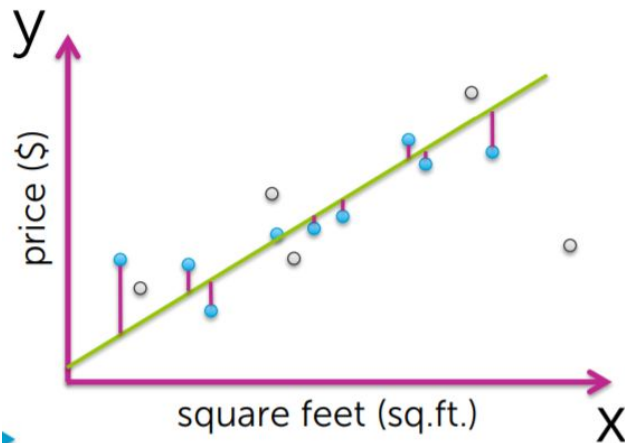

Introduction to Machine Learning: Regression

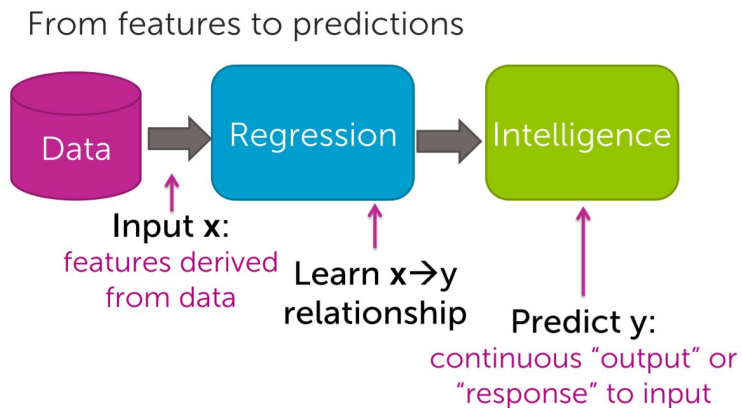
By Shrishty Chandra

What is regression ?

Regression is a **model**, e.g. a model that assumes a linear relationship between the input variables (x) and the single output variable (y). More specifically, that y can be calculated from a combination of the input variables (x).



Train Data + Test Data





Usecase

Predicting House Prices: Given features of the house, predict the price of it.

→ **$f(x=\text{sq.ft}, y=\$)$**

We will define a function which takes sq.ft of the house and predict the house price.

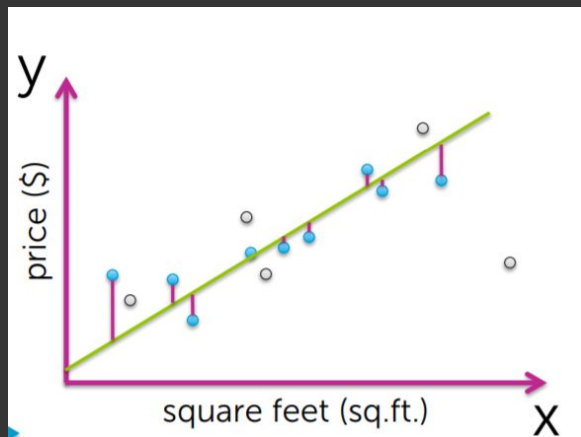
→ **$f(x1=\text{sq.ft}, x2=\text{\#bedrooms}, y=\$)$**

