_samples=50, centers=2,
                  random_state=0, cluster_std=0.60)

# 拟合支持向量机模型
clf = SVC(kernel='linear')
clf.fit(X, y)

# 创建新的点用于预测
X2, _ = make_blobs(n_samples=80, centers=2,
                   random_state=0, cluster_std=0.80)
X2 = X2[50:]

# 预测新数据点的标签
y2 = clf.predict(X2)

Classification Example Figure 1

分类例子图1

In [7]: # 绘制数据点
fig, ax = plt.subplots(figsize=(8, 6))
point_style = dict(cmap='Paired', s=50)
ax.scatter(X[:, 0], X[:, 1], c=y, **point_style)

# 格式化工具表
format_plot(ax, 'Input Data')
ax.axis([-1, 4, -2, 7])

fig.savefig('figures/05.01-classification-1.png')
```


Classification Example Figure 2

分类例子图2

```
In [8]: # 获得描述模型的信息
xx = np.linspace(-1, 4, 10)
yy = np.linspace(-2, 7, 10)
xy1, xy2 = np.meshgrid(xx, yy)
Z = np.array([clf.decision_function(t)
              for t in zip(xy1.flat, xy2.flat)]).reshape(xy1.shape)

# 绘制决策边界
fig, ax = plt.subplots(figsize=(8, 6))
line_style = dict(levels = [-1.0, 0.0, 1.0],
                      linestyle = ['dashed', 'solid', 'dashed'],
                      colors = 'gray', linewidth=1)
ax.scatter(X[:, 0], X[:, 1], c=y, **point_style)
ax.contour(xy1, xy2, Z, **line_style)

# 格式化工具表
format_plot(ax, 'Model Learned from Input Data')
ax.axis([-1, 4, -2, 7])

fig.savefig('figures/05.01-classification-2.png')
```


Classification Example Figure 3

分类例子图3

```
In [9]: # 绘制结果
fig, ax = plt.subplots(1, 2, figsize=(16, 6))
fig.subplots_adjust(left=0.0625, right=0.95, wspace=0.1)

ax[0].scatter(X2[:, 0], X2[:, 1], c='gray', **point_style)
ax[0].axis([-1, 4, -2, 7])

ax[1].scatter(X2[:, 0], X2[:, 1], c=y2, **point_style)
ax[1].contour(xy1, xy2, Z, **line_style)
ax[1].axis([-1, 4, -2, 7])

format_plot(ax[0], 'Unknown Data')
format_plot(ax[1], 'Predicted Labels')

fig.savefig('figures/05.01-classification-3.png')
```


Regression Example Figures

回归例子图表

The following code generates the figures from the regression section.

下面的代码生成回归小节的图表。

```
In [10]: from sklearn.linear_model import LinearRegression

# 创建数据点用于回归
rng = np.random.RandomState(1)
X = rng.randn(200, 2)
y = np.dot(X, [-2, 1]) + 0.1 * rng.randn(X.shape[0])

# 拟合回归模型
model = LinearRegression()
model.fit(X, y)

# 创建新数据点进行预测
X2 = rng.randn(100, 2)

# 预测标签
y2 = model.predict(X2)

Regression Example Figure 1

回归例子图1

In [11]: # 绘制数据点
fig, ax = plt.subplots()
points = ax.scatter(X[:, 0], X[:, 1], c=y, s=50,
                   cmap='viridis')

# 格式化工具表
format_plot(ax, 'Input Data')
ax.axis([-1, 4, -3, 3])

fig.savefig('figures/05.01-regression-1.png')
```


Regression Example Figure 2

回归例子图2

```
In [12]: from mpl_toolkits.mplot3d.art3d import Line3DCollection

points = np.hstack([X, y[:, None]]).reshape(-1, 1, 3)
segments = np.hstack([points, points])
segments[:, 0, 2] = -8

# 在3D中绘制点
fig = plt.figure()
ax = fig.add_subplot(111, projection='3d')
ax.scatter(X[:, 0], X[:, 1], c=y, s=35,
           cmap='viridis')
ax.add_collection3d(Line3DCollection(segments, colors='gray', alpha=0.2))
ax.scatter(X[:, 0], X[:, 1], -8 + np.zeros(X.shape[0]), c=y, s=10,
           cmap='viridis')

# 格式化工具表
ax.patch.set_facecolor('white')
ax.view_init(elev=20, azim=-70)
ax.set_zlim3d(-8, 8)
ax.xaxis.set_major_formatter(plt.NullFormatter())
ax.yaxis.set_major_formatter(plt.NullFormatter())
ax.zaxis.set_major_formatter(plt.NullFormatter())

# 隐藏坐标轴 (是否有更好的办法?)
ax.w_xaxis.line.set_visible(False)
ax.w_yaxis.line.set_visible(False)
ax.w_zaxis.line.set_visible(False)
for tick in ax.w_xaxis.get_ticklines():
    tick.set_visible(False)
for tick in ax.w_yaxis.get_ticklines():
    tick.set_visible(False)
for tick in ax.w_zaxis.get_ticklines():
    tick.set_visible(False)

fig.savefig('figures/05.01-regression-2.png')
```


Regression Example Figure 3

回归例子图3

```
In [13]: from matplotlib.collections import LineCollection

# 绘制数据点
fig, ax = plt.subplots()
pts = ax.scatter(X[:, 0], X[:, 1], c=y, s=50,
                 cmap='viridis', zorder=2)

# 计算和绘制模型拟合面
xx, yy = np.meshgrid(np.linspace(-4, 4),
                     np.linspace(-3, 3))
Xfit = np.vstack([xx.ravel(), yy.ravel()]).T
yfit = model.predict(Xfit)
zz = yfit.reshape(xx.shape)
ax.pcolorfast([-4, 4], [-3, 3], zz, alpha=0.5,
              cmap='viridis', norm=pts.norm, zorder=1)

# 格式化工具表
format_plot(ax, 'Input Data with Linear Fit')
ax.axis([-4, 4, -3, 3])

fig.savefig('figures/05.01-regression-3.png')
```


Regression Example Figure 4

回归例子图4

```
In [14]: # 绘制模型拟合面
fig, ax = plt.subplots(1, 2, figsize=(16, 6))
fig.subplots_adjust(left=0.0625, right=0.95, wspace=0.1)

ax[0].scatter(X2[:, 0], X2[:, 1], c='gray', s=50)
ax[0].axis([-4, 4, -3, 3])

ax[1].scatter(X2[:, 0], X2[:, 1], c=y2, s=50,
              cmap='viridis', norm=pts.norm)
ax[1].axis([-4, 4, -3, 3])

# 格式化工具表
format_plot(ax[0], 'Unknown Data')
format_plot(ax[1], 'Predicted Labels')

fig.savefig('figures/05.01-regression-4.png')
```


Clustering Example Figures

[图表所在正文](#)

The following code generates the figures from the clustering section.

下面的代码生成聚类小节的图表。

```
In [15]: from sklearn.datasets.samples_generator import make_blobs

# 创建50个独立的点
X, y = make_blobs(n_samples=100, centers=4,
                  random_state=42, cluster_std=1.5)

# 拟合均值聚类模型
model = KMeans(4, random_state=0)
y = model.fit_predict(X)

Clustering Example Figure 1

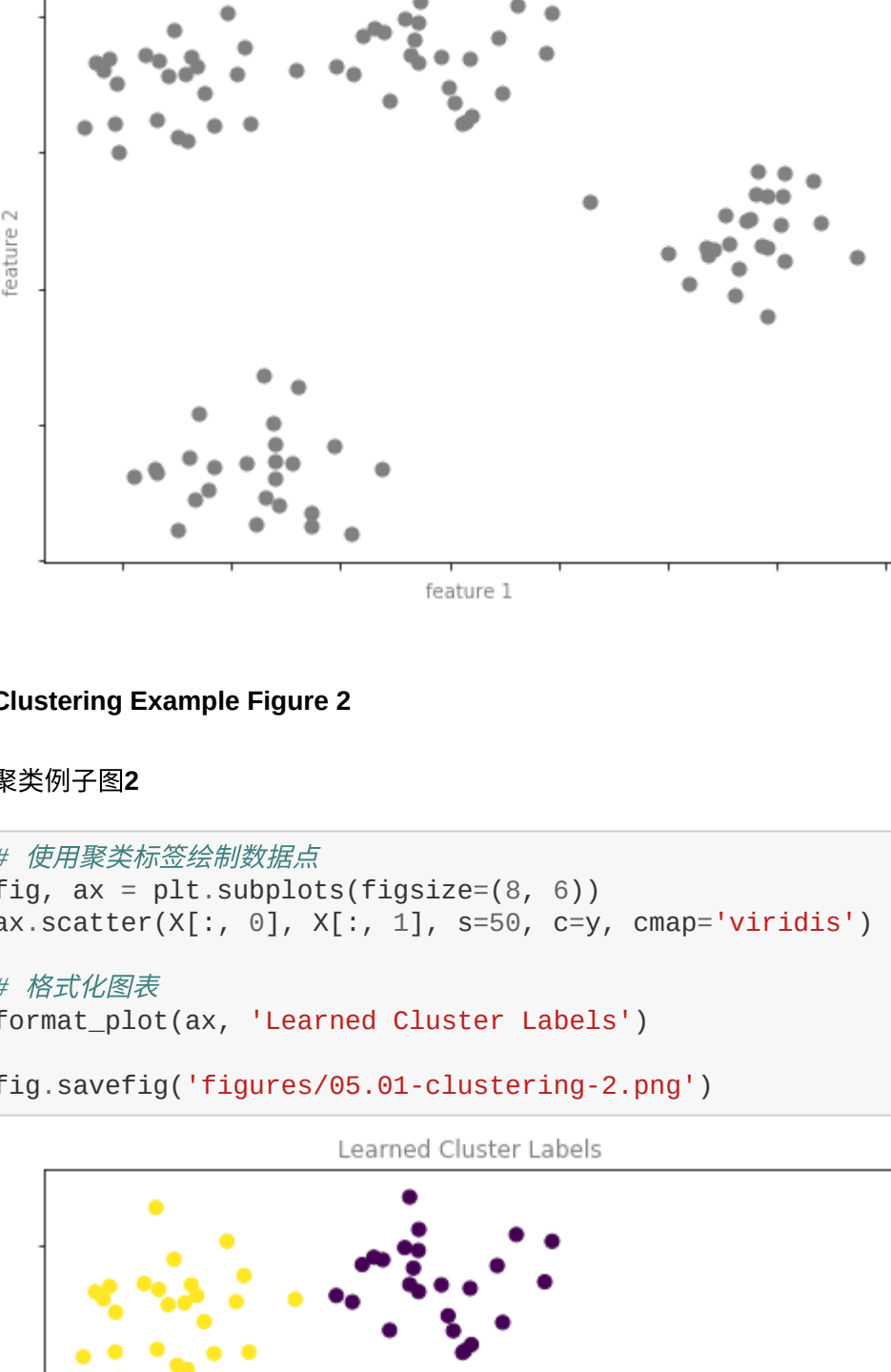
聚类例子图1
```



```
In [16]: # 绘制输入数据
fig, ax = plt.subplots(figsize=(8, 6))
ax.scatter(X[:, 0], X[:, 1], s=50, color='gray')

