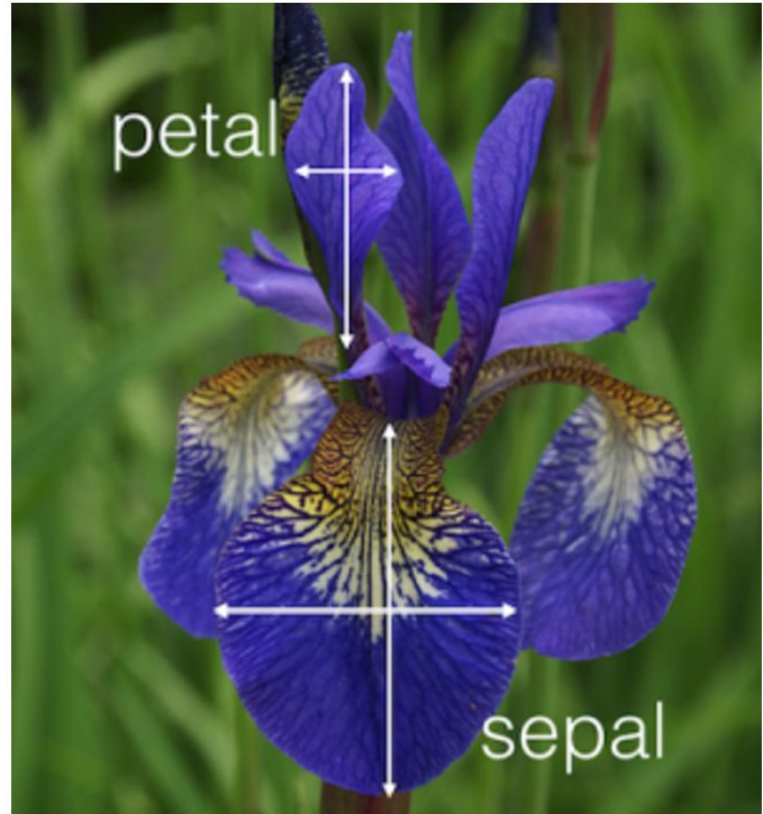


Hands on with scikit learn

Agenda

- What is the famous iris dataset, and how does it relate to machine learning?
- How do we load the iris dataset into scikit-learn?
- How do we describe a dataset using machine learning terminology?
- What are scikit-learn's four key requirements for working with data?

Introducing the iris dataset



- 50 samples of 3 different species of iris (150 samples total)
- Measurements: sepal length, sepal width, petal length, petal width

```
# import load_iris function from datasets module
```

```
from sklearn.datasets import load_
```

```
# save "bunch" object containing iris dataset and its attributes
```

```
iris = load_iris()
```

```
type(iris)
```

```
# print the iris data
```

```
print(iris.data)
```

```
# print the names of the four features
```

```
print(iris.feature_names)
```

```
# print integers representing the species of each observation
```

```
print(iris.target)
```

```
# print the encoding scheme for species: 0 = setosa, 1 = versicolor, 2 = virginica
```

```
print(iris.target_names)
```

Requirements for working with data in scikit-learn

1. Features and response are **separate objects**
2. Features and response should be **numeric**
3. Features and response should be **NumPy arrays**
4. Features and response should have **specific shapes**

check the types of the features and response

```
print(type(iris.data))
```

```
print(type(iris.target))
```

check the shape of the features (first dimension = number of observations, second dimensions = number of features)

```
print(iris.data.shape)
```

check the shape of the response (single dimension matching the number of observations)

```
print(iris.target.shape)
```