Given multiple features, predict the price of the house

→ **Feature selection and performance (Not talking but Imp)**

There are 100,000 features in a house, how to choose the correct features

Regression Model



$$y_i = f(x_i) + e_i$$



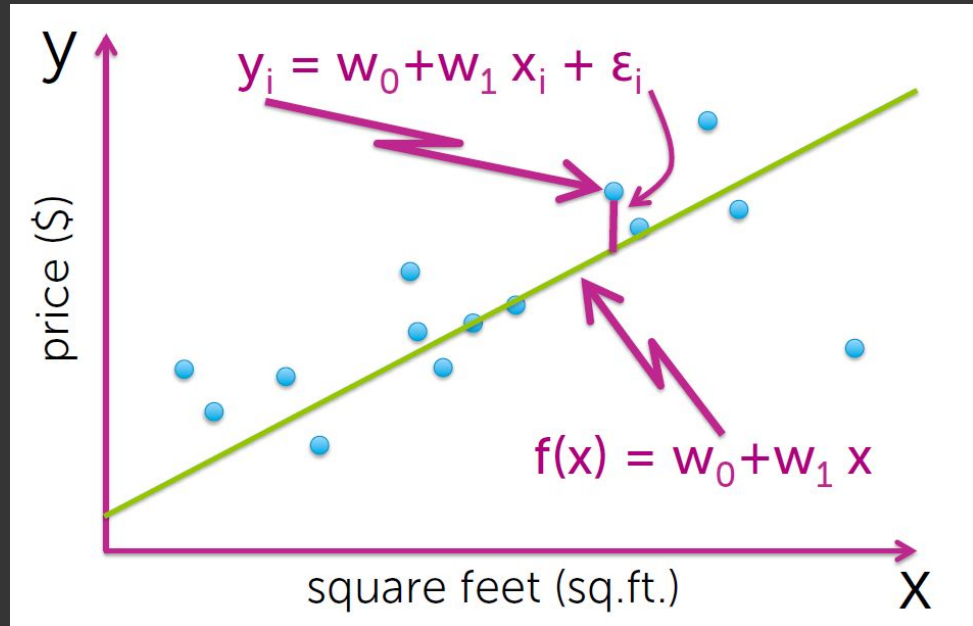
Tip

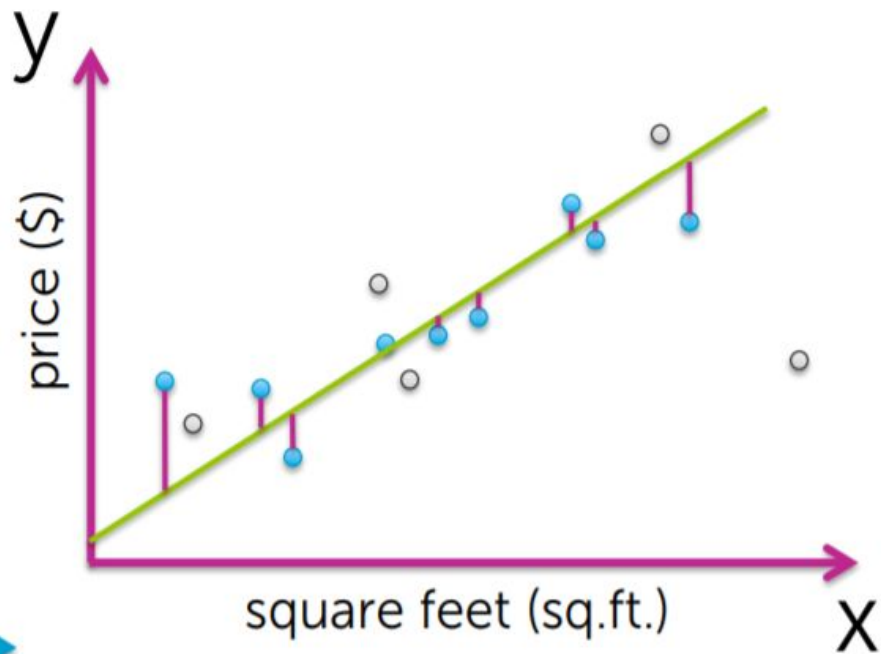
Essentially all models are wrong but some are useful.

George Box, 1987

Simple Linear regression model

- W_0 and W_1 are regression coefficients
- ε_i is the error.
- $f(x)$ is the fitted line in the data.





Residual Sum of Squares

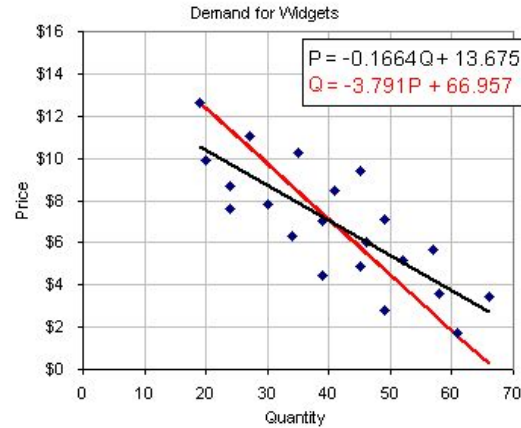
Sum of squares of the errors in predictions.

$$RSS = \sum_{i=1}^n (y_i - f(x_i))^2$$

$$RSS(w_0, w_1) = \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

Fitting the best line

We need to minimize the cost over all possible $\mathbf{w}_0, \mathbf{w}_1$, to fit the best line.