# 格式化图表
format_plot(ax, 'Input Data')

fig.savefig('figures/05_01-clustering-1.png')
```



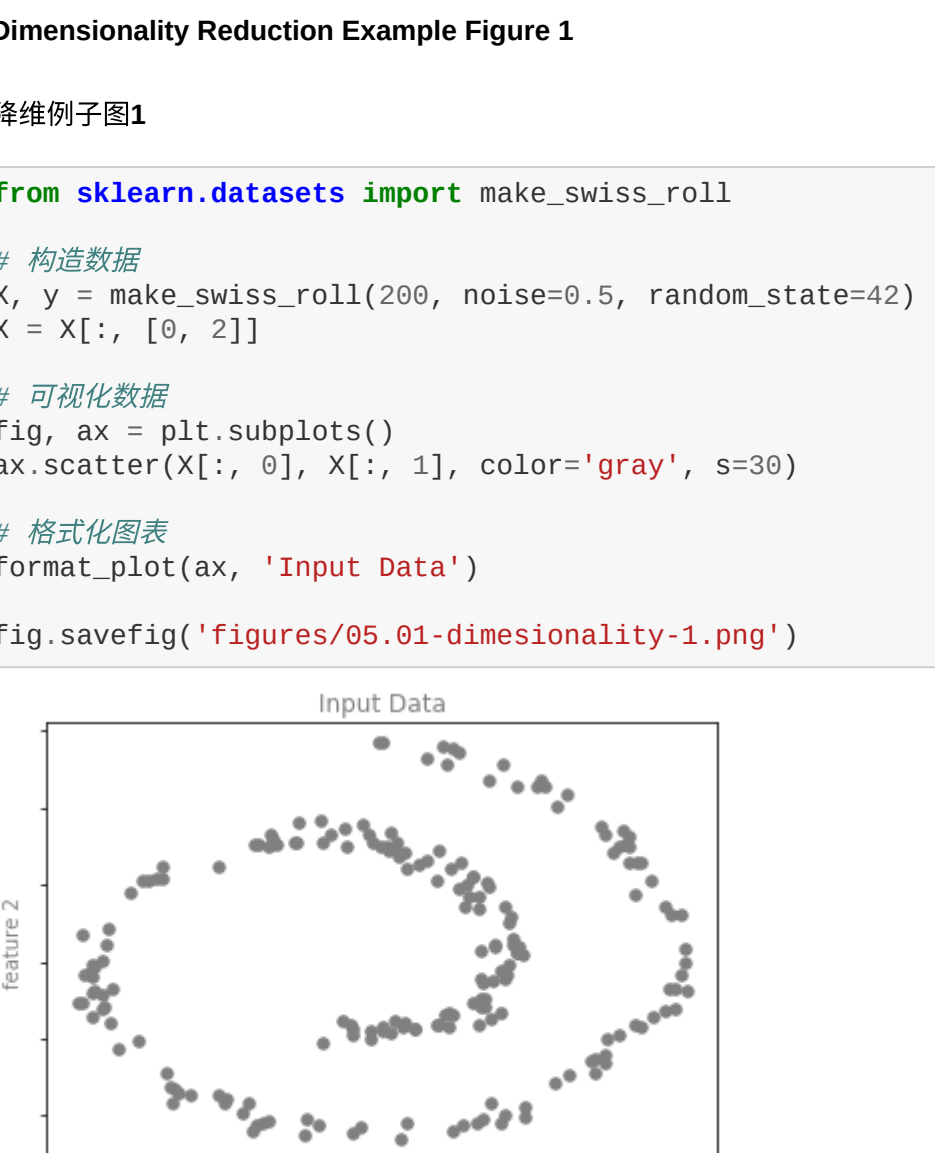
Clustering Example Figure 2

聚类例子图2

```
In [17]: # 使用聚类标签绘制数据点
fig, ax = plt.subplots(figsize=(8, 6))
ax.scatter(X[:, 0], X[:, 1], s=50, c=y, cmap='viridis')

# 格式化图表
format_plot(ax, 'Learned Cluster Labels')

fig.savefig('figures/05_01-clustering-2.png')
```



Dimensionality Reduction Example Figures

降维例子图表

图表所在正文

The following code generates the figures from the dimensionality reduction section.

下面代码生成降维小节图表。

Dimensionality Reduction Example Figure 1

降维例子图1

```
In [18]: from sklearn.datasets import make_swiss_roll

# 构造数据
X, y = make_swiss_roll(200, noise=0.5, random_state=42)
X = X[:, [0, 2]]

# 可视化数据
fig, ax = plt.subplots()
ax.scatter(X[:, 0], X[:, 1], color='gray', s=30)

# 格式化图表
format_plot(ax, 'Input Data')

fig.savefig('figures/05_01-dimensionality-1.png')
```



Dimensionality Reduction Example Figure 2

降维例子图2

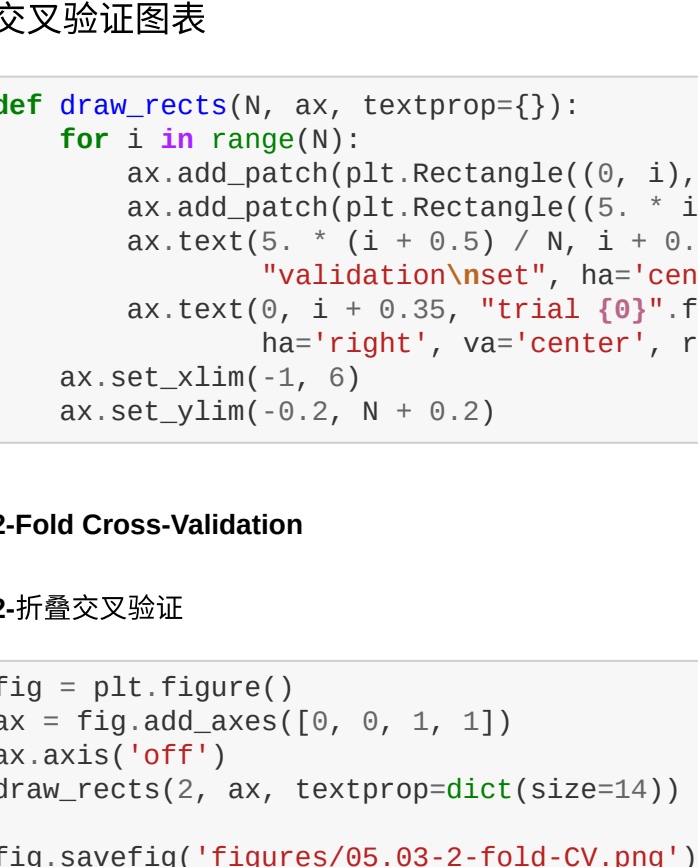
```
In [19]: from sklearn.manifold import Isomap

model = Isomap(n_neighbors=8, n_components=1)
y_fit = model.fit_transform(X).ravel()

# 可视化数据
fig, ax = plt.subplots()
pts = ax.scatter(X[:, 0], X[:, 1], c=y_fit, cmap='viridis', s=30)
cb = fig.colorbar(pts, ax=ax)

# 格式化图表
format_plot(ax, 'Learned Latent Parameter')
cb.set_ticks([])
cb.set_label('Latent Variable', color='gray')

fig.savefig('figures/05_01-dimensionality-2.png')
```



Introducing Scikit-Learn

Scikit-Learn 介绍

Features and Labels Grid

特征和标签网格

The following is the code generating the diagram showing the features matrix and target array.

下面代码生成特征矩阵和目标数组部分的图表。

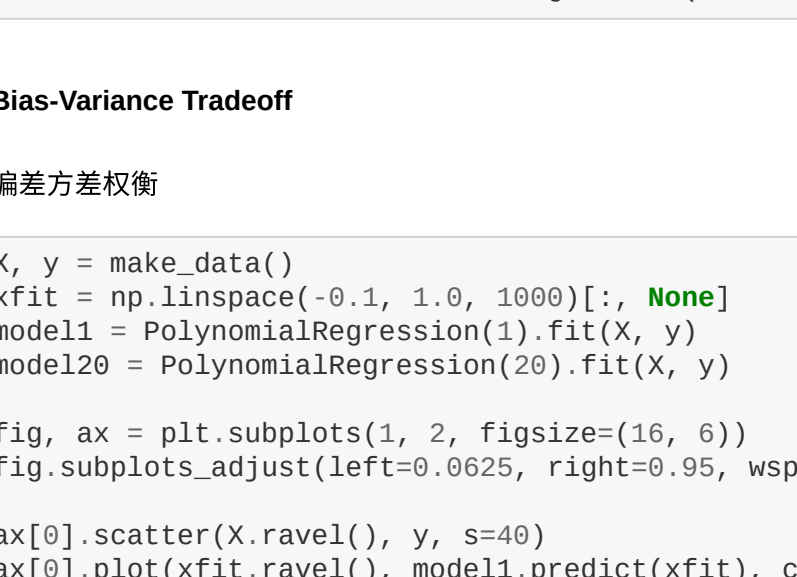
```
In [20]: fig = plt.figure(figsize=(6, 4))
ax = fig.add_axes([0, 0, 1, 1])
ax.axis('off')
ax.axis('equal')