store feature matrix in "X"

```
X = iris.data
```

store response vector in "y"

```
y = iris.target
```

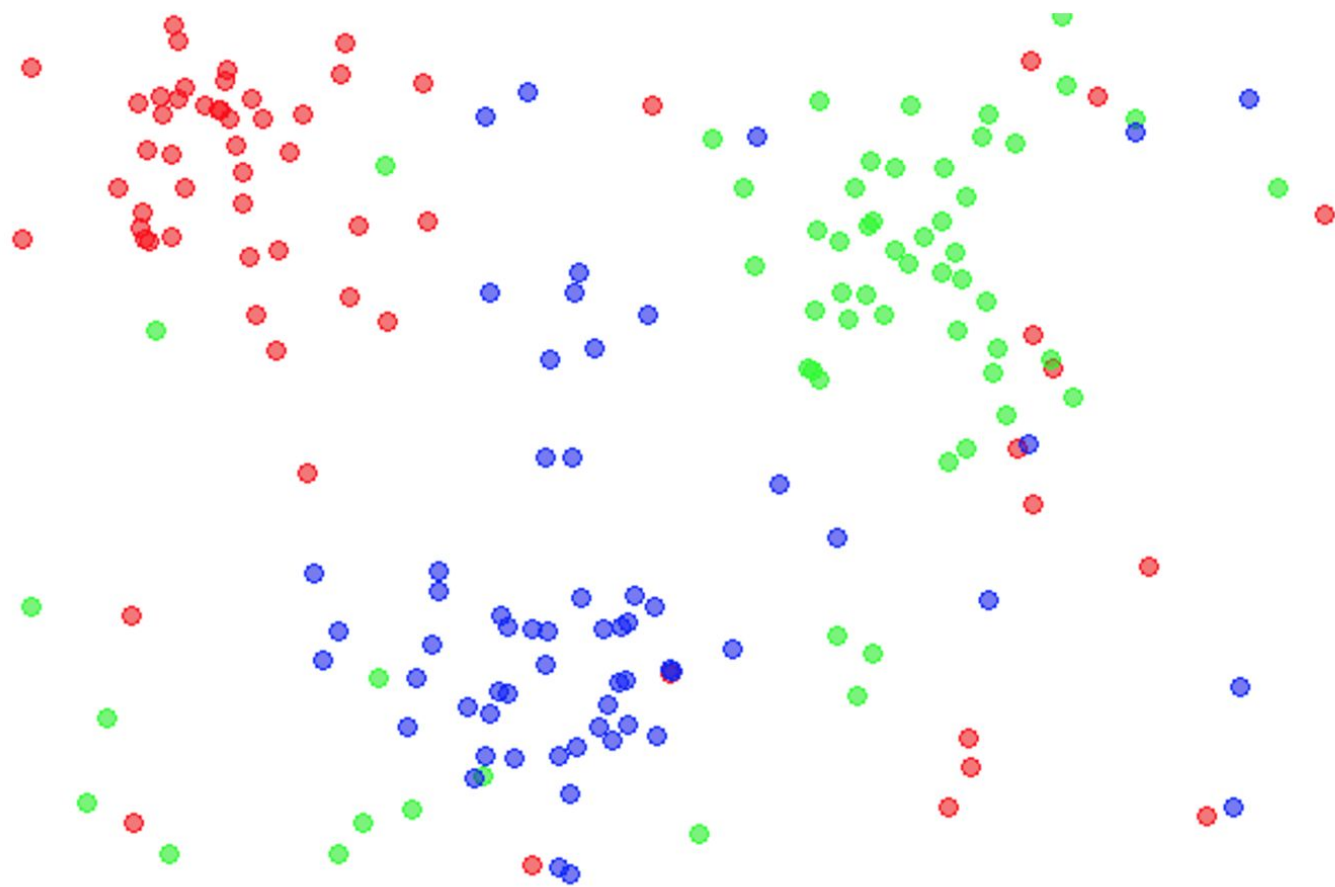
Agenda

- What is the **K-nearest neighbors** classification model?
- What are the four steps for **model training and prediction** in scikit-learn?
- How can I apply this pattern to **other machine learning models**?

