

Linear Regression Analysis

- Least Sum of Squares Approach -

Tseng-Ching James Shen, Ph.D.

Motivation of Linear Regression

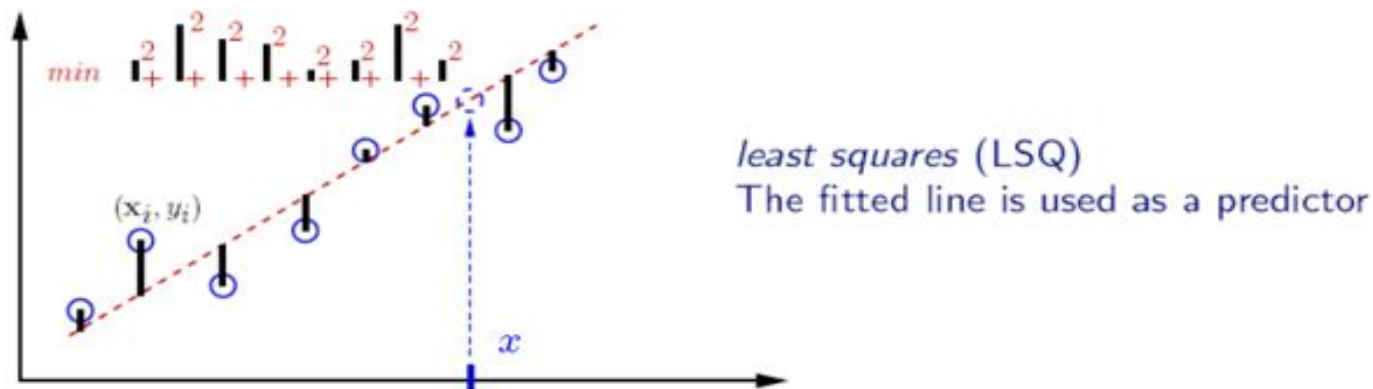
- Given two random variables x , and y , we use the following linear equation to model the correlation between two random variables

$$y = w x + b$$

- The coefficient w and correlation are related in the following way:
 - $w > 0$, positive correlation
 - $w = 0$, no correlation
 - $w < 0$, negative correlation

Motives of Linear Regression (cont.)

- The parameters w and b are unknown and can be determined by fitting the model with the dataset
- Fit model by minimizing the sum of squared errors (i.e. training errors)



Linear Regression Model

- Features are random variables x_1, x_2, \dots, x_d
- The model can be generalized to model response y is also a random variable as follows

$$y = w_0 + w_1 x_1 + \dots + w_d x_d$$

- It can be also expressed in terms of vector notation

$$y = h(x) = \mathbf{x}^T \mathbf{w} + w_0$$

Linear Regression Model (cont.)

- Given a set of observed dataset points, we need to determine the coefficients which are best to fit the dataset points $(\mathbf{x}_i, y_i) \mid = 1, \dots, N$
- For dataset point (\mathbf{x}_i, y_i) (*revised version*), the estimated y_i

$$h(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{w}$$

$$\mathbf{x}_i = (1, x_{i0}, x_{i1}, \dots, x_{id}), \mathbf{w} = (w_0, w_1, \dots, w_d)$$

- Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$

$$(h(\mathbf{x}_1), \dots, h(\mathbf{x}_N))^T = \mathbf{X} \mathbf{w}$$

Optimization Problem

- We can use training errors as the objective function to be minimized

$$\frac{1}{N} \sum_{i=1}^N (h(\mathbf{x}_i) - y_i)^2$$

- Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^\top$

$$(h(\mathbf{x}_1), \dots, h(\mathbf{x}_N))^\top = \mathbf{X} \mathbf{w}$$

- The objective function can be re-written as

$$\mathcal{E}(\mathbf{w}) = \frac{1}{N} (\mathbf{X} \mathbf{w} - \mathbf{y})^\top (\mathbf{X} \mathbf{w} - \mathbf{y})$$

Closed Form Solution

- Expand the objective function

$$\mathcal{E}(\mathbf{w}) = \frac{1}{N} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \frac{2}{N} \mathbf{w}^T \mathbf{X}^T \mathbf{y} + \frac{1}{N} \mathbf{y}^T \mathbf{y}$$

