

Chapter 6 – Decision Trees

This notebook contains all the sample code and solutions to the exercises in chapter 6.

 Open in Colab

 Open in Kaggle

Setup

This project requires Python 3.7 or above:

```
In [1]: import sys
        assert sys.version_info >= (3, 7)
```

It also requires Scikit-Learn \geq 1.0.1:

```
In [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:

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

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

And let's create the `images/decision_trees` 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:

```
In [4]: from pathlib import Path

        IMAGES_PATH = Path() / "images" / "decision_trees"
        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)
```

Regression

Let's prepare a simple quadratic training set:

Code example:

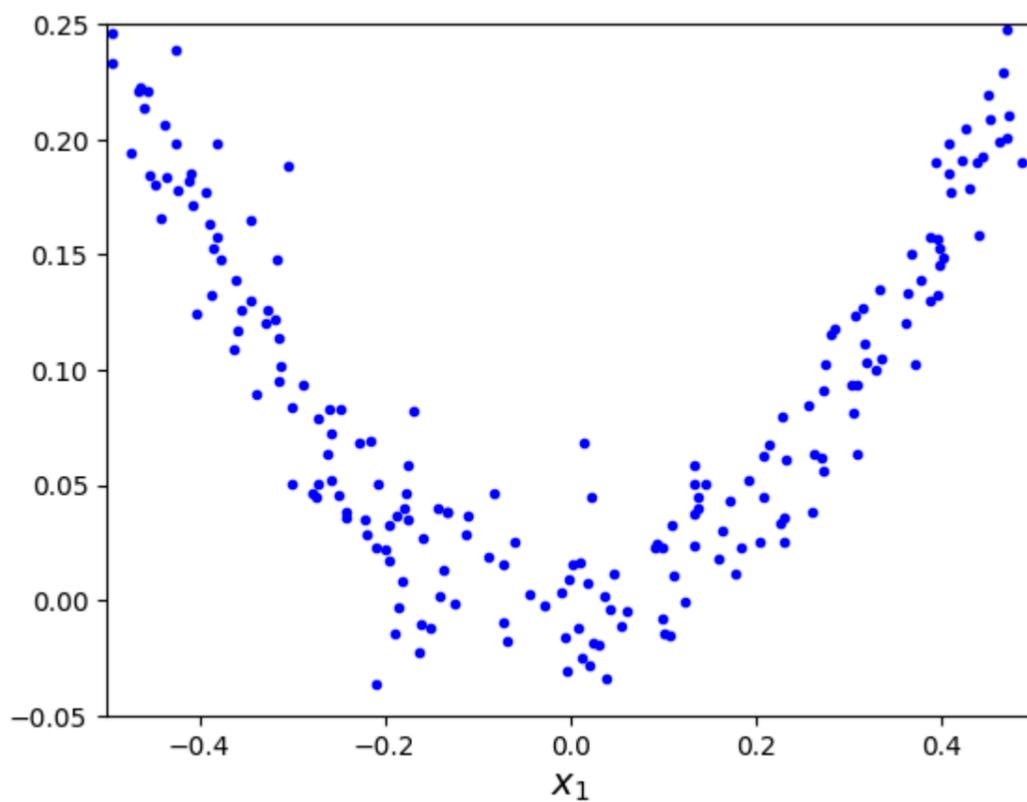
```
In [5]: from sklearn.tree import DecisionTreeRegressor
        import numpy as np

        np.random.seed(42)
        X_quad = np.random.rand(200, 1) - 0.5 # a single random input feature
        y_quad = X_quad ** 2 + 0.025 * np.random.randn(200, 1)

        axes = [-0.5, 0.5, -0.05, 0.25]
        x1 = np.linspace(axes[0], axes[1], 500).reshape(-1, 1)
        plt.axis(axes)
        plt.xlabel("$x_1$")
        plt.plot(X_quad, y_quad, "b.")