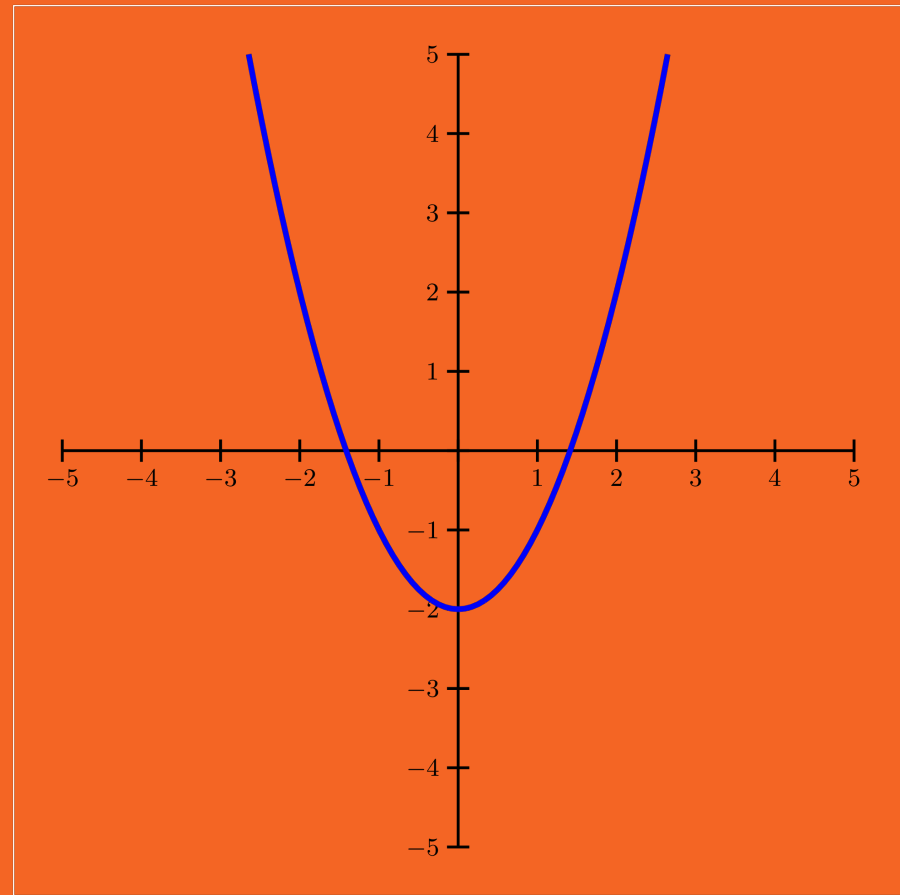


RSS is a quadratic equation with respect to W_0 and also W_1 .

When we are calculating W_0 such that it minimizes RSS we will consider all other variables as constant

Similar for W_1

So where do you think this function minimum here ?



Two ways to calculate the minimum

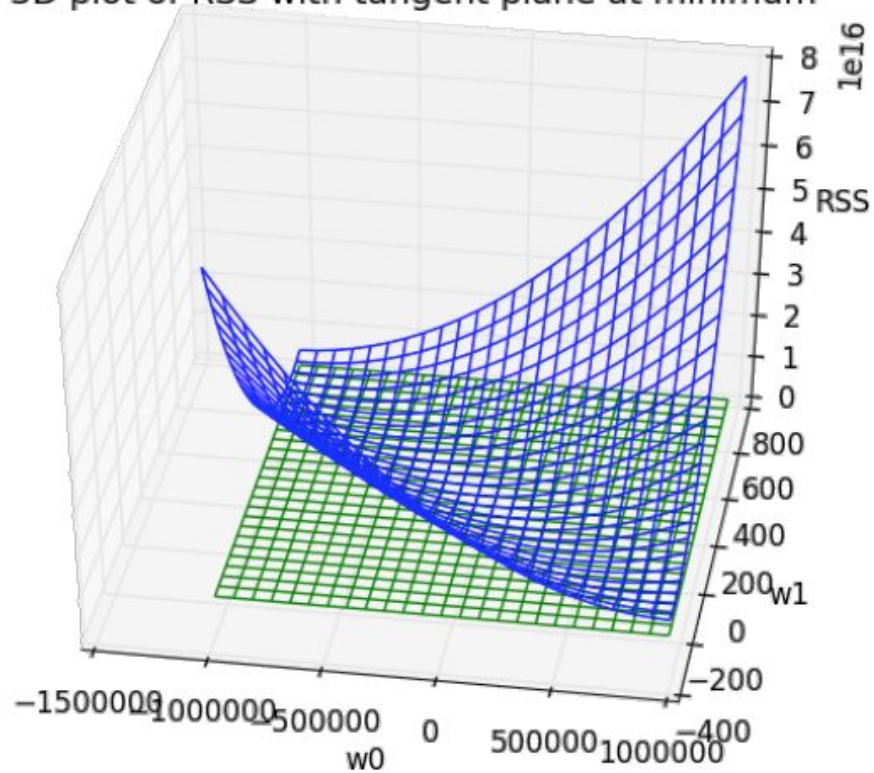
- Find the point where the derivative is zero
- Gradient descent
 - η is the step size here, represents how fast move towards the optimum w
 - Common choices of η are
 - $\eta_t = \alpha/t$
 - $\eta_t = \alpha/\text{sqrt}(t)$