# 绘制特征矩阵
ax.vlines(range(6), ymin=0, ymax=9, lw=1)
ax.hlines(range(10), xmin=0, xmax=5, lw=1)
font_prop = dict(size=12, family='monospace')
ax.text(-1, -1, "Feature Matrix (X$X$)", size=14)
ax.text(0.1, -0.3, r"$n\_features \rightarrow$", **font_prop)
ax.text(-0.1, 0.1, r"$\rightarrow n\_samples$", rotation=90,
       va='top', ha='right', **font_prop)

# 绘制目标向量
ax.vlines(range(8, 10), ymin=0, ymax=9, lw=1)
ax.hlines(range(10), xmin=0, xmax=9, lw=1)
ax.text(7, -1, "Target Vector (y$y$)", size=14)
ax.text(7, 0.1, r"$\rightarrow n\_samples$", rotation=90,
       va='top', ha='right', **font_prop)

ax.set_xlim(10, -2)

fig.savefig('figures/05_02-samples-features.png')
```



Hyperparameters and Model Validation

超参数和模型验证

Cross-Validation Figures

交叉验证图表

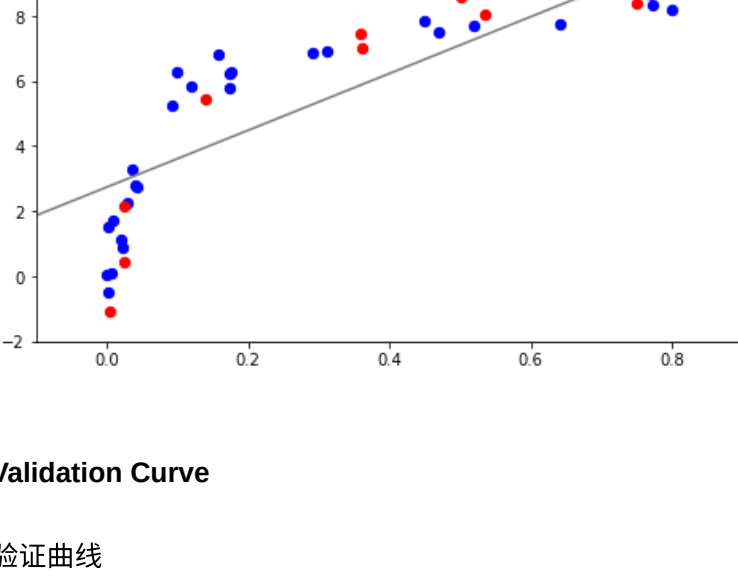
```
In [21]: def draw_rects(N, ax, textprop={}):
    for i in range(N):
        ax.add_patch(plt.Rectangle((0, 1), 5, 0.7, fc='white'))
        ax.add_patch(plt.Rectangle((5, -2), 4, 0.7, fc='lightgray'))
        ax.text(5, -1 + 0.5 / N, 1 + 0.35,
               "Validation\inset", ha='center', va='center', **textprop)
        ax.text(0, 1 + 0.35, "trial {0}".format(N - i),
               ha='right', va='center', rotation=90, **textprop)
    ax.set_xlim(-3, 6)
    ax.set_ylim(-0.2, N + 0.2)
```

2-Fold Cross-Validation

2-折交叉验证

```
In [22]: fig = plt.figure()
ax = fig.add_axes([0, 0, 1, 1])
ax.axis('off')
draw_rects(2, ax, textprop=dict(size=14))

fig.savefig('figures/05_03-2-fold-CV.png')
```

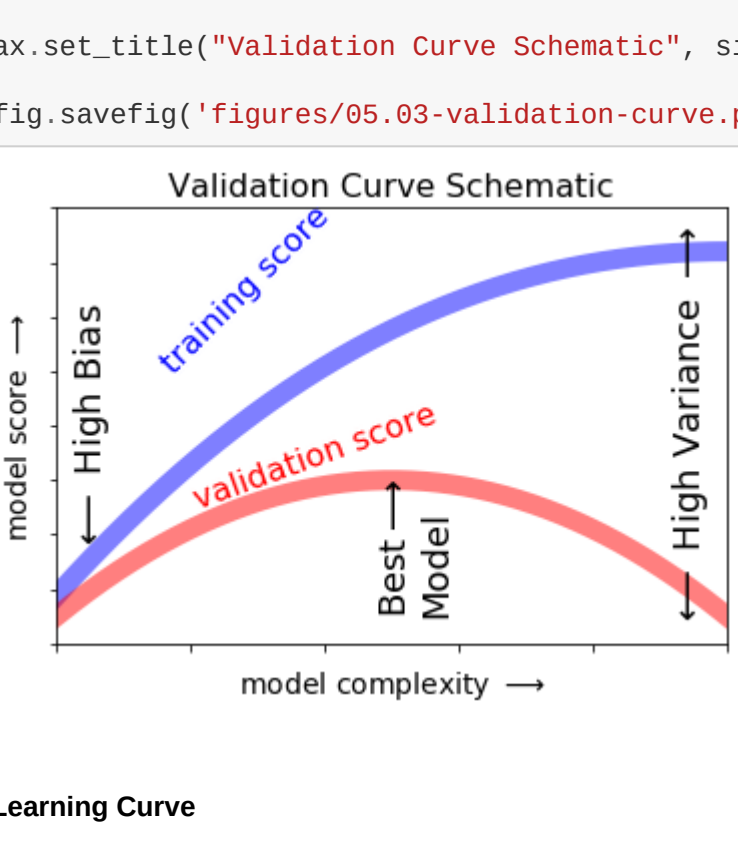


5-Fold Cross-Validation

5-折交叉验证

```
In [23]: fig = plt.figure()
ax = fig.add_axes([0, 0, 1, 1])
ax.axis('off')
draw_rects(5, ax, textprop=dict(size=10))

fig.savefig('figures/05_03-5-fold-CV.png')
```



Overfitting and Underfitting

过拟合与欠拟合

```
In [24]: import numpy as np

def make_data(N=30, err=0.8, rseed=1):
    # 随机产生数据样本
    rng = np.random.RandomState(rseed)
    X = rng.rand(N, 1) ** 2
    y = 10 - 1. / (X.ravel() + 0.1)
    if err > 0:
        y += err * rng.randn(N)
    return X, y
```

```
In [25]: from sklearn.preprocessing import PolynomialFeatures
from sklearn.linear_model import LinearRegression
from sklearn.pipeline import make_pipeline

def PolynomialRegression(degree=2, **kwargs):
    return make_pipeline(PolynomialFeatures(degree),
                        LinearRegression(**kwargs))
```

Bias-Variance Tradeoff

偏差方差权衡

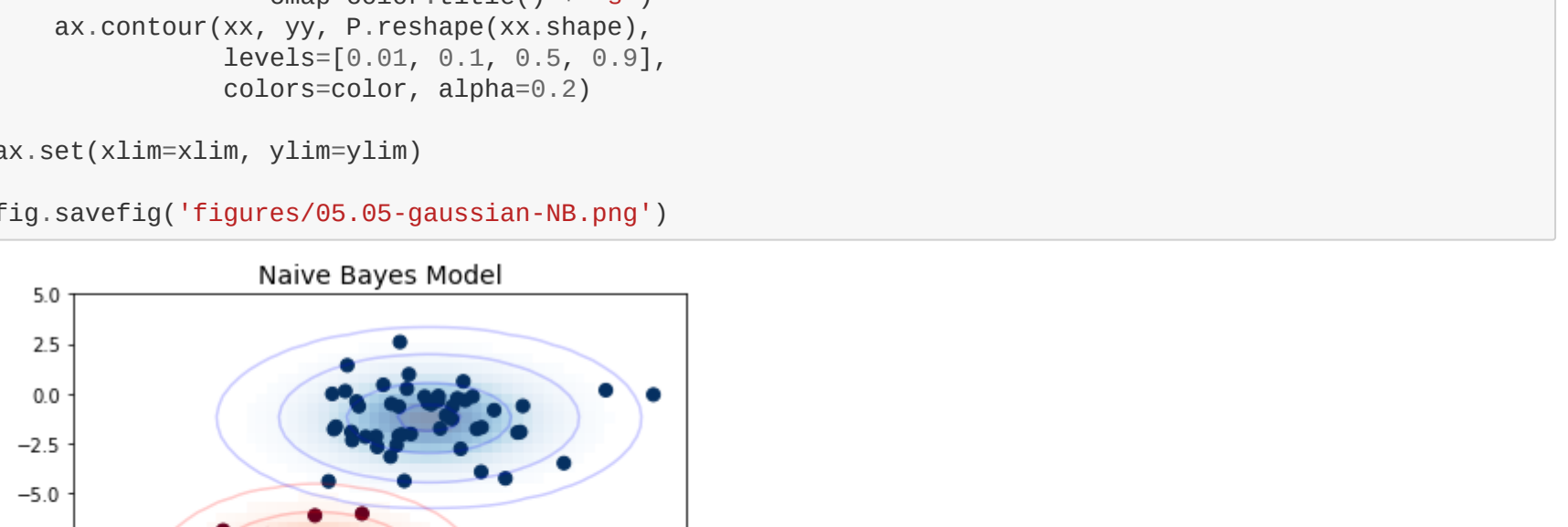
```
In [26]: X, y = make_data()
xfit = np.linspace(-0.1, 1.0, 1000)[1:, None]
model1 = PolynomialRegression(1).fit(X, y)
model20 = PolynomialRegression(20).fit(X, y)

fig, ax = plt.subplots(1, 2, figsize=(16, 6))
ax.subplots_adjust(left=0.0625, right=0.95, wspace=0.1)

ax[0].scatter(X.ravel(), y, s=40)
ax[0].plot(xfit.ravel(), model1.predict(xfit), color='gray')
ax[0].axis([-0.1, 1.0, -2, 14])
ax[0].set_title('High-bias model: Underfits the data', size=14)
ax[0].text(0.02, 0.98, "training score: $R^2$ = {0:.2f}".format(model1.score(X, y)),
         ha='left', va='top', transform=ax[0].transAxes, size=14, color='blue')
ax[0].text(0.02, 0.91, "validation score: $R^2$ = {0:.2f}".format(model1.score(X2, y2)),
         ha='left', va='top', transform=ax[0].transAxes, size=14, color='red')