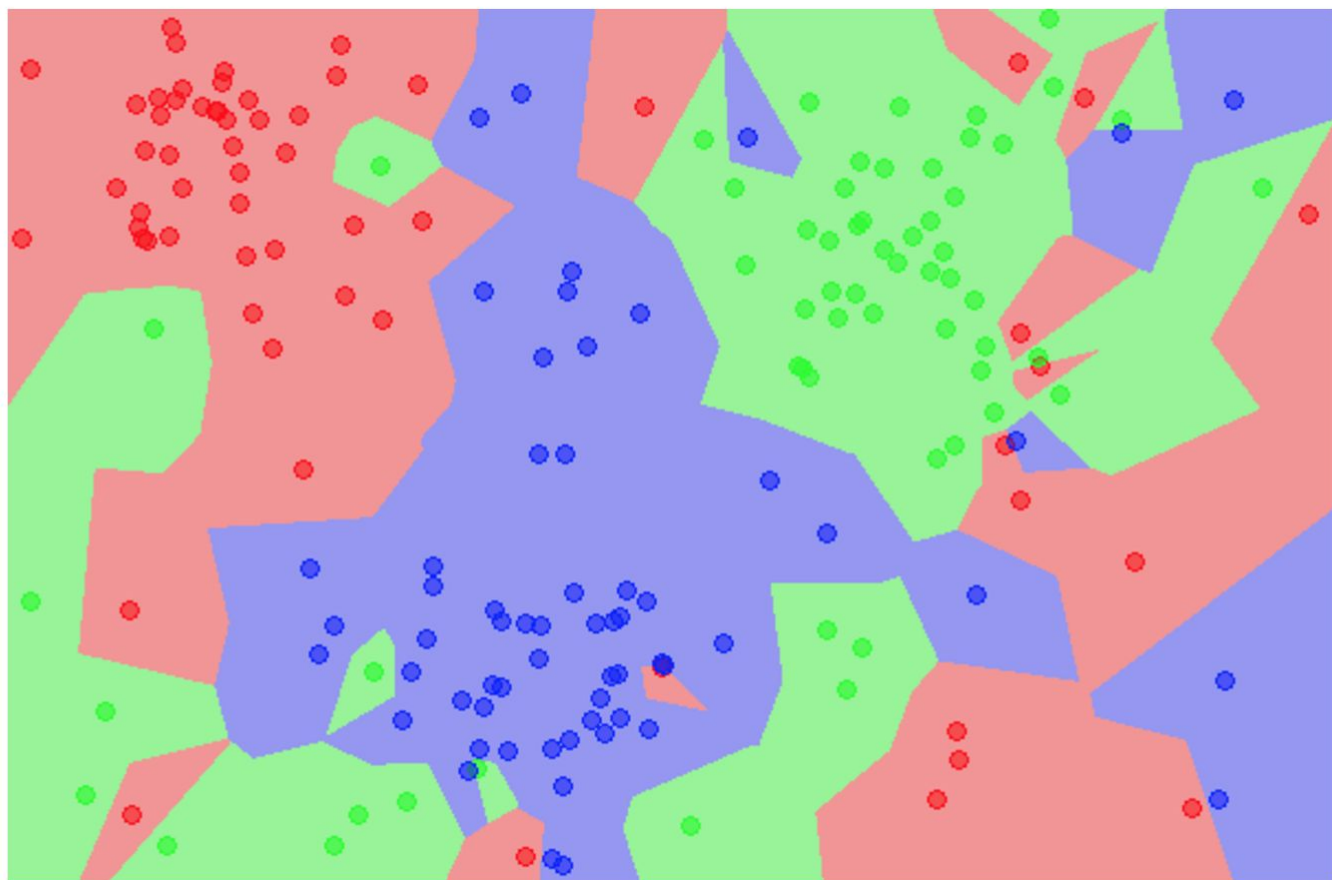
K-nearest neighbors (KNN) classification

1. Pick a value for K.
2. Search for the K observations in the training data that are "nearest" to the measurements of the unknown iris.
3. Use the most popular response value from the K nearest neighbors as the predicted response value for the unknown iris.

