```python
"""
Example 2.1 Linear Regression
@author: Austin Downey
"""

import IPython as IP
IP.get_ipython().run_line_magic('reset', '-sf')

import numpy as np
import matplotlib.pyplot as plt
import sklearn as sk
from sklearn import datasets
from sklearn.linear_model import LinearRegression

plt.close('all')

#%% load data

ames = sk.datasets.fetch_openml(name="house_prices", as_frame=True,parser='auto')
target = ames['target'].values
data = ames['data']
YrSold = data['YrSold'].values  # Year Sold (YYYY)
MoSold = data['MoSold'].values  # Month Sold (MM)
OverallCond = data['OverallCond'].values # OverallCond: Rates the overall condition of
the house
GrLivArea = data['GrLivArea'].values # Above grade (ground) living area square feet
BedroomAbvGr = data['BedroomAbvGr'].values # Bedrooms above grade (does NOT include
basement bedrooms)

# Ask a home buyer to describe their dream house, and they probably won't begin
# with the height of the basement ceiling or the proximity to an east-west railroad.
# But this playground competition's dataset proves that much more influences price
# negotiations than the number of bedrooms or a white-picket fence.

# With 79 explanatory variables describing (almost) every aspect of residential
# homes in Ames, Iowa, this competition challenges you to predict the final price
# of each home.

# Plot a few of the interesting features vs the target (price). In particular,
# let's plot the number of rooms vs. the price.
plt.figure()
plt.plot(GrLivArea,target,'o',markersize=2)
plt.xlabel('Above grade (ground) living area square feet')
plt.ylabel('price (USD)')
plt.grid(True)
plt.tight_layout()

#%% Build a model for the data
X = GrLivArea
Y = target
model_X = np.linspace(0,5000)

theta_1 = 0
theta_2 = 100
model_Y_manual = theta_1 + theta_2*model_X

plt.figure()
plt.plot(X,Y,'o',markersize=2,label='data')
plt.plot(model_X,model_Y_manual,'--',label='manual fit')
plt.xlabel('Above grade (ground) living area square feet')
plt.ylabel('price (USD)')
plt.grid(True)
#plt.xlim([3.5,9])
#plt.ylim([0,50000])
plt.legend(framealpha=1)
plt.tight_layout()
```

```python
66    # add a dimension to the data as math is easier in 2d arrays and sk learn only
67    # takes 2d arrays
68    X = np.expand_dims(X,axis=1)
69    Y = np.expand_dims(Y,axis=1)
70    model_X = np.expand_dims(model_X,axis=1)
71
72    #%% compute the linear regression solution using the closed form solution
73
74    # compute
75    X_b = np.ones((X.shape[0],2))
76    X_b[:,1] = X.T # add x0 = 1 to each instance
77
78    theta_closed_form = np.linalg.inv(X_b.T@X_b)@X_b.T@Y
79
80    model_y_closed_form = theta_closed_form[0] + theta_closed_form[1]*3000
81    model_Y_closed_form = theta_closed_form[0] + theta_closed_form[1]*model_X
82
83    plt.figure()
84    plt.plot(X,Y,'o',markersize=3,label='data points')
85    plt.xlabel('Above grade (ground) living area square feet')
86    plt.ylabel('price (USD)')
87    plt.plot(3000,model_y_closed_form,'dr',markersize=10,zorder=10,
88             label='inferred data point')
89    plt.plot(model_X,model_Y_closed_form,'-',label='normal equation')
90    plt.grid(True)
91    plt.legend()
92    plt.tight_layout()
93
94    #%% compute the linear regression solution using gradient descent
95
96    eta = 0.00000001 # learning rate
97    n_iterations = 100
98    m = X.shape[0]
99    theta_gradient_descent = np.random.randn(2,1) # random initialization
100   for iteration in range(n_iterations):
101       gradients = 2/m * X_b.T.dot(X_b.dot(theta_gradient_descent) - Y)
102       theta_gradient_descent = theta_gradient_descent - eta * gradients
103
104   print(theta_gradient_descent)
105
106   model_Y_gradient_descent =  theta_gradient_descent[0] \
107       + theta_gradient_descent[1]*model_X
108
109   plt.figure()
110   plt.plot(X,Y,'o',markersize=3,label='data points')
111   plt.xlabel('Above grade (ground) living area square feet')
112   plt.ylabel('price (USD)')
113   plt.plot(model_X,model_Y_closed_form,'-',label='normal equation')
114   plt.plot(model_X,model_Y_gradient_descent,':',label='gradient descent')
115   plt.grid(True)
116   plt.legend()
117   plt.tight_layout()
118
119   #%% compute the linear regression solution using sk-learn
120
121   # build and train a closed from linear regression model in sk-learn
122   model_LR = sk.linear_model.LinearRegression()
123   model_LR.fit(X,Y[:,0])
124   model_Y_sk_LR = model_LR.predict(model_X)
125
126   # build and train a Stochastic Gradient Descent linear regression model in sk-learn.
127   # Note that in running the model, the best way to do this would be to use a pipeline
128   # =with feature scaling. However, here we just set 'eta0' to a low value, this
129   # is done only for educational # purposes and is not the ideal methodology in
130   # terms of system robustness.
131   model_SGD = sk.linear_model.SGDRegressor(learning_rate='constant',eta0=0.00000001)
132   model_SGD.fit(X,Y[:,0])
```

```python
133    model_Y_sk_SGD = model_SGD.predict(model_X)
134
135    # plot the modeled results
136    plt.figure()
137    plt.plot(X,Y,'o',markersize=2,label='data')
138    plt.plot(model_X,model_Y_closed_form,'-',label='normal equation')
139    plt.plot(model_X,model_Y_gradient_descent,'--',label='gradient descent')
140    plt.plot(model_X,model_Y_sk_LR,':',label='sklearn normal equation')
141    plt.plot(model_X,model_Y_sk_SGD,'-.',label='sklearn stochastic gradient descent')
142
143    plt.xlabel('Above grade (ground) living area square feet')
144    plt.ylabel('price (USD)')
145    plt.grid(True)
146    plt.legend(framealpha=1)
147    plt.tight_layout()
```