ax[1].scatter(X.ravel(), y, s=40)
ax[1].plot(xfit.ravel(), model20.predict(xfit), color='gray')
ax[1].axis([-0.1, 1.0, -2, 14])
ax[1].set_title('High-variance model: Overfits the data', size=14)
ax[1].text(0.02, 0.98, "training score: $R^2$ = {0:.2f}".format(model20.score(X, y)),
         ha='left', va='top', transform=ax[1].transAxes, size=14, color='blue')
ax[1].text(0.02, 0.91, "validation score: $R^2$ = {0:.2f}".format(model20.score(X2, y2)),
         ha='left', va='top', transform=ax[1].transAxes, size=14, color='red')

fig.savefig('figures/05_03-bias-variance.png')
```



Bias-Variance Tradeoff Metrics

偏差方差权衡指标

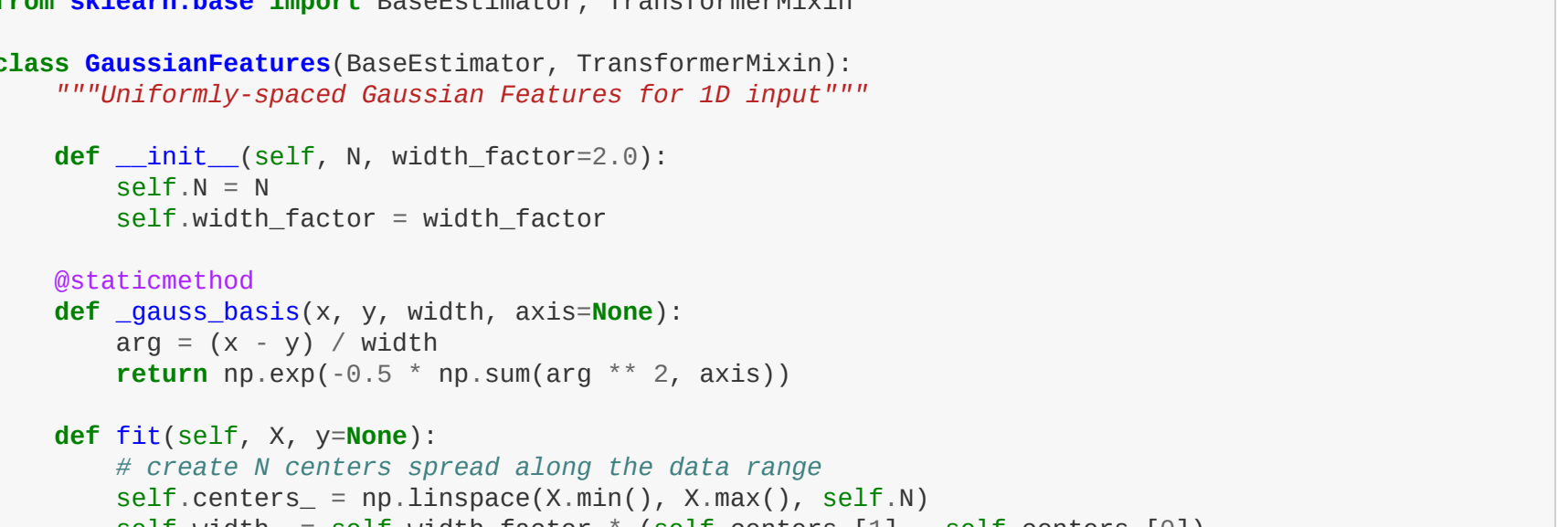
```
In [27]: fig = plt.subplots(1, 2, figsize=(16, 6))
fig.subplots_adjust(left=0.0625, right=0.95, wspace=0.1)

X2, y2 = make_data(10, rseed=42)

ax[0].scatter(X.ravel(), y, s=40, c='blue')
ax[0].plot(xfit.ravel(), model1.predict(xfit), color='gray')
ax[0].axis([-0.1, 1.0, -2, 14])
ax[0].set_title('High-bias model: Underfits the data', size=14)
ax[0].text(0.02, 0.98, "training score: $R^2$ = {0:.2f}".format(model1.score(X, y)),
         ha='left', va='top', transform=ax[0].transAxes, size=14, color='blue')
ax[0].text(0.02, 0.91, "validation score: $R^2$ = {0:.2f}".format(model1.score(X2, y2)),
         ha='left', va='top', transform=ax[0].transAxes, size=14, color='red')

ax[1].scatter(X.ravel(), y, s=40, c='blue')
ax[1].plot(xfit.ravel(), model20.predict(xfit), color='gray')
ax[1].axis([-0.1, 1.0, -2, 14])
ax[1].set_title('High-variance model: Overfits the data', size=14)
ax[1].text(0.02, 0.98, "training score: $R^2$ = {0:.2f}".format(model20.score(X, y)),
         ha='left', va='top', transform=ax[1].transAxes, size=14, color='blue')
ax[1].text(0.02, 0.91, "validation score: $R^2$ = {0:.2f}".format(model20.score(X2, y2)),
         ha='left', va='top', transform=ax[1].transAxes, size=14, color='red')

fig.savefig('figures/05_03-bias-variance-2.png')
```



Validation Curve

验证曲线

```
In [28]: X = np.linspace(0, 1, 1000)
y1 = (X - 0.5) ** 2
y2 = y1 - 0.33 + np.exp(-X - 1)

fig, ax = plt.subplots()
ax.plot(X, y2, lw=10, alpha=0.5, color='blue')
ax.plot(X, y1, lw=10, alpha=0.5, color='red')

ax.text(0.15, 0.2, "training score", rotation=45, size=16, color='blue')
ax.text(0.2, -0.05, "validation score", rotation=20, size=16, color='red')

ax.text(0.02, 0.1, r"$\rightarrow$ High Bias", size=18, rotation=90, ha='center')
ax.text(0.48, 0.1, r"$\rightarrow$ High Variance", size=18, rotation=90, ha='center')
ax.text(0.48, -0.32, "Best Fit", size=18, rotation=90, va='center')

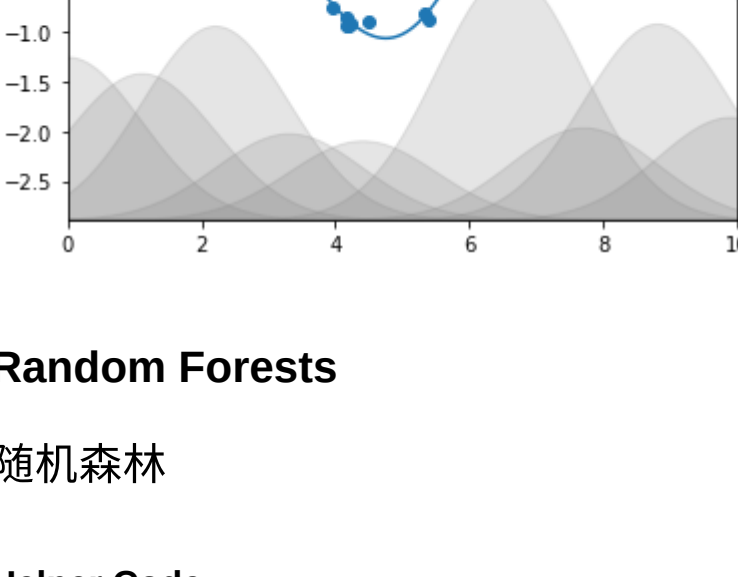
ax.set_xlim(0, 1)
ax.set_ylim(-0.3, 0.5)

ax.set_xlabel(r"model complexity $\rightarrow$", size=14)
ax.set_ylabel(r"model score $\rightarrow$", size=14)

ax.xaxis.set_major_formatter(plt.NullFormatter())
ax.yaxis.set_major_formatter(plt.NullFormatter())

ax.set_title("Validation Curve Schematic", size=16)

fig.savefig('figures/05_03-validation-curve.png')
```



Learning Curve

学习曲线

```
In [29]: N = np.linspace(0, 1, 1000)
y1 = 0.75 + 0.2 * np.exp(-4 * N)
y2 = 0.7 - 0.6 * np.exp(-4 * N)

fig, ax = plt.subplots()
ax.plot(X, y1, lw=10, alpha=0.5, color='blue')
ax.plot(X, y2, lw=10, alpha=0.5, color='red')

ax.text(0.2, 0.8, "training score", rotation=10, size=16, color='blue')
ax.text(0.2, 0.5, "validation score", rotation=30, size=16, color='red')

ax.text(0.02, 0.98, "Good Fit", size=18, rotation=90, ha='right', va='center')
ax.text(0.02, 0.57, r"$\rightarrow$ High Variance", size=18, rotation=90, va='center')
ax.text(0.48, 0.1, r"$\rightarrow$ High Bias", size=18, rotation=90, va='center')

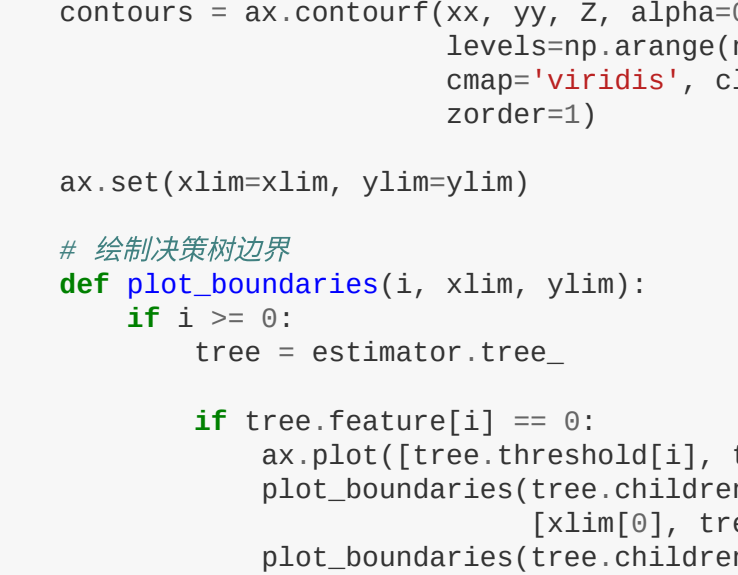
ax.set_xlim(0, 1)
ax.set_ylim(0, 1)

ax.set_xlabel(r"training set size $\rightarrow$", size=14)
ax.set_ylabel(r"model score $\rightarrow$", size=14)

ax.xaxis.set_major_formatter(plt.NullFormatter())
ax.yaxis.set_major_formatter(plt.NullFormatter())

ax.set_title("Learning Curve Schematic", size=16)

fig.savefig('figures/05_03-learning-curve.png')
```



Gaussian Naive Bayes

高斯朴素贝叶斯

Gaussian Naive Bayes Example

高斯朴素贝叶斯例子

图表所在正文