- Find the solution which satisfies the following equation

$$\frac{\partial \mathcal{E}}{\partial \mathbf{w}} = 0$$

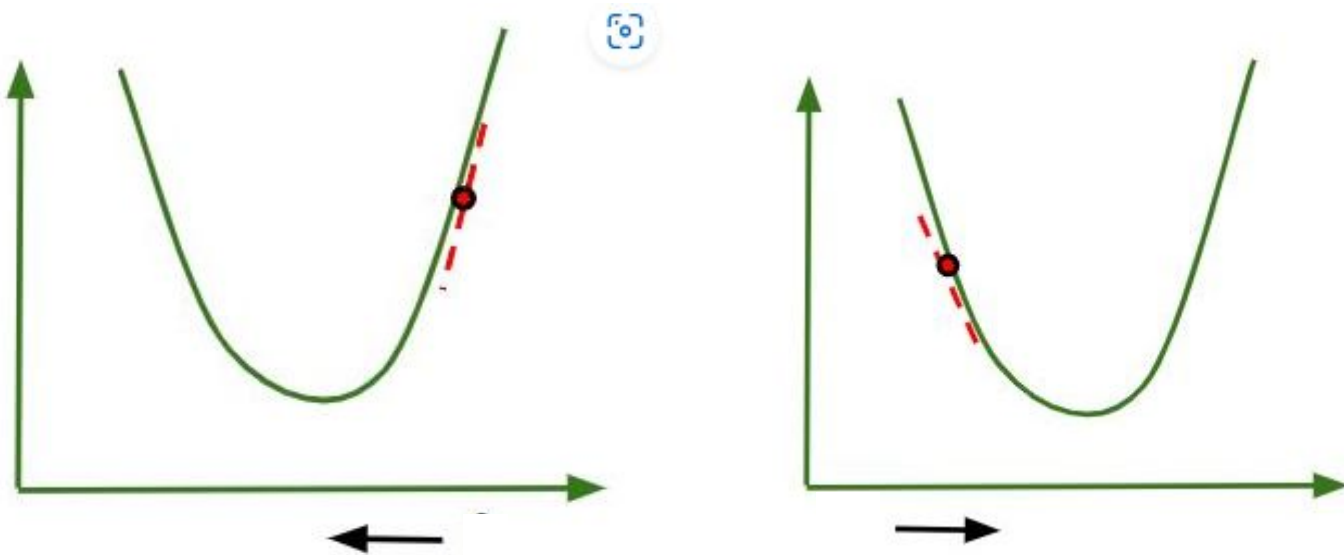
which leads to $\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$

- Therefore, the solution is $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Use Python NumPy to implement the solution

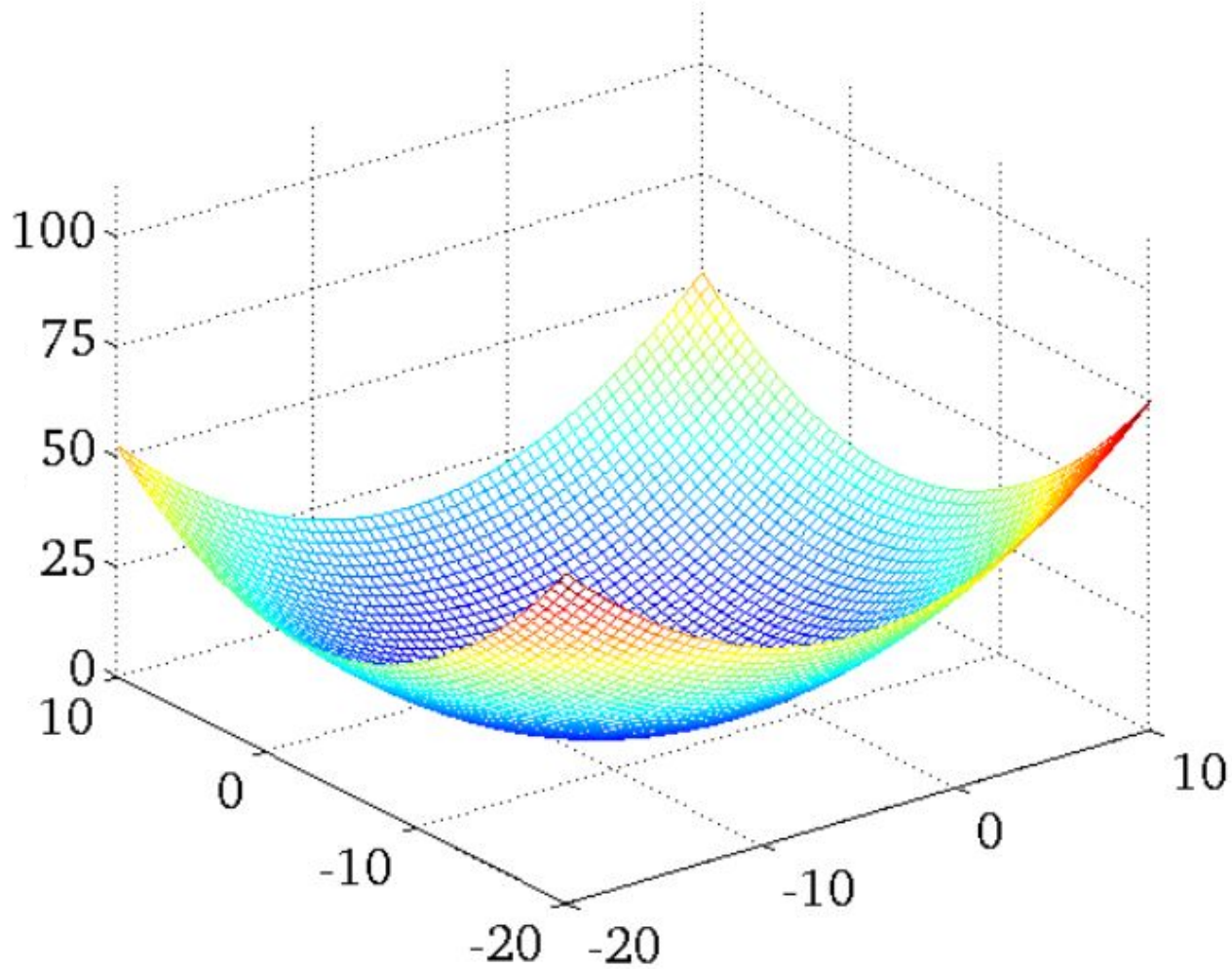
```
alpha = np.dot((np.dot(np.linalg.inv(np.dot(A.T,A)),A.T)),y)
```