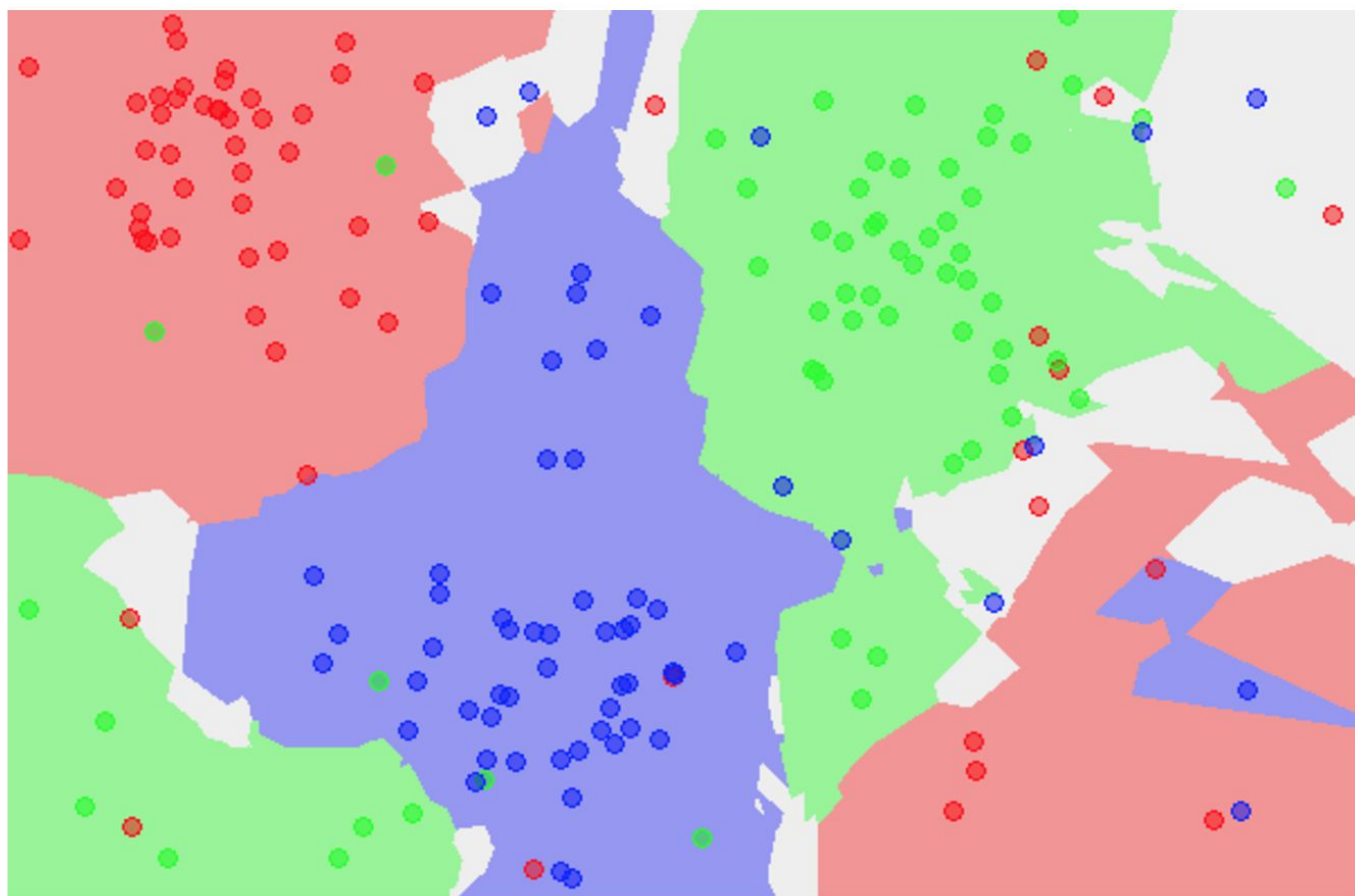
Example training data



KNN classification map (K=1)



KNN classification map (K=5)



```
# import load_iris function from datasets module
```

```
from sklearn.datasets import load_iris
```

```
# save "bunch" object containing iris dataset and its attributes
```

```
iris = load_iris()
```

```
# store feature matrix in "X"
```

```
X = iris.data
```

```
# store response vector in "y"
```

```
y = iris.target
```

```
# print the shapes of X and y
```

```
print(X.shape)
```

```
print(y.shape)
```

```
from sklearn.neighbors import KNeighborsClassifier
```

```
knn = KNeighborsClassifier(n_neighbors=1)
```

```
print(knn)
```

```
knn.fit(X, y)
```

```
knn.predict([[3, 5, 4, 2]])
```

```
X_new = [[3, 5, 4, 2], [5, 4, 3, 2]]
```

```
knn.predict(X_new)
```

Using a different value for K

instantiate the model (using the value K=5)

```
knn = KNeighborsClassifier(n_neighbors=5)
```

fit the model with data

```
knn.fit(X, y)
```

predict the response for new observations

```
knn.predict(X_new)
```

Using a different classification model

import the class

```
from sklearn.linear_model import LogisticRegression
```

instantiate the model (using the default parameters)

```
logreg = LogisticRegression()
```

fit the model with data

```
logreg.fit(X, y)
```

predict the response for new observations

```
logreg.predict(X_new)
```

Review

- Classification task: Predicting the species of an unknown iris
- Used three classification models: KNN ($K=1$), KNN ($K=5$), logistic regression
- Need a way to choose between the models

Agenda

- How do I choose **which model to use** for my supervised learning task?
- How do I choose the **best tuning parameters** for that model?
- How do I estimate the **likely performance of my model** on out-of-sample data?