```
In [30]: from sklearn.datasets import make_blobs
X, y = make_blobs(100, 2, centers=2, random_state=2, cluster_std=1.5)

fig, ax = plt.subplots()
ax.scatter(X[:, 0], X[:, 1], c=y, s=50, cmap='RdBu')
ax.set_title('Naive Bayes Model', size=14)

xlim = (-8, 8)
ylim = (-15, 5)

xg = np.linspace(xlim[0], xlim[1], 60)
yg = np.linspace(ylim[0], ylim[1], 40)
xx, yy = np.meshgrid(xg, yg)
Xgrid = np.vstack([xx.ravel(), yy.ravel()]).T

for label, color in enumerate(['red', 'blue']):
    mask = (y == label)
    mu, std = X[mask].mean(0), X[mask].std(0)
    P = np.exp(-0.5 * (Xgrid - mu) ** 2 / std ** 2).prod(1)
    Pm = np.ma.masked_array(P, P < 0.03)
    ax.pcolormesh(xg, yg, Pm.reshape((xg.shape, yg.shape)), alpha=0.5,
                  cmap=cm.get_cmap('RdBu'))
    ax.contour(xx, yy, P.reshape((xx.shape, yy.shape)),
               levels=[0.01, 0.1, 0.5, 0.9],
               colors=color, alpha=0.2)

ax.set(xlim=xlim, ylim=ylim)

fig.savefig('figures/05_05-gaussian-NB.png')
```



Linear Regression

线性回归

Gaussian Basis Functions

高斯基本函数

图表所在正文

```
In [31]: from sklearn.pipeline import make_pipeline
from sklearn.linear_model import LinearRegression

from sklearn.base import BaseEstimator, TransformerMixin

class GaussianFeatures(BaseEstimator, TransformerMixin):
    """Uniformly-scaled Gaussian Features for 1D input"""

    def __init__(self, N, width_factor=2.0):
        self.N = N
        self.width_factor = width_factor

    @staticmethod
    def _gauss_basis(x, y, width, axis=None):
        arg = (x - y) / width
        return np.exp(-0.5 * np.sum(arg ** 2, axis))

    def fit(self, X, y=None):
        # create N centers spread along the data range
        self.centers = np.linspace(X.min(), X.max(), self.N)
        self.width_factor = (self.centers[1] - self.centers[0])
        return self

    def transform(self, X):
        return self._gauss_basis(X[:, :, np.newaxis], self.centers_,
                                self.width_, axis=1)

rng = np.random.RandomState(1)
x = 10 * rng.rand(50)
y = np.sin(x) + 0.1 * rng.randn(50)
xfit = np.linspace(0, 10, 1000)

gauss_model = make_pipeline(GaussianFeatures(10, 1.0),
                           LinearRegression())
gauss_model.fit(xfit, np.newaxis[:, y])
yfit = gauss_model.predict(xfit[:, np.newaxis])

gf = gauss_model.named_steps['gaussianfeatures']
lm = gauss_model.named_steps['linearregression']

fig, ax = plt.subplots()

for i in range(10):
    selector = np.zeros(xfit)
    selector[i] = 1
    Xfit = gf.transform(xfit[:, None]) * selector
    yfit = lm.predict(Xfit)
    ax.fill_between(xfit, yfit.min(), yfit, color='gray', alpha=0.2)

ax.scatter(x, y)
ax.plot(xfit, gauss_model.predict(xfit[:, np.newaxis]))
ax.set_xlim(0, 10)
ax.set_ylim(yfit.min(), 1.5)

fig.savefig('figures/05_06-gaussian-basis.png')
```


Random Forests

随机森林

Helper Code

工具代码

The following will create a module `helpers_05_08.py` which contains some tools used in [In-Denith Decision Trees](#) and [Random Forests](#).

下面代码创建模块 `helpers_05_08.py`，包含一些在[深入：决策树和随机森林](#)中用到的工具。

```
In [32]: %file helpers_05_08.py

import numpy as np
import matplotlib.pyplot as plt
from sklearn.tree import DecisionTreeClassifier
from ipywidgets import interact
```

```
def visualize_tree(estimator, X, y, boundaries=True,
                  xlim=None, ylim=None, ax=None):
    ax = ax or plt.gca()

    # 绘制训练数据
    ax.scatter(X[:, 0], X[:, 1], c=y, s=30, cmap='viridis',
              zorder=3)
    ax.axis('tight')
    ax.axis('equal')
    if xlim is None:
        xlim = ax.get_xlim()
    if ylim is None:
        ylim = ax.get_ylim()

    # 拟合评估器
    estimator.fit(X, y)
    xx, yy = np.meshgrid(np.linspace(*xlim, num=200),
                        np.linspace(*ylim, num=200))
    Z = estimator.predict(np.c_[xx.ravel(), yy.ravel()])

    # 将结果放到颜色图表中
    n_classes = len(np.unique(y))
    Z = Z.reshape(xx.shape)
    contours = ax.contourf(xx, yy, Z, alpha=0.3,
                          levels=np.arange(n_classes + 1) - 0.5,
                          cmap=cm.get_cmap('viridis'), clim=(y.min(), y.max()),
                          zorder=1)

    ax.set(xlim=xlim, ylim=ylim)

    # 绘制决策树边界
    def plot_boundaries(i, xlim, ylim):
        if i > 0:
            tree = estimator.tree_

            if tree.feature[i] == 0:
                ax.plot([tree.threshold[i], tree.threshold[i]], ylim, '-k', zorder=2)
                plot_boundaries(tree.children_left[i], xlim,
                               [xlim[0], tree.threshold[i]], ylim)
                plot_boundaries(tree.children_right[i], xlim,
                               [tree.threshold[i], xlim[1]], ylim)
            elif tree.feature[i] == 1:
                ax.plot([tree.threshold[i], tree.threshold[i]], ylim, '-k', zorder=2)
                plot_boundaries(tree.children_left[i], xlim,
                               [ylim[0], tree.threshold[i]])
                plot_boundaries(tree.children_right[i], xlim,
                               [tree.threshold[i], ylim[1]])

    if boundaries:
        plot_boundaries(0, xlim, ylim)

def plot_tree_interactive(X, y):
    def interactive_tree(depth=5):
        clf = DecisionTreeClassifier(max_depth=depth, random_state=0)
        visualize_tree(clf, X, y)

    return interact(interactive_tree, depth=[1, 5])

def randomized_tree_interactive(X, y):
    xlim = [X[:, 0].min(), X[:, 0].max()]
    ylim = [X[:, 1].min(), X[:, 1].max()]

    def fit_randomized_tree(random_state=0):
        clf = DecisionTreeClassifier(max_depth=15)
        i = np.arange(len(y))
        rng = np.random.RandomState(random_state)
        rng.shuffle(i)
        visualize_tree(clf, X[i[None]], y[i[None]], boundaries=False,
                      xlim=xlim, ylim=ylim)

    interact(fit_randomized_tree, random_state=[0, 100])
```

Overwriting `helpers_05_08.py`

Decision Tree Example

决策树例子


```
In [33]: fig = plt.figure(figsize=(10, 4))
ax = fig.add_axes([0, 0, 0.8, 1], frameon=False, xticks=[], yticks=[])
ax.set_title('Example Decision Tree: Animal Classification', size=24)

def text(ax, x, y, t, size=20, **kwargs):
    ax.text(x, y, t, size=size, **kwargs)
    hax=dict(ha='center', va='center', size=size,
            bbox=dict(boxstyle='round', ec='k', fc='w'), **kwargs)

text(ax, 0.5, 0.9, "How big is the animal?", 20)
text(ax, 0.3, 0.6, "Does the animal have horns?", 18)
text(ax, 0.7, 0.6, "Does the animal have two legs?", 18)
text(ax, 0.12, 0.3, "Are the horns longer than 10cm?", 14)
text(ax, 0.18, 0.3, "Is the animal wearing a collar?", 14)
text(ax, 0.62, 0.3, "Does the animal have wings?", 14)
text(ax, 0.88, 0.3, "Does the animal have a tail?", 14)