Gradient Descent

- The objective function is a differential and convex function
- Gradient descent is an iterative process
 - Pick a starting point
 - At each iteration move to the next point by applying the slope of current point (derivative of the function)



3-D Convex Function



Derivative of Function

- The derivative of the function is

$$\frac{\partial}{\partial \mathbf{w}} \mathcal{L}(\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} \frac{1}{N} \sum_{i=1}^N (h(\mathbf{x}_i) - y_i)^2$$

- Since $\frac{\partial}{\partial \mathbf{w}} \mathcal{L}(\mathbf{w}) = (\frac{\partial}{\partial w_0} \mathcal{L}(\mathbf{w}), \dots, \frac{\partial}{\partial w_d} \mathcal{L}(\mathbf{w}))$, we will find $\frac{\partial}{\partial w_j} \mathcal{L}(\mathbf{w})$ $j = 0, 1, 2, \dots, d$

- The equation for $\frac{\partial}{\partial w_j} \mathcal{L}(\mathbf{w})$

$$\frac{\partial}{\partial w_j} \mathcal{L}(\mathbf{w}) = \frac{\partial}{\partial w_j} \frac{1}{N} \sum_{i=1}^N (\sum_{k=0}^d w_k x_{ik} - y_i)^2$$

Derivation of Function

- Apply the Derivative Chain Rule:

$$\frac{\partial}{\partial w_j} \mathcal{L}(\mathbf{w}) = \frac{2}{N} \sum_{i=1}^N \left(\sum_{k=0}^{k=d} w_d x_{id} - y_i \right) \frac{\partial}{\partial w_j} \left(\sum_{k=0}^{k=d} w_d x_{id} - y_i \right)$$

- We will get the following equation

$$\frac{\partial}{\partial w_j} \mathcal{L}(\mathbf{w}) = \frac{2}{N} \sum_{i=1}^N \left(\sum_{k=0}^{k=d} w_d x_{id} - y_i \right) x_{ij}$$

Gradient Descent Algorithm

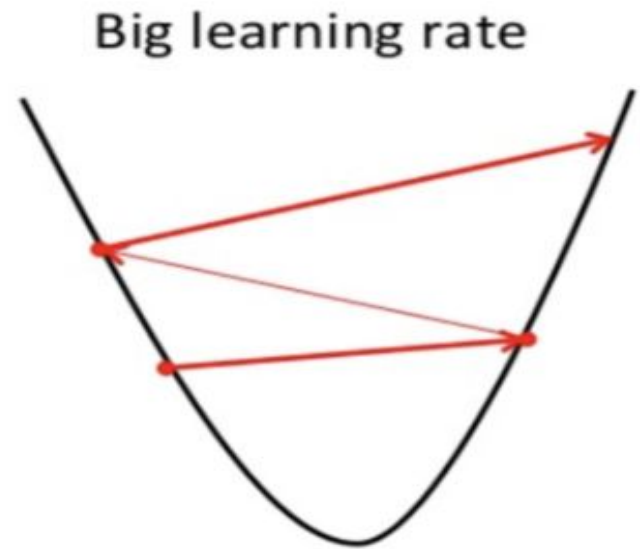
- Initialize \mathbf{w}
- Update w_j at each iteration using the following equation until \mathbf{w} converges

$$w_j \leftarrow w_j + \alpha \left(\frac{2}{N} \sum_{i=1}^N \left(\sum_{k=0}^d w_k x_{ik} - y_i \right) x_{ij} \right)$$