Hill Descent algorithm

Algorithm:

while not converged
 $w^{(t+1)} \leftarrow w^{(t)} - \eta \left. \frac{dg}{dw} \right|_{w^{(t)}}$

3D plot of RSS with tangent plane at minimum



Multidimensional view of gradients

$$g(w) = 5W_0 + 10W_0W_1 + 2W_1^2$$

$$\partial g / \partial w_0 = 5 + 10W_1$$

$$\partial g / \partial w_1 = 4W_1 + 10W_0$$

$$\nabla g(w) = [\partial g / \partial w_0, \partial g / \partial w_1]$$

So if,

$$\text{RSS}(w_0, w_1) = \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

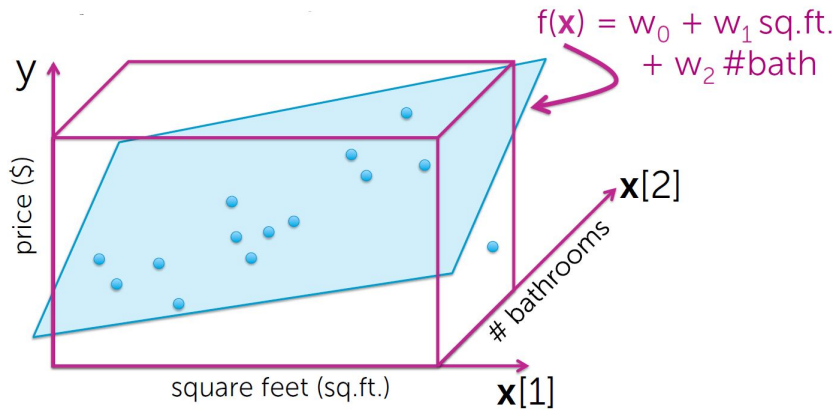
Then,

$$\nabla \text{RSS}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \\ -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] x_i \end{bmatrix}$$

After setting the gradient to zero. We will be able to calculate w_0 and w_1

Comparing the derivative and gradient approach

- Most ML problems cannot solve $\text{gradient} = 0$
- Even if solving $\text{gradient} = 0$ is feasible, gradient descent can be more efficient
- But gradient descent depends on step size and convergence criteria



Multiple Regression

Linear regression with multiple features

There are many possible inputs

- Sqft
- #bedrooms
- #bathrooms etc

General notation

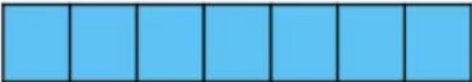

- Output: y (Scalar)
- Inputs: $x = (x[1], x[2], x[3], x[4], \dots)$
 - D -dimensional vector
 - D is number of features
- $x[j]$ = j th input (scalar)
- $h_j(x)$ = j th feature (scalar)
- x_i = input of i th data point (vector)
- $x_i[j]$ = j th input of i th data point (scalar)

Model:

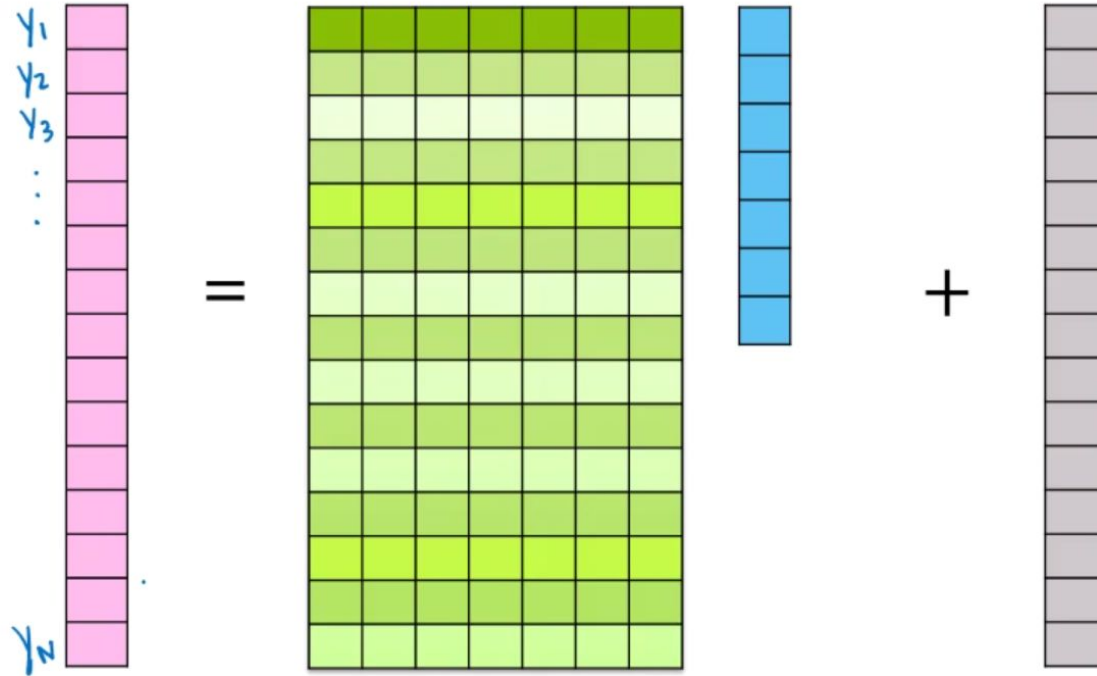
$$\begin{aligned} y_i &= w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i) + \varepsilon_i \\ &= \sum_{j=0}^D w_j h_j(\mathbf{x}_i) + \varepsilon_i \end{aligned}$$

In Matrices

$$y_i = \sum_{j=0}^D w_j h_j(\mathbf{x}_i) + \varepsilon_i$$

$y_i =$   $+$ ε_i

Calculating all house prices at once



Calculating RSS in Matrix

residual ₁	residual ₂	residual ₃	...	residual _N
-----------------------	-----------------------	-----------------------	-----	-----------------------

$$\begin{aligned}\text{RSS}(\mathbf{w}) &= \sum_{i=1}^N (y_i - h(\mathbf{x}_i)^T \mathbf{w})^2 \\ &= (\mathbf{y} - \mathbf{H}\mathbf{w})^T (\mathbf{y} - \mathbf{H}\mathbf{w})\end{aligned}$$

residual ₁
residual ₂
residual ₃
...
residual _N

The Gradient

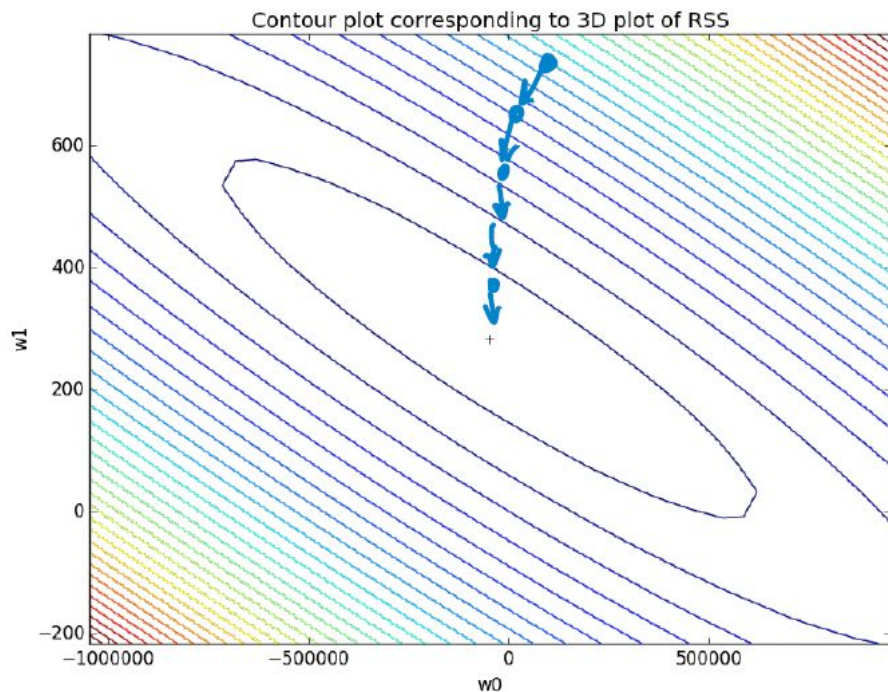
$$\begin{aligned}\nabla_{\text{RSS}}(\mathbf{w}) &= \nabla [(\mathbf{y} - \mathbf{H}\mathbf{w})^T (\mathbf{y} - \mathbf{H}\mathbf{w})] \\ &= -2\mathbf{H}^T (\mathbf{y} - \mathbf{H}\mathbf{w})\end{aligned}$$