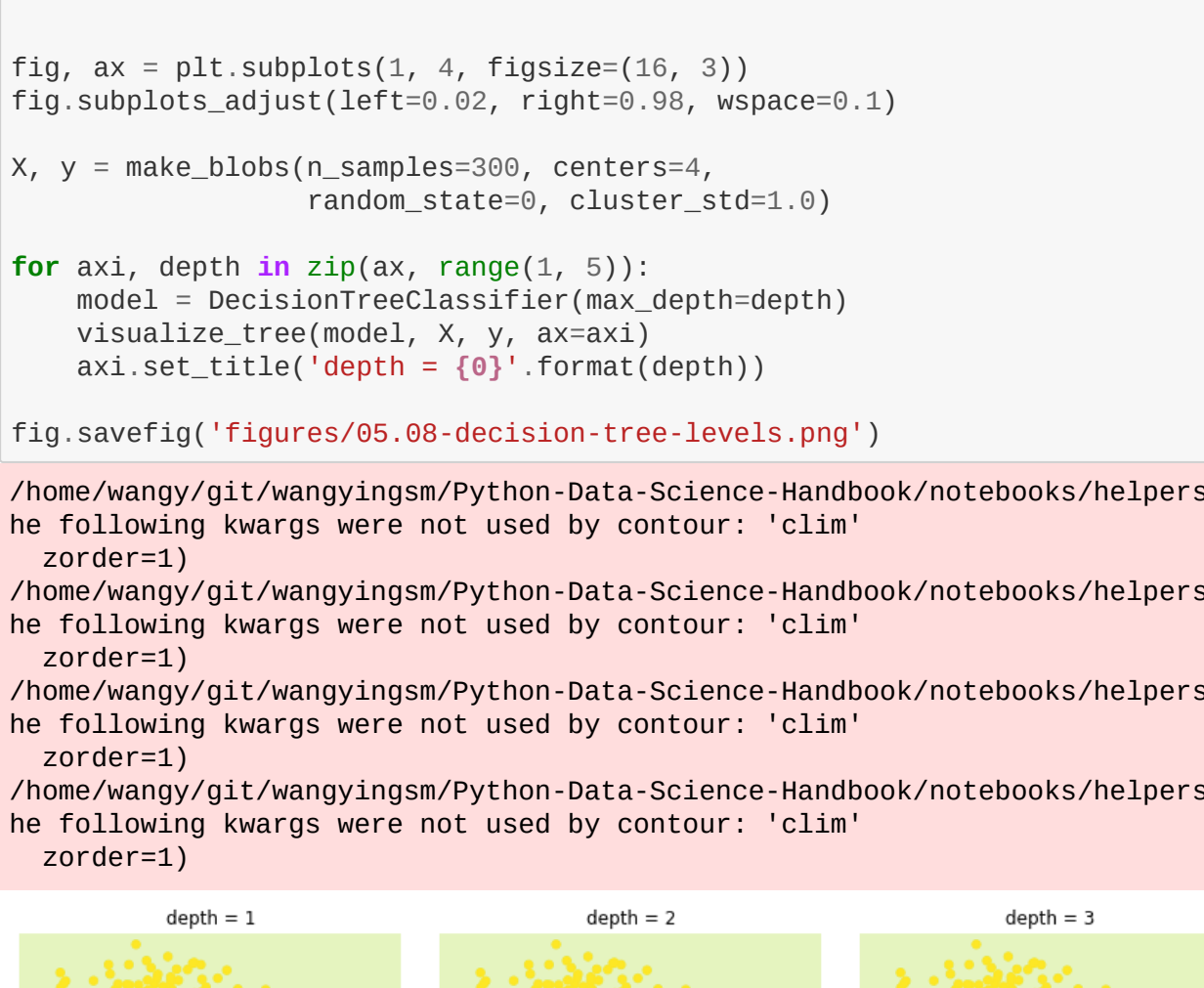
text(ax, 0.4, 0.75, "> 1m", 12, alpha=0.4)
text(ax, 0.6, 0.75, "< 1m", 12, alpha=0.4)

text(ax, 0.21, 0.45, "yes", 12, alpha=0.4)
text(ax, 0.34, 0.45, "no", 12, alpha=0.4)

text(ax, 0.66, 0.45, "yes", 12, alpha=0.4)
text(ax, 0.79, 0.45, "no", 12, alpha=0.4)

ax.plot([0.3, 0.5, 0.7], [0.6, 0.9, 0.6], '-k')
ax.plot([0.12, 0.3, 0.38], [0.3, 0.6, 0.3], '-k')
ax.plot([0.62, 0.7, 0.88], [0.3, 0.6, 0.3], '-k')
ax.plot([0.8, 0.12, 0.3, 0.28], [0.6, 0.3, 0.1], '-k')
ax.plot([0.18, 0.38, 0.48], [0.6, 0.3, 0.1], '-k')
ax.plot([0.52, 0.62, 0.72], [0.6, 0.3, 0.1], '-k')
ax.plot([0.8, 0.88, 1.0], [0.6, 0.3, 0.1], '-k')
ax.axis([0, 1, 0, 1])
fig.savefig('figures/05_08-decision-tree.png')
```

Example Decision Tree: Animal Classification



Decision Tree Levels

决策树层次

```
In [34]: from helpers_05_08 import visualize_tree
from sklearn.tree import DecisionTreeClassifier
from sklearn.datasets import make_blobs

fig, ax = plt.subplots(1, 4, figsize=(16, 3))
fig.subplots_adjust(left=0.02, right=0.98, wspace=0.1)

X, y = make_blobs(n_samples=300, centers=4,
                  random_state=0, cluster_std=1.0)

for axi, depth in zip(ax, range(1, 5)):
    model = DecisionTreeClassifier(max_depth=depth)
    visualize_tree(model, X, y, ax=axi)
    axi.set_title('depth = {}'.format(depth))

fig.savefig('figures/05_08-decision-tree-levels.png')
```

/home/wangy/git/wangyingsm/Python-Data-Science-Handbook/notebooks/helpers_05_08.py:34: UserWarning: The following kwargs were not used by contour: 'clim'

/home/wangy/git/wangyingsm/Python-Data-Science-Handbook/notebooks/helpers_05_08.py:34: UserWarning: The following kwargs were not used by contour: 'clim'

/home/wangy/git/wangyingsm/Python-Data-Science-Handbook/notebooks/helpers_05_08.py:34: UserWarning: The following kwargs were not used by contour: 'clim'

/home/wangy/git/wangyingsm/Python-Data-Science-Handbook/notebooks/helpers_05_08.py:34: UserWarning: The following kwargs were not used by contour: 'clim'



Decision Tree Overfitting

决策树过拟合

```
In [35]: model = DecisionTreeClassifier()

fig, ax = plt.subplots(1, 2, figsize=(16, 6))
fig.subplots_adjust(left=0.0625, right=0.95, wspace=0.1)
visualize_tree(model, X[::2], y[::2], boundaries=False, ax=ax[0])
visualize_tree(model, X[1::2], y[1::2], boundaries=False, ax=ax[1])

fig.savefig('figures/05_08-decision-tree-overfitting.png')
```

/home/wangy/git/wangyingsm/Python-Data-Science-Handbook/notebooks/helpers_05_08.py:34: UserWarning: The following kwargs were not used by contour: 'clim'

/home/wangy/git/wangyingsm/Python-Data-Science-Handbook/notebooks/helpers_05_08.py:34: UserWarning: The following kwargs were not used by contour: 'clim'



Principal Component Analysis

主成分分析

Principal Components Rotation

主成分旋转

```
In [36]: from sklearn.decomposition import PCA

In [37]: def draw_vector(v0, v1, ax=None):
ax = ax or plt.gca()
arrowprops=dict(arrowstyle='>',
                linewidth=2,
                shrinkA=0, shrinkB=0)
ax.annotate('', v1, v0, arrowprops=arrowprops)

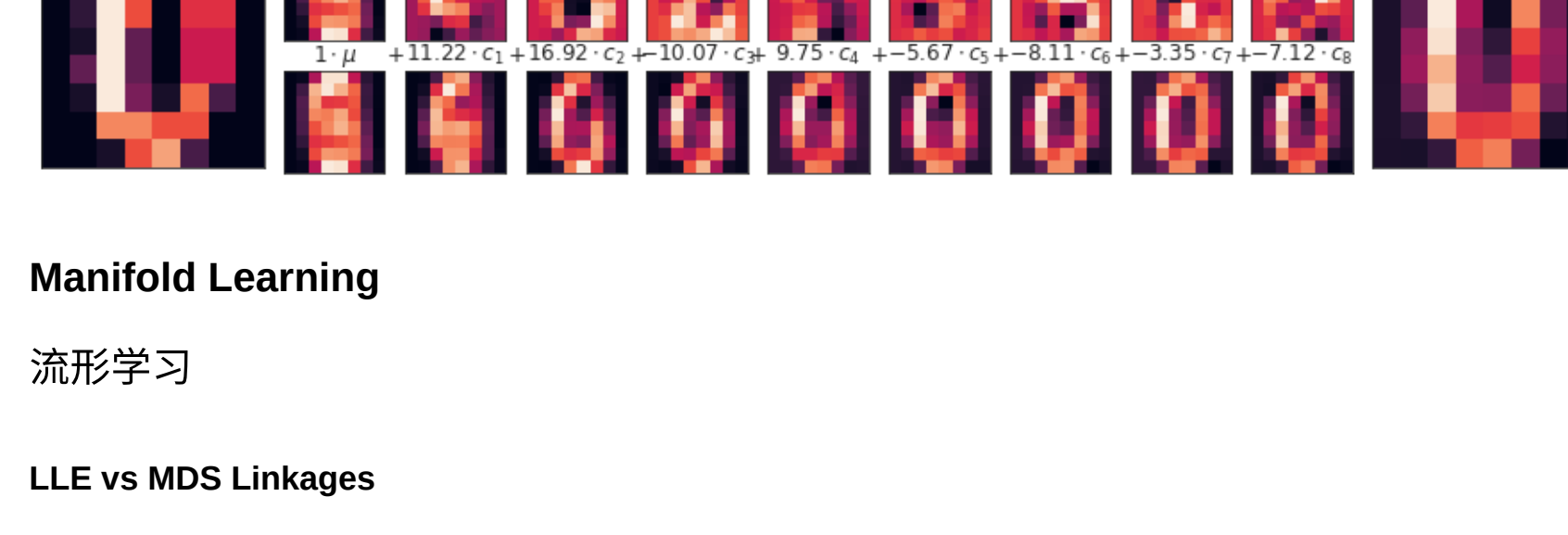
In [38]: rng = np.random.RandomState(1)
X = np.dot(rng.rand(2, 2), rng.rand(2, 200)).T
pca = PCA(n_components=2, whiten=True)
pca.fit(X)

fig, ax = plt.subplots(1, 2, figsize=(16, 6))
fig.subplots_adjust(left=0.0625, right=0.95, wspace=0.1)

# 绘制数据
ax[0].scatter(X[:, 0], X[:, 1])
for length, vector in zip(pca.explained_variance_, pca.components_):
    v = vector * 3 * np.sqrt(length)
    draw_vector(pca.mean_, pca.mean_ + v, ax=ax[0])
ax[0].axis('equal')
ax[0].set_xlabel='x', ylabel='y', title='input'

# 绘制主成分
X_pca = pca.transform(X)
ax[1].scatter(X_pca[:, 0], X_pca[:, 1], alpha=0.2)
draw_vector([0, 0], [0, 3], ax=ax[1])
draw_vector([0, 0], [3, 0], ax=ax[1])
ax[1].axis('equal')
ax[1].set_xlabel='component 1', ylabel='component 2',
        title='principal components',
        xlim=(-5, 5), ylim=(-3, 3.1))

fig.savefig('figures/05_09-PCA-rotation.png')
```



Digits Pixel Components

手写数字像素成分