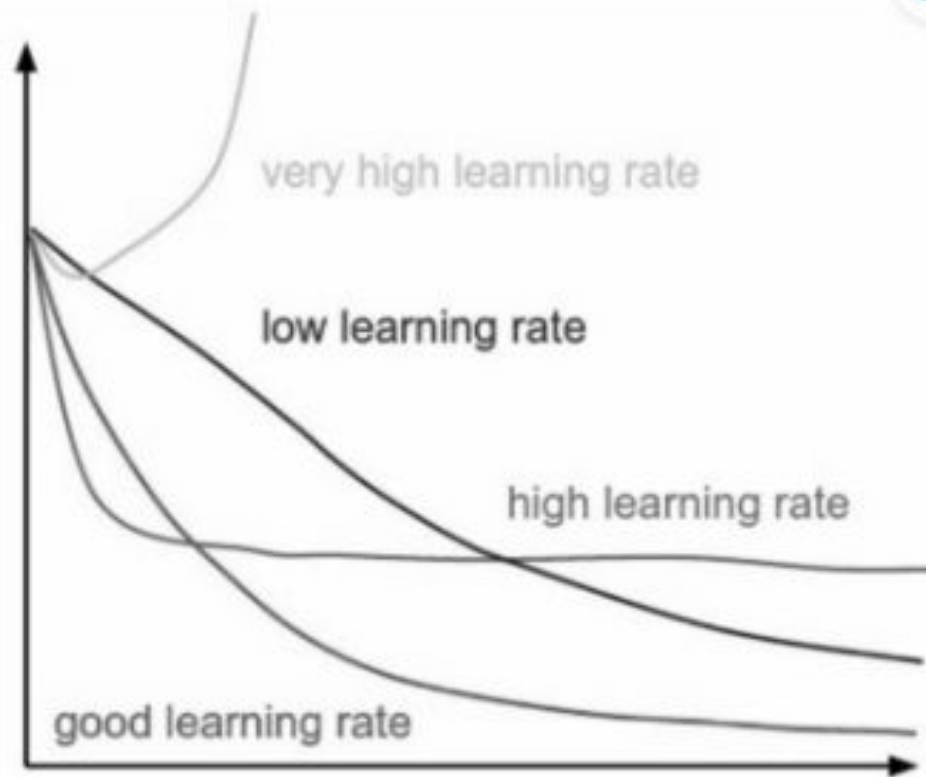
where α is the *learning rate*

The Learning Rate

- If the learning rate is too high, we might **OVERSHOOT** the minima and keep bouncing, without reaching the minima
- If the learning rate is too small, the training might turn out to be too long



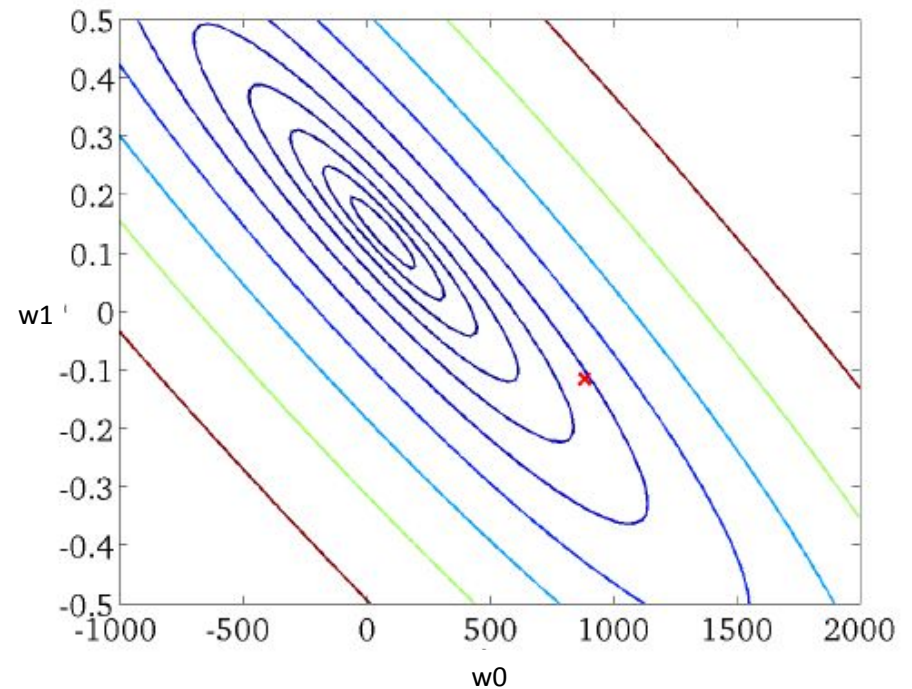
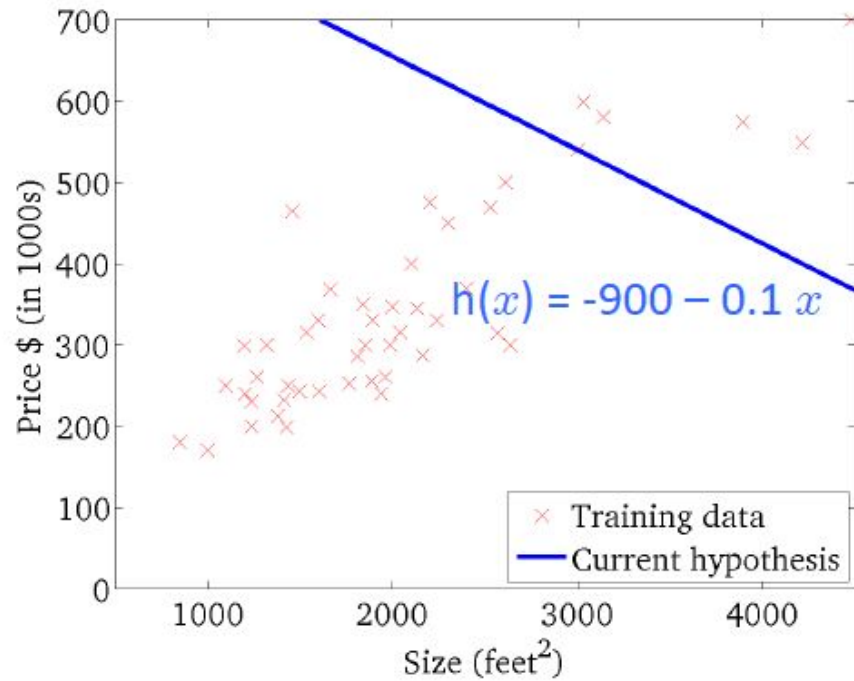
Learning Rate (cont.)