Evaluation procedure #1: Train and test on the entire dataset

```
# read in the iris data
```

```
from sklearn.datasets import load_iris
```

```
iris = load_iris()
```

```
# create X (features) and y (response)
```

```
X = iris.data
```

```
y = iris.target
```

Evaluation procedure #1: Train and test on the entire dataset

```
# read in the iris data
```

```
from sklearn.datasets import load_iris
```

```
iris = load_iris()
```

```
# create X (features) and y (response)
```

```
X = iris.data
```

```
y = iris.target
```

Logistic regression

import the class

```
from sklearn.linear_model import LogisticRegression
```

instantiate the model (using the default parameters)

```
logreg = LogisticRegression()
```

fit the model with data

```
logreg.fit(X, y)
```

predict the response values for the observations in X

```
logreg.predict(X)
```

store the predicted response values

```
y_pred = logreg.predict(X)
```

check how many predictions were generated

```
len(y_pred)
```

compute classification accuracy for the logistic regression model

```
from sklearn import metrics
```

```
print(metrics.accuracy_score(y, y_pred))
```

KNN (K=5)

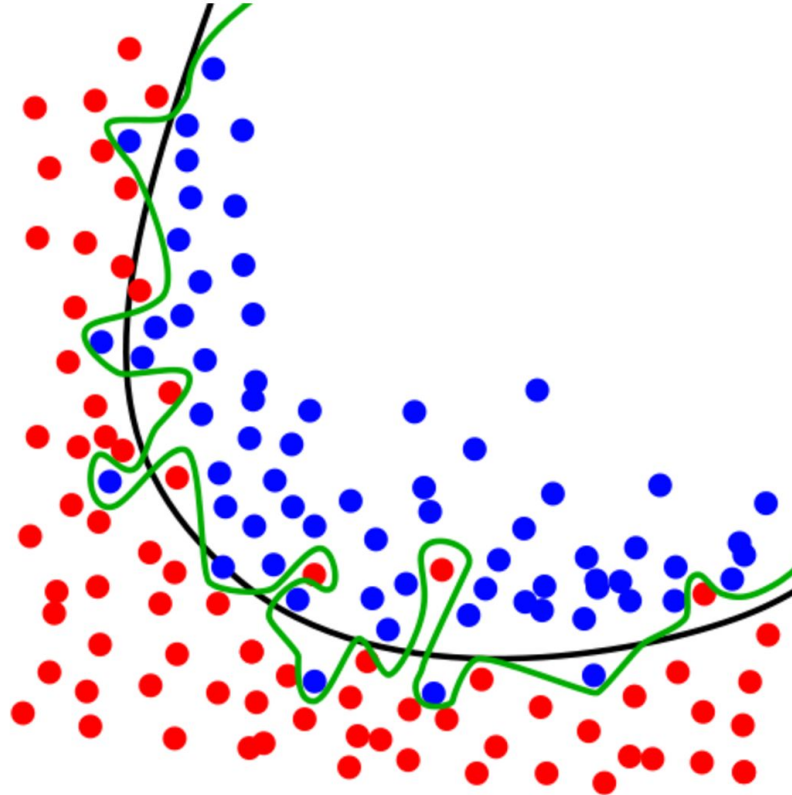
```
from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X, y)
y_pred = knn.predict(X)
print(metrics.accuracy_score(y, y_pred))
```

KNN (K=1)

```
knn = KNeighborsClassifier(n_neighbors=1)
knn.fit(X, y)
y_pred = knn.predict(X)
print(metrics.accuracy_score(y, y_pred))
```

Problems with training and testing on the same data

- Goal is to estimate likely performance of a model on **out-of-sample data**
- But, maximizing training accuracy rewards **overly complex models** that won't necessarily generalize
- Unnecessarily complex models **overfit** the training data