```
In [39]: def plot_pca_components(x, coefficients=None, mean=0, components=None,
                             imshowshape=(8, 8), n_components=8, fontsize=12,
                             show_mean=True):
    if coefficients is None:
        coefficients = x
    if components is None:
        components = np.eye(len(coefficients), len(x))
    mean = np.zeros_like(x) + mean

    fig = plt.figure(figsize=(1.2 * 5 + n_components, 1.2 * 2))
    g = plt.GridSpec(2, 4 + bool(show_mean), left=0.05, hspace=0.3)

    def show(i, j, x, title=None):
        ax = fig.add_subplot(g[i, j], xticks=[], yticks=[])
        ax.imshow(x.reshape(imshowshape), interpolation='nearest')
        if title:
            ax.set_title(title, fontsize=fontsize)

    show(slice(2), slice(2), x, "True")
    approx = mean.copy()

    counter = 2
    if show_mean:
        show(0, 2, np.zeros_like(x) + mean, r'$\mu$')
        show(1, 2, approx, r'$\hat{\mu}$ dot $\mu$')
        counter += 1

    for i in range(n_components):
        approx = approx + coefficients[i] * components[i]
        show(0, 1 + counter, components[i], r'$c_{\{0\}}$')
        show(1, 1 + counter, approx, r'$c_{\{0\}}$')
        r"$S_{\{0\}}^2$ dot $c_{\{0\}}$")
        if show_mean or i > 0:
            plt.gca().text(0, 1.05, r'$\hat{\mu}$', ha='right', va='bottom',
                           transform=plt.gca().transAxes, fontsize=fontsize)

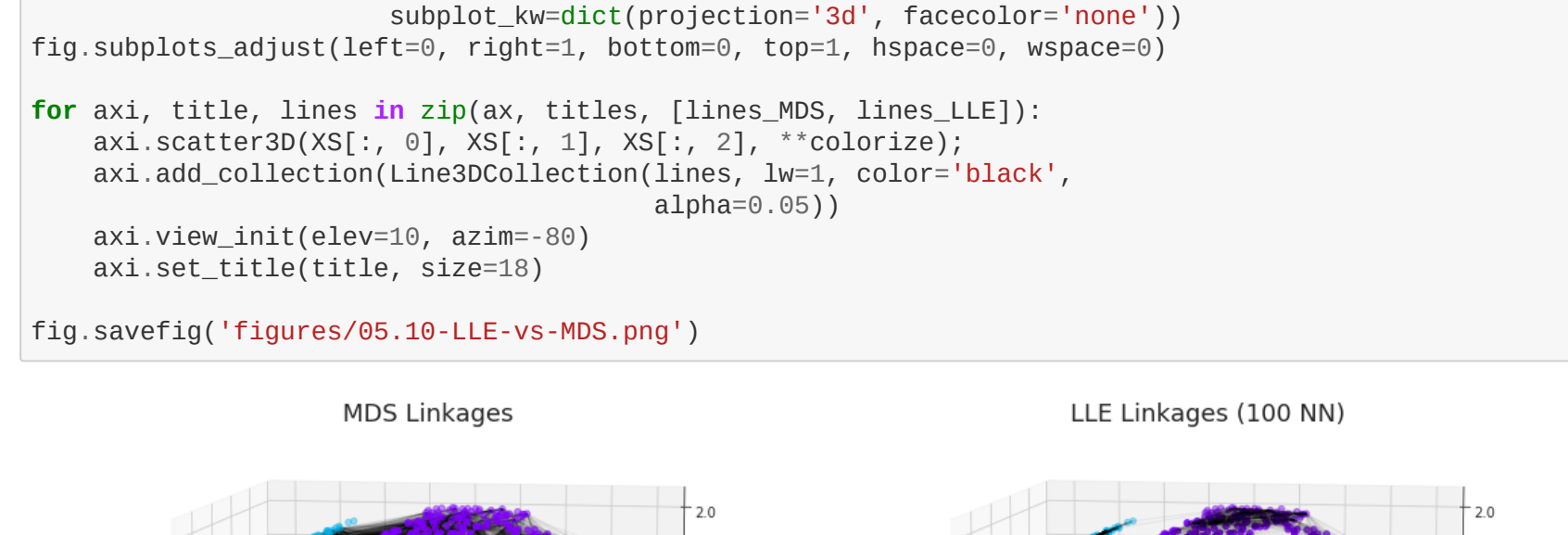
    show(slice(2), slice(-2, None), approx, "Approx")
    return fig

In [40]: from sklearn.datasets import load_digits

digits = load_digits()
sns.set_style('white')

fig = plot_pca_components(digits.data[10],
                          pca.components_)

fig.savefig('figures/05_09-digits-pixel-components.png')
```

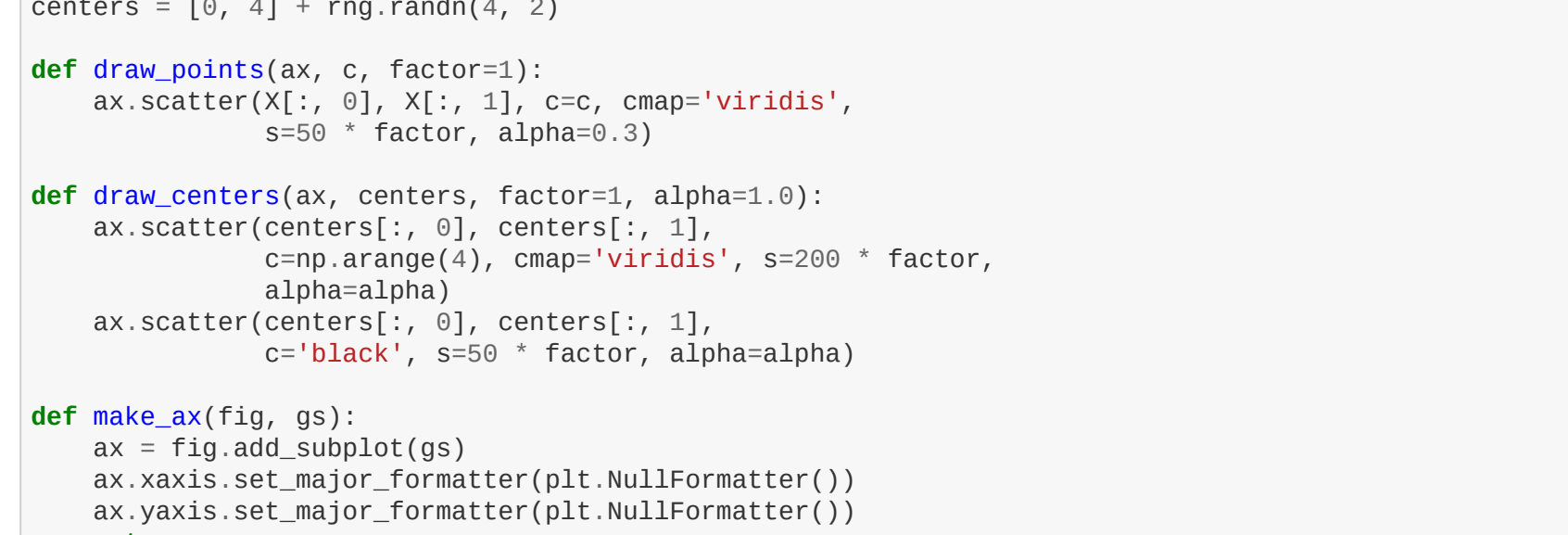


Digits PCA Components

手写数字PCA成分

```
In [41]: pca = PCA(n_components=8)
Xproj = pca.fit_transform(digits.data)
sns.set_style('white')
fig = plot_pca_components(digits.data[10], Xproj[10],
                          pca.components_)

fig.savefig('figures/05_09-digits-pca-components.png')
```



Manifold Learning

流形学习

LLE vs MDS Linkages

LLE对比MDS连接图

```
In [42]: def make_hello(N=1000, rseed=42):
# 创建'HELLO'图像，存储到hello.png文件中
fig, ax = plt.subplots(figsize=(4, 1))
fig.subplots_adjust(left=0, right=1, bottom=0, top=1)
ax.text('off')
ax.text(0.98, 0.4, 'HELLO', va='center', ha='center', weight='bold', size=85)
ax.axis([0, 1, 0, 1])
plt.close(fig)

# 打开图像文件，在上面随机绘制数据点
from matplotlib.image import imread
data = imread('hello.png')[::-1, :, 0].T
rng = np.random.RandomState(rseed)
X = rng.rand(N, 2)
i, j = (X * data.shape).astype(int).T
mask = (data[i, j] < 1)
X[:, 0] = (data.shape[0] / data.shape[1]) * X[:, 0]
X = X[N:]
return X[np.argsort(X[:, 0])]
```

```
In [43]: def make_hello_s_curve(X):
t = (X[:, 0] - 2) * 0.75 * np.pi
x = np.sin(t)
y = np.sin(t)
z = np.sign(t) * (np.cos(t) - 1)
return np.vstack((x, y, z)).T

X = make_hello(1000)
XS = make_hello_s_curve(X)
colorize = dict(c=X[:, 0], cmap=plt.cm.get_cmap('rainbow', 5))

In [44]: from mpl_toolkits.mplot3d.art3d import Line3DCollection
from sklearn.neighbors import NearestNeighbors