Loss value v.s. the number of iterations

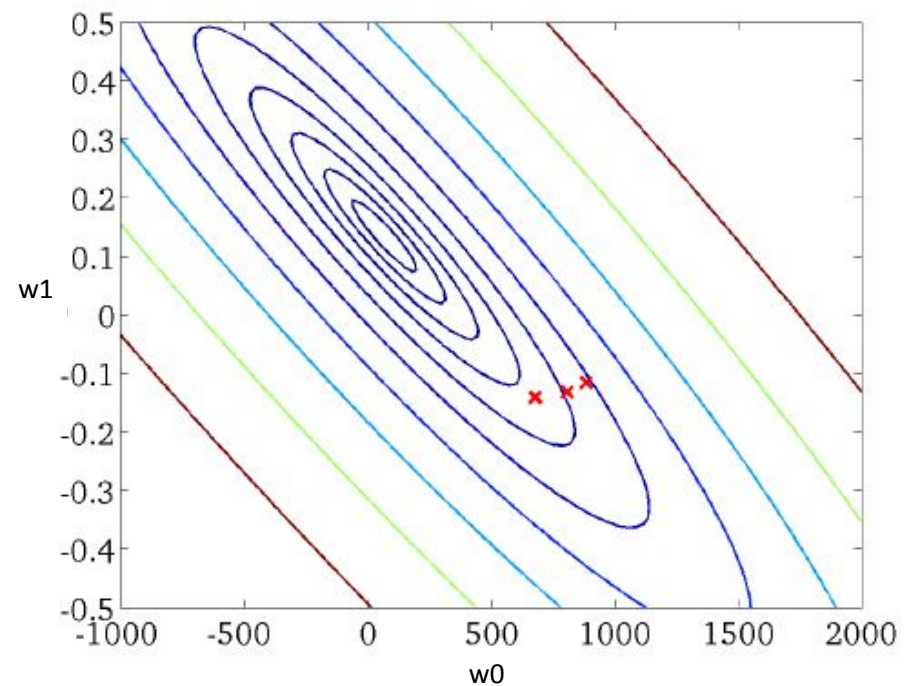
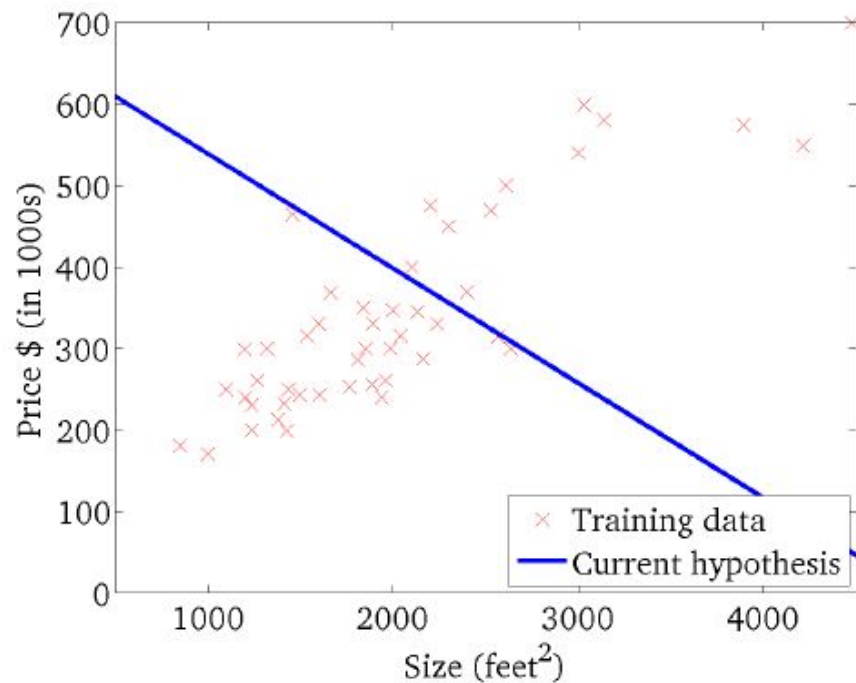
Gradient Descent