Evaluation procedure #2: Train/test split

```
# print the shapes of X and y
```

```
print(X.shape)
```

```
print(y.shape)
```

```
# STEP 1: split X and y into training and testing sets
```

```
from sklearn.cross_validation import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4, random_state=4)
```


X_train

X_test

feature 1	feature 2	response
1	2	2
3	4	12
5	6	30
7	8	56
9	10	90

y_train

y_test

print the shapes of the new X objects

```
print(X_train.shape)
```

```
print(X_test.shape)
```

print the shapes of the new y objects

```
print(y_train.shape)
```

```
print(y_test.shape)
```

STEP 2: train the model on the training set

```
logreg = LogisticRegression()
```

```
logreg.fit(X_train, y_train)
```

STEP 3: make predictions on the testing set

```
y_pred = logreg.predict(X_test)
```

compare actual response values (y_test) with predicted response values (y_pred)

```
print(metrics.accuracy_score(y_test, y_pred))
```

Repeat for KNN

```
knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)
print(metrics.accuracy_score(y_test, y_pred))
```

```
knn = KNeighborsClassifier(n_neighbors=1)
knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)
print(metrics.accuracy_score(y_test, y_pred))
```

Can we locate even better K value

```
# try K=1 through K=25 and record testing accuracy
```

```
k_range = list(range(1, 26))
```

```
scores = []
```

```
for k in k_range:
```

```
    knn = KNeighborsClassifier(n_neighbors=k)
```

```
    knn.fit(X_train, y_train)
```

```
    y_pred = knn.predict(X_test)
```

```
    scores.append(metrics.accuracy_score(y_test, y_pred))
```

```
# import Matplotlib (scientific plotting library)
```

```
import matplotlib.pyplot as plt
```

```
# plot the relationship between K and testing accuracy
```

```
plt.plot(k_range, scores)
```

```
plt.xlabel('Value of K for KNN')
```

```
plt.ylabel('Testing Accuracy')
```

Making predictions on out-of-sample data

instantiate the model with the best known parameters

```
knn = KNeighborsClassifier(n_neighbors=11)
```

train the model with X and y (not X_train and y_train)

```
knn.fit(X, y)
```

make a prediction for an out-of-sample observation

```
knn.predict([[3, 5, 4, 2]])
```

Downsides of train/test split?

- Provides a **high-variance estimate** of out-of-sample accuracy
- **K-fold cross-validation** overcomes this limitation
- But, train/test split is still useful because of its **flexibility and speed**

Agenda

- How do I use the **pandas library** to read data into Python?
- How do I use the **seaborn library** to visualize data?
- What is **linear regression**, and how does it work?
- How do I **train and interpret** a linear regression model in scikit-learn?
- What are some **evaluation metrics** for regression problems?
- How do I choose **which features to include** in my model?

Reading data using pandas

```
# conventional way to import pandas
```

```
import pandas as pd
```

```
# read CSV file directly from a URL and save the results
```

```
data = pd.read_csv('http://www-bcf.usc.edu/~gareth/ISL/Advertising.csv', index_col=0)
```

```
# display the first 5 rows
```

```
data.head()
```

```
# display the last 5 rows
```

```
data.tail()
```

```
# check the shape of the DataFrame (rows, columns)
```

```
data.shape
```

What are the features?

- **TV:** advertising dollars spent on TV for a single product in a given market (in thousands of dollars)
- **Radio:** advertising dollars spent on Radio
- **Newspaper:** advertising dollars spent on Newspaper

What is the response?

- **Sales:** sales of a single product in a given market (in thousands of items)

What else do we know?

- Because the response variable is continuous, this is a **regression** problem.
- There are 200 **observations** (represented by the rows), and each observation is a single market.