—

$$\nabla_{\text{RSS}(\mathbf{w})} = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) = 0$$

Approach 1

- Set the gradient to zero
- Then solve for \mathbf{W}



Approach 2

- Use Gradient descent to calculate the optimum w

while not converged

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \underbrace{\nabla \text{RSS}(\mathbf{w}^{(t)})}_{-2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w})}$$

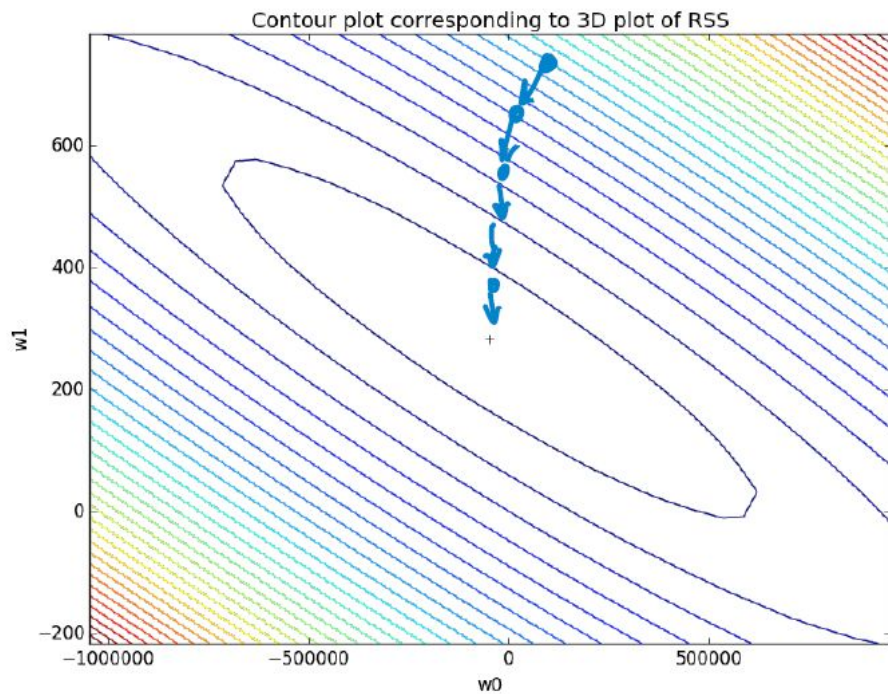


Reference

- <https://www.coursera.org/learn/ml-regression/home/welcome>

Other good resources:

- Machine Learning Recipes:
<https://www.youtube.com/watch?v=cKxRvEZd3Mw>
- Tensorflow,
[https://www.youtube.com/watch?v=g-EvyKpZjmQTens
orflow,](https://www.youtube.com/watch?v=g-EvyKpZjmQTensorflow)
<https://www.youtube.com/watch?v=g-EvyKpZjmQ>
- Neural Network,
[https://www.youtube.com/watch?v=NfnWJUyUJJYU&lis
t=PLwQyVqI_3POsyBPRNUU_ryNfXzgfkW2p&index=1](https://www.youtube.com/watch?v=NfnWJUyUJJYU&list=PLwQyVqI_3POsyBPRNUU_ryNfXzgfkW2p&index=1)
- AndrewNG Machine Learning:
<https://see.stanford.edu/Course/CS229>



Lets See the code