Initial Iteration



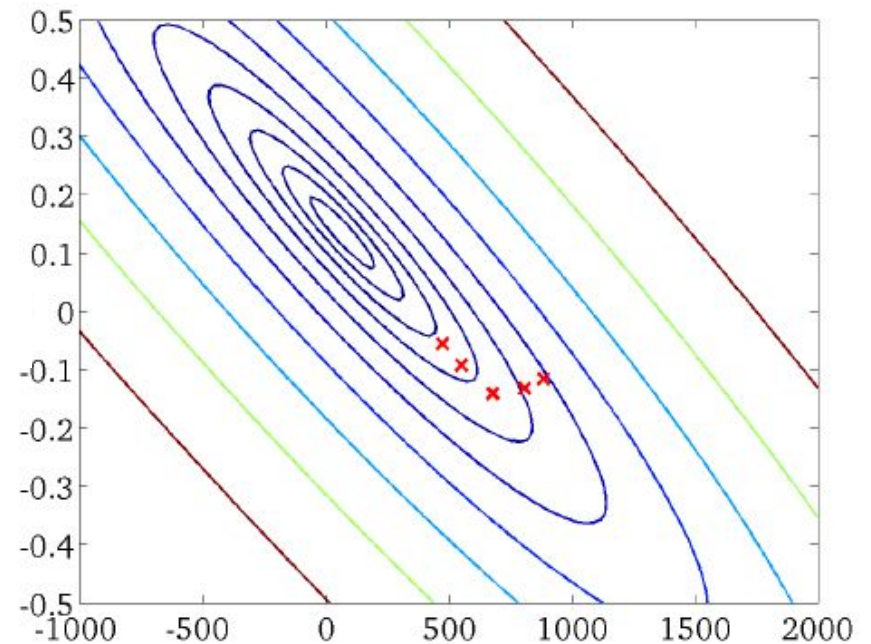
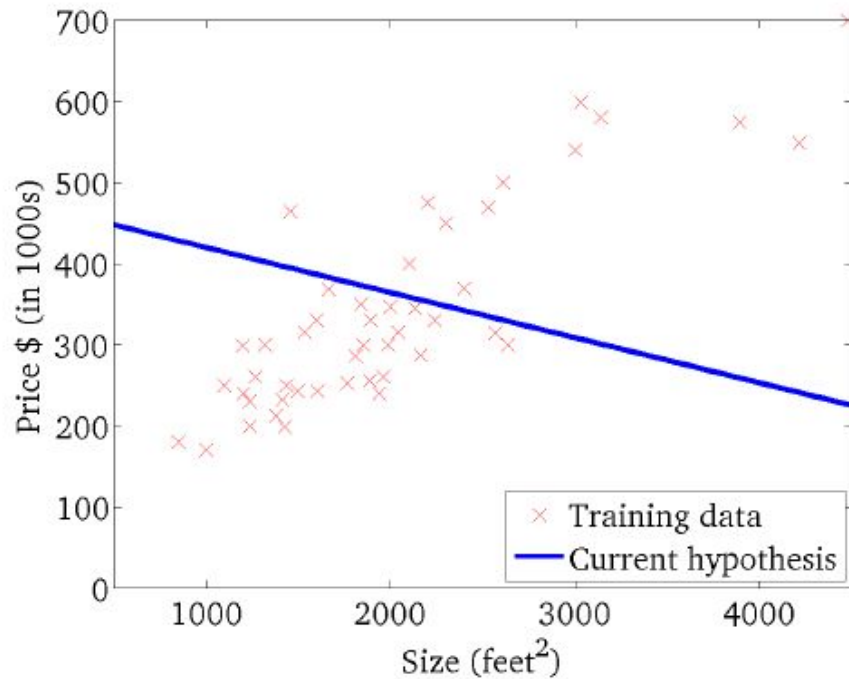
Gradient Descent

The 3rd Iteration



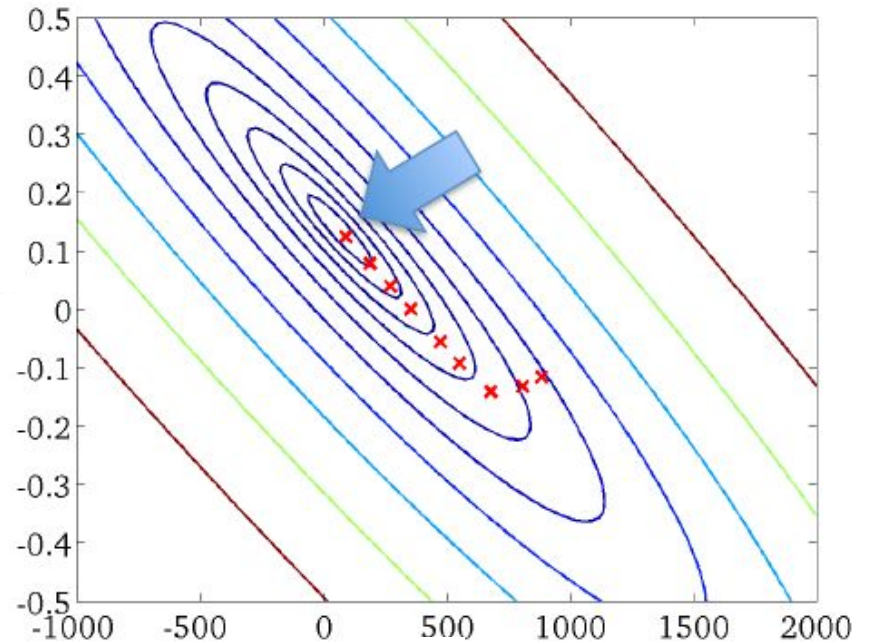
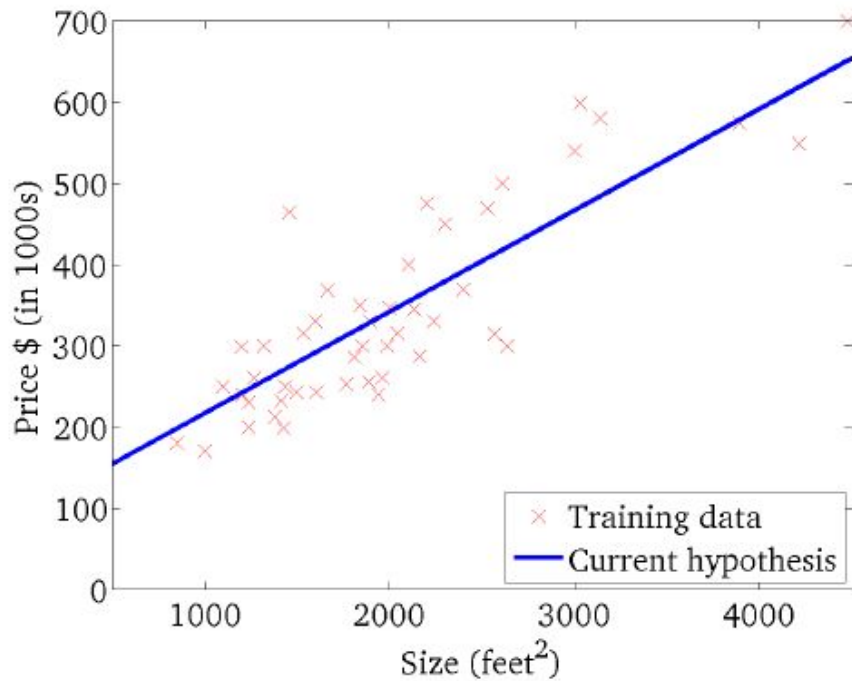
Gradient Descent