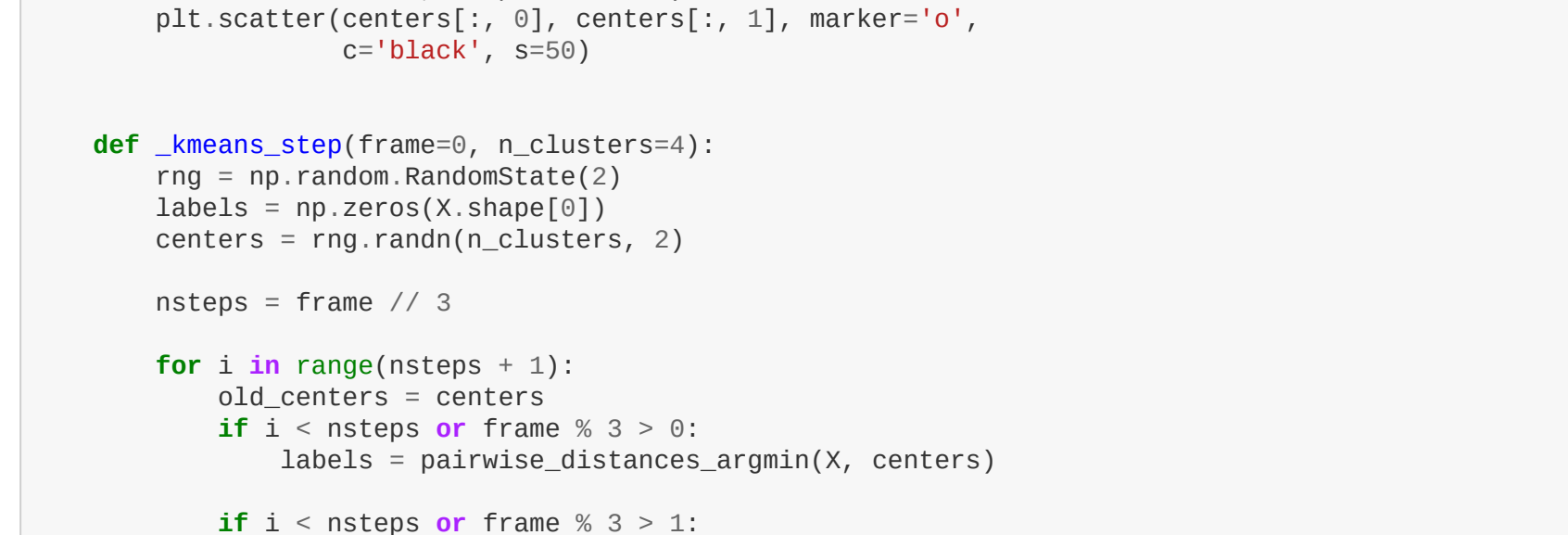
# 构建MDS的连接
rng = np.random.RandomState(42)
ind = rng.permutation(len(X))
lines_MDS = [XS[ind[i], XS[ind[j]] for i in ind[100:1000] for j in ind[100:200]]]

# 构建LLE的连接
nbrs = NearestNeighbors(n_neighbors=100).fit(XS)
lines_LLE = [XS[ind[i], XS[ind[j]] for i in range(100) for j in nbrs[i]]]
titles = ['MDS Linkages', 'LLE Linkages (100 NN)']

# 绘制结果
fig, ax = plt.subplots(1, 2, figsize=(16, 6))
subplot_kw=dict(projection='3d', facecolor='none')
fig.subplots_adjust(left=0, right=1, bottom=0, top=1, hspace=0)

for axi, title, lines in zip(ax, titles, [lines_MDS, lines_LLE]):
    axi.scatter3D(XS[:, 0], XS[:, 1], XS[:, 2], **colorize);
    axi.add_collection(Line3DCollection(lines, lw=1, color='black',
                                       alpha=0.05))
    axi.view_init(elev=10, azim=-80)
    axi.set_title(title, size=18)

fig.savefig('figures/05_10-LLE-vs-MDS.png')
```



K-Means

K均值

Expectation-Maximization

期望最大化

图表所在正文

The following figure shows a visual depiction of the Expectation-Maximization approach to K Means:

下面图表展示了K均值的期望最大化算法的可可视化说明:

```
In [45]: from sklearn.datasets.samples_generator import make_blobs
from sklearn.metrics import pairwise_distances_argmin

X, y_true = make_blobs(n_samples=300, centers=4,
                      cluster_std=0.60, random_state=0)

rng = np.random.RandomState(42)
centers = [0, 4] + rng.random(4, 2)

def draw_points(ax, c, factor=1):
    ax.scatter(X[:, 0], X[:, 1], c=c, cmap='viridis',
               s=50 * factor, alpha=0.5)

def draw_centers(ax, centers, factor=1, alpha=1.0):
    ax.scatter(centers[:, 0], centers[:, 1],
               c=np.arange(4), cmap='viridis', s=200 * factor,
               alpha=alpha)
    ax.scatter(centers[:, 0], centers[:, 1],
               c='black', s=50 * factor, alpha=alpha)

def make_ax(fig, gs):
    ax = fig.add_subplot(gs)
    ax.xaxis.set_major_formatter(plt.NullFormatter())
    ax.yaxis.set_major_formatter(plt.NullFormatter())
    return ax

fig = plt.figure(figsize=(15, 4))
gs = plt.GridSpec(4, 15, left=0.02, right=0.98, bottom=0.05, top=0.95, wspace=0.2, hspace=0.2)
ax0 = make_ax(fig, gs[0, :4])
ax2 = make_ax(fig, gs[2, 5: 2 * 5 + 2 * 1])

# E-step 期望步骤
y_pred = pairwise_distances_argmin(X, centers)
draw_points(ax1, y_pred)
draw_centers(ax1, centers)

# M-step 最大化步骤
new_centers = np.array([X[y_pred == i].mean(0) for i in range(4)])
draw_points(ax2, y_pred)
draw_centers(ax2, centers, alpha=0.3)
for i in range(4):
    ax2.annotate('', new_centers[i], centers[i],
                 arrowprops=dict(arrowstyle='>', linewidth=1))

# 完成迭代
centers = new_centers
ax1.text(0.95, 0.95, "E-Step", transform=ax1.transAxes, ha='right', va='top', size=14)
ax2.text(0.95, 0.95, "M-Step", transform=ax2.transAxes, ha='right', va='top', size=14)

# Final E-step 最终期望步骤
y_pred = pairwise_distances_argmin(X, centers)
axf = make_ax(fig, gs[-1, :4])
draw_points(axf, y_pred, factor=2)
draw_centers(axf, centers, factor=2)
axf.text(0.98, 0.98, "Final Clustering", transform=axf.transAxes,
        ha='right', va='top', size=16)

fig.savefig('figures/05_11-expectation-maximization.png')
```



Interactive K-Means

交互式K均值

The following script uses IPython's interactive widgets to demonstrate the K-means algorithm interactively. Run this within the IPython notebook to explore the expectation maximization algorithm for computing K Means.

下面脚本使用IPython的交互式组件展示K均值算法。在IPython notebook中运行它们来研究计算K均值的期望最大化算法。

```
In [46]: from matplotlib.inline
import matplotlib.pyplot as plt
import seaborn; seaborn.set() # seaborn风格
import numpy as np

from ipywidgets import interact
from sklearn.metrics import pairwise_distances_argmin
from sklearn.datasets.samples_generator import make_blobs

def plot_kmeans_interactive(min_clusters=3, max_clusters=6):
    X, y = make_blobs(n_samples=300, centers=4,
                    random_state=0, cluster_std=0.60)

    def plot_points(X, labels, n_clusters):
        plt.scatter(X[:, 0], X[:, 1], c=labels, s=50, cmap='viridis',
                    vmin=0, vmax=n_clusters - 1)

    def plot_centers(centers):
        plt.scatter(centers[:, 0], centers[:, 1], marker='o',
                    c=np.arange(centers.shape[0]),
                    cmap='viridis')
        plt.scatter(centers[:, 0], centers[:, 1], marker='o',
                    c='black', s=50)

    def kmeans_step(frame=0, n_clusters=4):
        rng = np.random.RandomState(2)
        labels = np.zeros(X.shape[0])
        centers = rng.random(n_clusters, 2)

        nsteps = frame // 3

        for i in range(nsteps + 1):
            old_centers = centers
            if i < nsteps or frame % 3 > 0:
                labels = pairwise_distances_argmin(X, centers)

            if i < nsteps or frame % 3 > 1:
                centers = np.array([X[labels == j].mean(0)
                                   for j in range(n_clusters)])
                nans = np.isnan(centers)
                centers[nans] = old_centers[nans]

            # 绘制数据和聚类中心点
            plot_points(X, labels, n_clusters)
            plot_centers(centers)

            # 在第三步时更新聚类中心点
            if frame % 3 == 2:
                for i in range(n_clusters):
                    plt.annotate('', centers[i], old_centers[i],
                                arrowprops=dict(arrowstyle='>', linewidth=1))
            plot_centers(centers)

        plt.xlim(-4, 4)
        plt.ylim(-2, 10)

        if frame % 3 == 1:
            plt.text(3, 9.5, "1. Reassign points to nearest centroid",
                    ha='right', va='top', size=14)
        elif frame % 3 == 2:
            plt.text(3, 9.5, "2. Update centroids to cluster means",
                    ha='right', va='top', size=14)

        return interact(kmeans_step, frame=[0, 60],
                       n_clusters=[min_clusters, max_clusters])

plot_kmeans_interactive();
```



Gaussian Mixture Models

高斯混合模型

Covariance Type

图表所在正文

```
In [47]: from sklearn.mixture import GaussianMixture as GMM
from matplotlib.patches import Ellipse

def draw_an_ellipse(position, covariance, ax=None, **kwargs):
    """Draw an ellipse with a given position and covariance"""
    ax = ax or plt.gca()

    # 将协方差转换到主坐标轴
    if covariance.shape == (2, 2):
        U, S, Vt = np.linalg.svd(covariance)
        angle = np.degrees(np.arctan2(U[1, 0], U[0, 0]))
        width, height = 2 * np.sqrt(S)
    else:
        angle = 0
        width, height = 2 * np.sqrt(covariance)

    # 绘制椭圆
    for nsig in range(1, 4):
        ax.add_patch(Ellipse(position, nsig * width, nsig * height,
                             angle, **kwargs))

fig, ax = plt.subplots(1, 3, figsize=(14, 4))
fig.subplots_adjust(wspace=0.05)

rng = np.random.RandomState(5)
X = np.dot(rng.rand(500, 2), rng.rand(2, 2))

for i, cov_type in enumerate(['diag', 'spherical', 'full']):
    model = GMM(2, covariance_type=cov_type).fit(X)
    ax[i].axis('equal')
    ax[i].scatter(X[:, 0], X[:, 1], alpha=0.5)
    ax[i].set_xlim(-3, 3)
    ax[i].set_title(covariance_type+"(0)", format(cov_type),
                    size=14, family='monospace')
    if model.covariance_type == 'spherical':
        cov = np.eye(2) * model.covariances_
    else:
        cov = model.covariances_[0]
    draw_ellipse(model.means_[0], cov, ax[i], alpha=0.2)
    ax[i].axis.set_major_formatter(plt.NullFormatter())
    ax[i].yaxis.set_major_formatter(plt.NullFormatter())

fig.savefig('figures/05_12-covariance-type.png')
```