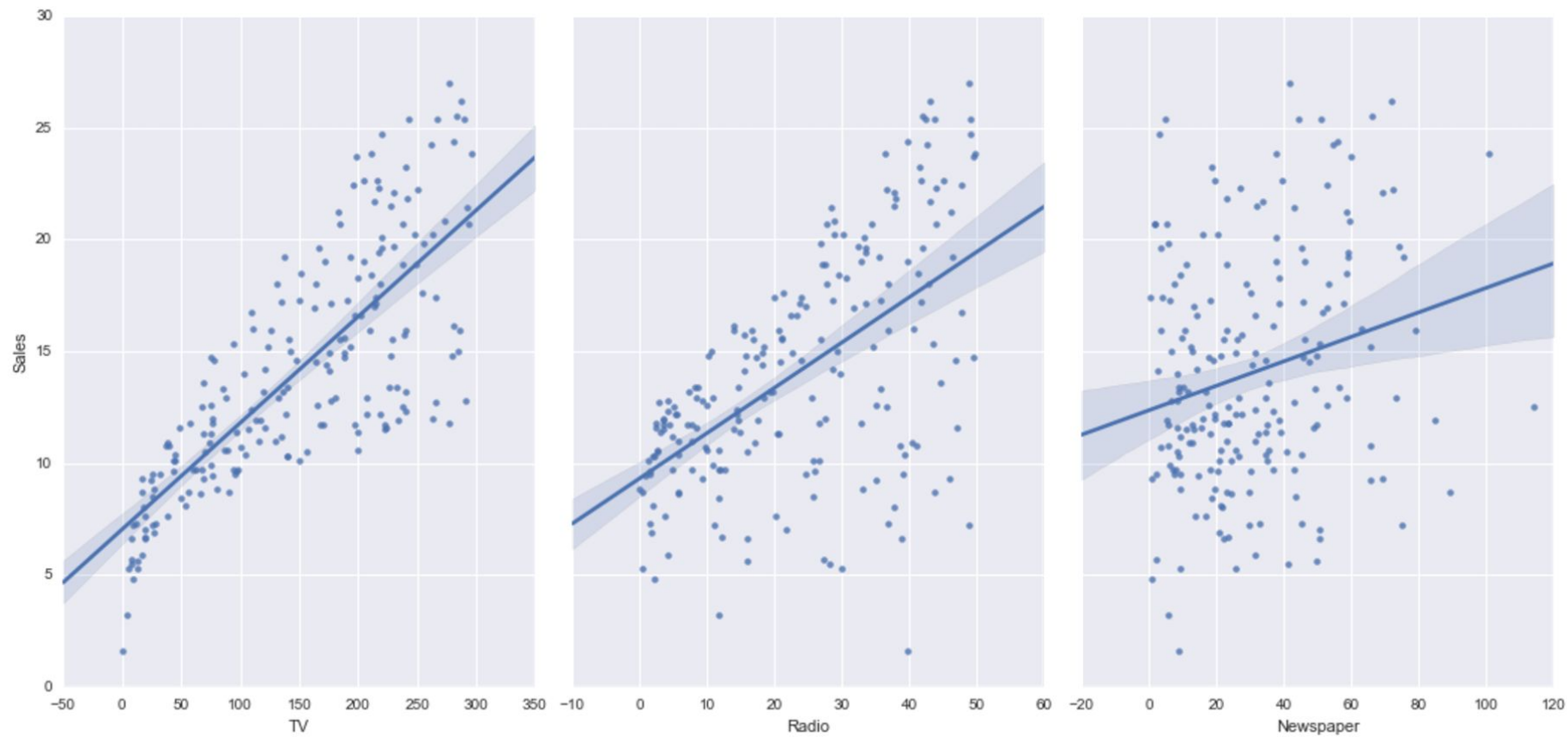
Visualizing data using seaborn

conventional way to import seaborn

```
import seaborn as sns
```

visualize the relationship between the features and the response using scatterplots

```
sns.pairplot(data, x_vars=['TV', 'Radio', 'Newspaper'], y_vars='Sales', size=7, aspect=0.7, kind='reg')
```



Linear regression

Pros: fast, no tuning required, highly interpretable, well-understood

Cons: unlikely to produce the best predictive accuracy (presumes a linear relationship between the features and response)

Form of linear regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

- y is the response
- β_0 is the intercept
- β_1 is the coefficient for x_1 (the first feature)
- β_n is the coefficient for x_n (the nth feature)

In this case:

$$y = \beta_0 + \beta_1 \times TV + \beta_2 \times Radio + \beta_3 \times Newspaper$$

The β values are called the **model coefficients**. These values are "learned" during the model fitting step using the "least squares" criterion. Then, the fitted model can be used to make predictions!