The 5th Iteration



Gradient Descent

The 9th Iteration



Example

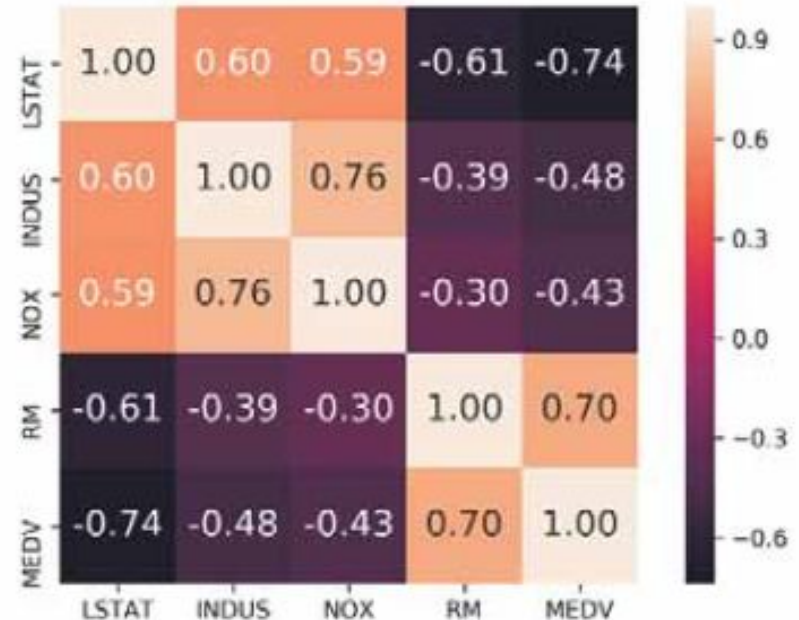
- Chapter 10 in Python Machine Learning, Sebastian Raschka Vahid Mirjalili Packt Publishing Ltd.
- Data preparation

```
>>> import pandas as pd
>>> df = pd.read_csv('https://raw.githubusercontent.com/rasbt/'
...                  'python-machine-learning-book-2nd-edition/'
...                  '/master/code/ch10/housing.data.txt',
...                  header=None,
...                  sep='\s+')
>>> df.columns = ['CRIM', 'ZN', 'INDUS', 'CHAS',
...               'NOX', 'RM', 'AGE', 'DIS', 'RAD',
...               'TAX', 'PTRATIO', 'B', 'LSTAT', 'MEDV']
>>> df.head()
```

Example (cont.)

- Correlation analysis

```
>>> import numpy as np
>>> cm = np.corrcoef(df[cols].values.T)
>>> sns.set(font_scale=1.5)
>>> hm = sns.heatmap(cm,
...                   cbar=True,
...                   annot=True,
...                   square=True,
...                   fmt='.2f',
...                   annot_kws={'size': 15},
...                   yticklabels=cols,
...                   xticklabels=cols)
>>> plt.show()
```



Example (cont.)

- Feature selection
- Feature scaling

```
>>> X = df[['RM']].values
>>> y = df['MEDV'].values
>>> from sklearn.preprocessing import StandardScaler
>>> sc_x = StandardScaler()
>>> sc_y = StandardScaler()
>>> X_std = sc_x.fit_transform(X)
>>> y_std = sc_y.fit_transform(y[:, np.newaxis]).flatten()
```

Example (cont.)