        tree_reg = DecisionTreeRegressor(max_depth=2, random_state=42)
        tree_reg.fit(X_quad, y_quad)
```

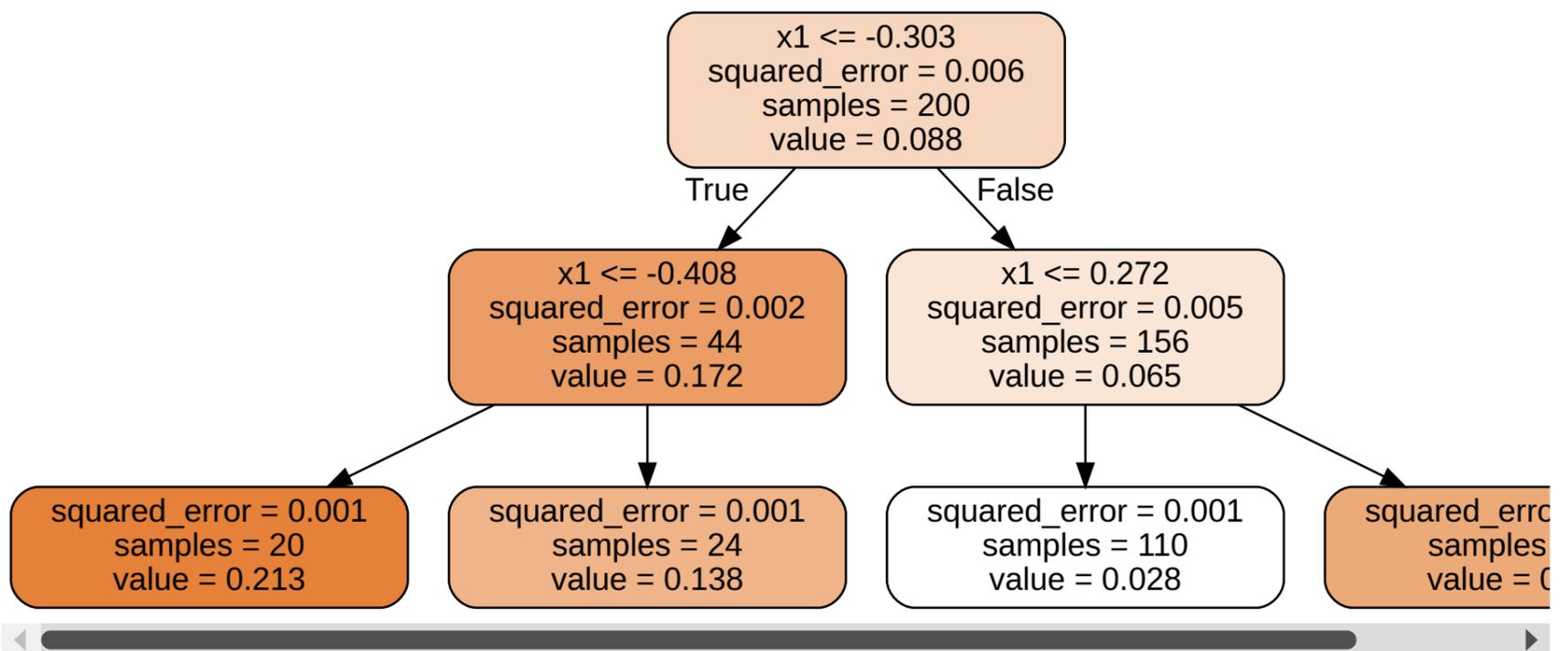
```
Out[5]: DecisionTreeRegressor ⓘ ?
        Parameters
```



```
In [6]: from sklearn.tree import export_graphviz
from graphviz import Source

# extra code - we've already seen how to use export_graphviz()
export_graphviz(
    tree_reg,
    out_file=str(IMAGE_PATH / "regression_tree.dot"),
    feature_names=["x1"],
    rounded=True,
    filled=True
)
Source.from_file(IMAGE_PATH / "regression_tree.dot")
```

Out[6]:



```
In [7]: tree_reg2 = DecisionTreeRegressor(max_depth=3, random_state=42)
tree_reg2.fit(X_quad, y_quad)
```

Out[7]: **DecisionTreeRegressor** ⓘ ?

- Parameters

```
In [8]: tree_reg.tree_.threshold
```

Out[8]: array([-0.30265072, -0.40830374, -2. , -2. , 0.27175756, -2. , -2.])

```
In [9]: tree_reg2.tree_.threshold
```

Out[9]: array([-0.30265072, -0.40830374, -0.45416115, -2. , -2. , -0.37022041, -2. , -2. , 0.27175756, -0.21270403, -2. , -2. , 0.40399227, -2. , -2.])

```
In [10]: # extra code - this cell generates and saves Figure 6-5

def plot_regression_predictions(tree_reg, X, y, axes=[-0.5, 0.5, -0.05, 0.25]):
    x1 = np.linspace(axes[0], axes[1], 500).reshape(-1, 1)
    y_pred = tree_reg.predict(x1)
```

```

plt.axis(axes)
plt.xlabel("$x_1$")
plt.plot(X, y, "b.")
plt.plot(x1, y_pred, "r.-", linewidth=2, label=r"$\hat{y}$")