Form of linear regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

- y is the response
- β_0 is the intercept
- β_1 is the coefficient for x_1 (the first feature)
- β_n is the coefficient for x_n (the nth feature)

In this case:

$$y = \beta_0 + \beta_1 \times TV + \beta_2 \times Radio + \beta_3 \times Newspaper$$

The β values are called the **model coefficients**. These values are "learned" during the model fitting step using the "least squares" criterion. Then, the fitted model can be used to make predictions!

Preparing X and y using pandas

- scikit-learn expects X (feature matrix) and y (response vector) to be NumPy arrays.
- However, pandas is built on top of NumPy.
- Thus, X can be a pandas DataFrame and y can be a pandas Series!

create a Python list of feature names

```
feature_cols = ['TV', 'Radio', 'Newspaper']
```

use the list to select a subset of the original DataFrame

```
X = data[feature_cols]
```

equivalent command to do this in one line

```
X = data[['TV', 'Radio', 'Newspaper']]
```

print the first 5 rows

```
X.head()
```

check the type and shape of X

```
print(type(X))
```

```
print(X.shape
```

```
)
```

select a Series from the DataFrame

```
y = data['Sales']
```

equivalent command that works if there are no spaces in the column name

```
y = data.Sales
```

print the first 5 values

```
y.head()
```

check the type and shape of y

```
print(type(y))
```

```
print(y.shape)
```

Splitting X and y into training and testing sets

```
from sklearn.cross_validation import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=1)
```

default split is 75% for training and 25% for testing

```
print(X_train.shape)  
print(y_train.shape)  
print(X_test.shape)  
print(y_test.shape)
```

Linear regression in scikit-learn

import model

```
from sklearn.linear_model import LinearRegression
```

instantiate

```
linreg = LinearRegression()
```

fit the model to the training data (learn the coefficients)

```
linreg.fit(X_train, y_train)
```

Interpreting model coefficients

```
# print the intercept and coefficients
```

```
print(linreg.intercept_)
```

```
print(linreg.coef_)
```

```
# pair the feature names with the coefficients
```

```
list(zip(feature_cols, linreg.coef_))
```

Making predictions

```
# make predictions on the testing set
```

```
y_pred = linreg.predict(X_test)
```

$$y = 2.88 + 0.0466 \times TV + 0.179 \times Radio + 0.00345 \times Newspaper$$

Model evaluation metrics for regression

Mean Absolute Error

define true and predicted response values

```
true = [100, 50, 30, 20]
```

```
pred = [90, 50, 50, 30]
```


Mean Absolute Error (MAE) is the mean of the absolute value of the errors:

$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Mean Absolute Error (MAE)

calculate MAE by hand

```
print((10 + 0 + 20 + 10)/4.)
```

calculate MAE using scikit-learn

```
from sklearn import metrics
```

```
print(metrics.mean_absolute_error(true, pred))
```

Mean Squared Error (MSE) is the mean of the squared errors:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Mean Squared Error (MSE)

calculate MSE by hand

```
print((10**2 + 0**2 + 20**2 + 10**2)/4.)
```

calculate MSE using scikit-learn

```
print(metrics.mean_squared_error(true, pred))
```

Root Mean Squared Error (RMSE) is the square root of the mean of the squared errors:

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Root Mean Squared Error (RMSE)

calculate RMSE by hand

```
import numpy as np
```

```
print(np.sqrt((10**2 + 0**2 + 20**2 + 10**2)/4.))
```

calculate RMSE using scikit-learn

```
print(np.sqrt(metrics.mean_squared_error(true, pred)))
```

#For our problem

```
print(np.sqrt(metrics.mean_squared_error(y_test, y_pred)))
```

Feature selection

create a Python list of feature names

```
feature_cols = ['TV', 'Radio']
```

use the list to select a subset of the original DataFrame

```
X = data[feature_cols]
```

select a Series from the DataFrame

```
y = data.Sales
```

split into training and testing sets

```
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=1)
```

fit the model to the training data (learn the coefficients)

```
linreg.fit(X_train, y_train)
```

make predictions on the testing set

```
y_pred = linreg.predict(X_test)
```

compute the RMSE of our predictions

```
print(np.sqrt(metrics.mean_squared_error(y_test, y_pred)))
```