- Implementation of Gradient Descent

```
class LinearRegressionGD(object):

    def __init__(self, eta=0.001, n_iter=20):
        self.eta = eta
        self.n_iter = n_iter

    def fit(self, X, y):
        self.w_ = np.zeros(1 + X.shape[1])
        self.cost_ = []

        for i in range(self.n_iter):
            output = self.net_input(X)
            errors = (y - output)
            self.w_[1:] += self.eta * X.T.dot(errors)
            self.w_[0] += self.eta * errors.sum()
            cost = (errors**2).sum() / 2.0
            self.cost_.append(cost)
        return self

    def net_input(self, X):
        return np.dot(X, self.w_[1:]) + self.w_[0]

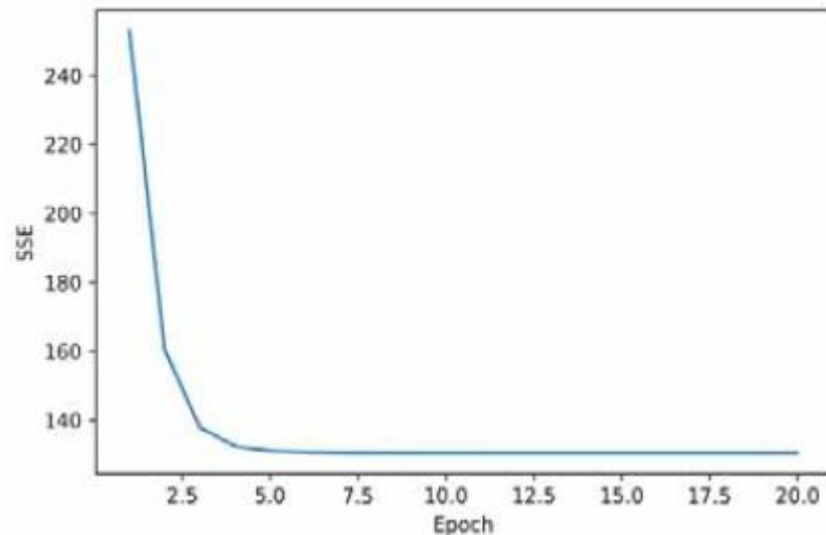
    def predict(self, X):
        return self.net_input(X)
```

Example (cont.)

- Usage of GD method

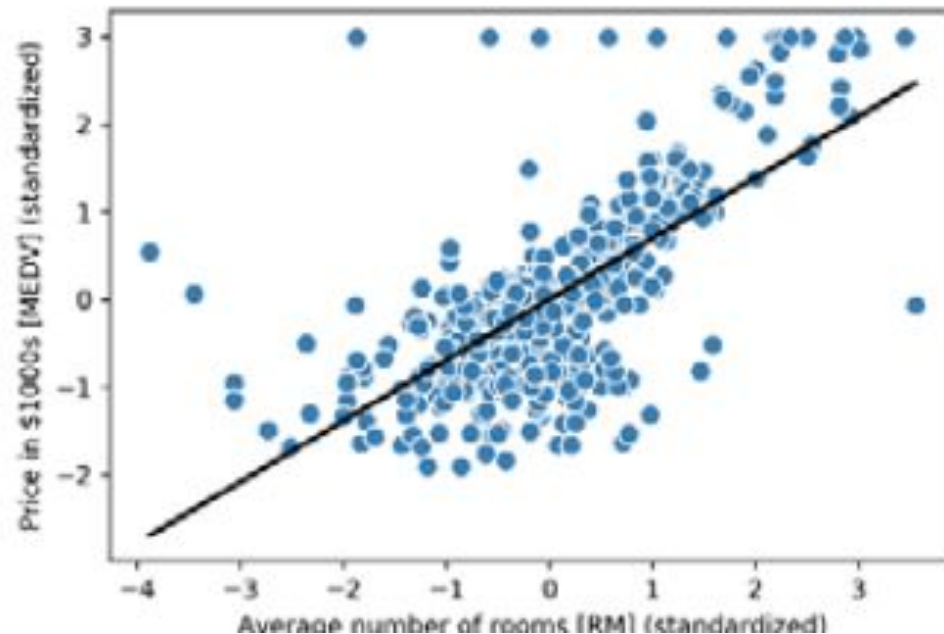
```
>>> lr = LinearRegressionGD()  
>>> lr.fit(X_std, y_std)
```

```
>>> sns.reset_orig() # resets matplotlib style  
>>> plt.plot(range(1, lr.n_iter+1), lr.cost_)  
>>> plt.ylabel('SSE')  
>>> plt.xlabel('Epoch')  
>>> plt.show()
```



Example (cont.)

```
>>> lin_regplot(X_std, y_std, lr)
>>> plt.xlabel('Average number of rooms [RM] (standardized)')
>>> plt.ylabel('Price in $1000s [MEDV] (standardized)')
>>> plt.show()
```



Example (cont.)

- Scikit-learn library

```
>>> from sklearn.linear_model import LinearRegression
>>> slr = LinearRegression()
>>> slr.fit(X, y)
>>> print('Slope: %.3f' % slr.coef_[0])
Slope: 9.102
>>> print('Intercept: %.3f' % slr.intercept_)
Intercept: -34.671
```

Gradient Descent vs Closed Form

Gradient Descent

- Requires multiple iterations
- Need to choose α
- Works well when n is large
- Can support incremental learning

Closed Form Solution

- Non-iterative
- No need for α
- Slow if n is large
 - Computing $(X^T X)^{-1}$ is roughly $O(n^3)$

Kernel Regression Analysis

- Linear functions might not be good enough for some datasets
- Kernel method could be used to fit dataset with nonlinear functions
- Apply the kernel method and trick to get

$$y_i = \sum_{j=1}^N \alpha_j \phi^\top(\mathbf{x}_j) \phi(\mathbf{x}_i)$$

- Kernel matrix $\mathbf{K}(i, j) = \phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j)$