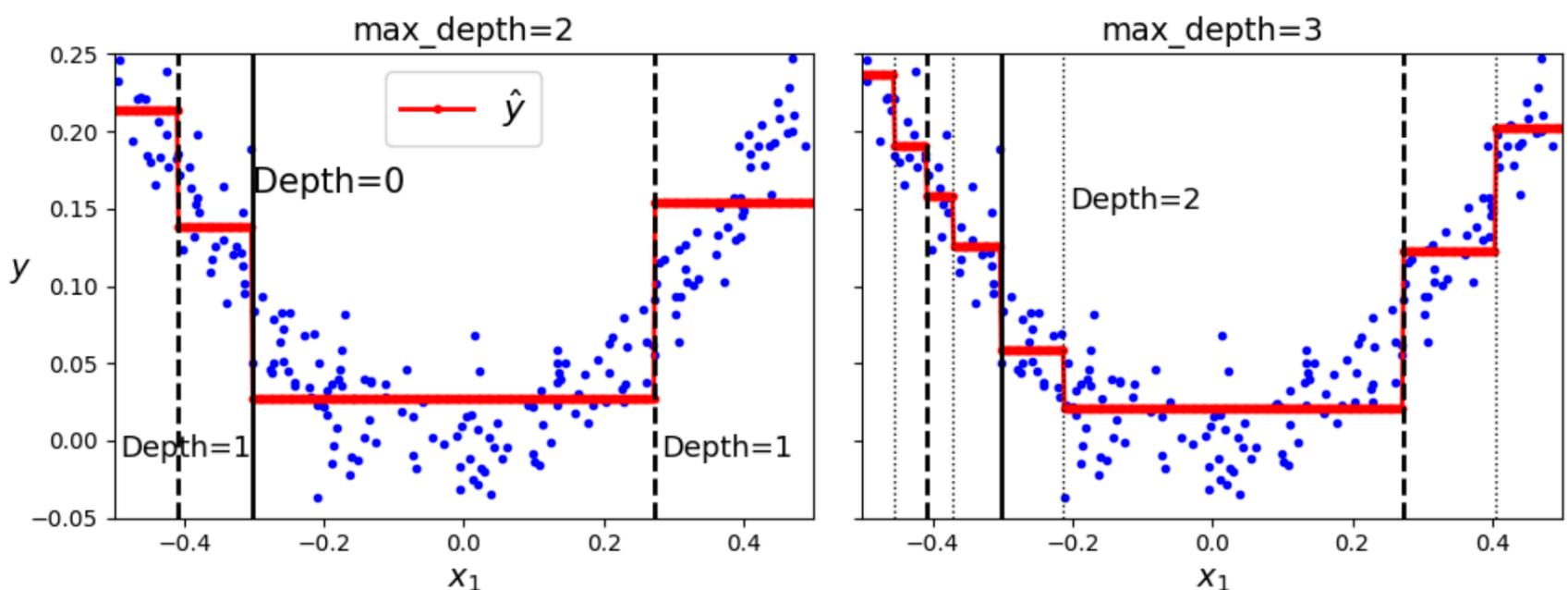
fig, axes = plt.subplots(ncols=2, figsize=(10, 4), sharey=True)
plt.sca(axes[0])
plot_regression_predictions(tree_reg, X_quad, y_quad)

th0, th1a, th1b = tree_reg.tree_.threshold[[0, 1, 4]]
for split, style in ((th0, "k-"), (th1a, "k--"), (th1b, "k--")):
    plt.plot([split, split], [-0.05, 0.25], style, linewidth=2)
plt.text(th0, 0.16, "Depth=0", fontsize=15)
plt.text(th1a + 0.01, -0.01, "Depth=1", horizontalalignment="center", fontsize=13)
plt.text(th1b + 0.01, -0.01, "Depth=1", fontsize=13)
plt.ylabel("$y$", rotation=0)
plt.legend(loc="upper center", fontsize=16)
plt.title("max_depth=2")

plt.sca(axes[1])
th2s = tree_reg2.tree_.threshold[[2, 5, 9, 12]]
plot_regression_predictions(tree_reg2, X_quad, y_quad)
for split, style in ((th0, "k-"), (th1a, "k--"), (th1b, "k--")):
    plt.plot([split, split], [-0.05, 0.25], style, linewidth=2)
for split in th2s:
    plt.plot([split, split], [-0.05, 0.25], "k:", linewidth=1)
plt.text(th2s[2] + 0.01, 0.15, "Depth=2", fontsize=13)
plt.title("max_depth=3")

save_fig("tree_regression_plot")
plt.show()

```



In [11]: # extra code – this cell generates and saves Figure 6–6

```

tree_reg1 = DecisionTreeRegressor(random_state=42)
tree_reg2 = DecisionTreeRegressor(random_state=42, min_samples_leaf=10)
tree_reg1.fit(X_quad, y_quad)
tree_reg2.fit(X_quad, y_quad)

x1 = np.linspace(-0.5, 0.5, 500).reshape(-1, 1)
y_pred1 = tree_reg1.predict(x1)
y_pred2 = tree_reg2.predict(x1)

fig, axes = plt.subplots(ncols=2, figsize=(10, 4), sharey=True)

plt.sca(axes[0])
plt.plot(X_quad, y_quad, "b.")
plt.plot(x1, y_pred1, "r.-", linewidth=2, label=r"$\hat{y}$")
plt.axis([-0.5, 0.5, -0.05, 0.25])
plt.xlabel("$x_1$")
plt.ylabel("$y$", rotation=0)
plt.legend(loc="upper center")
plt.title("No restrictions")

plt.sca(axes[1])
plt.plot(X_quad, y_quad, "b.")
plt.plot(x1, y_pred2, "r.-", linewidth=2, label=r"$\hat{y}$")
plt.axis([-0.5, 0.5, -0.05, 0.25])
plt.xlabel("$x_1$")
plt.title(f"min_samples_leaf={tree_reg2.min_samples_leaf}")

save_fig("tree_regression_regularization_plot")
plt.show()

```

