

# BITS F464: Machine Learning

# 20

## Linear Model for Regression



**Dr. Kamlesh Tiwari**

Assistant Professor, Department of CSIS,  
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA

March 10, 2021

**ONLINE** (Campus @ BITS-Pilani Jan-May 2021)

<http://ktiwari.in/ml>

## Regression

Consider following data

$x_1$	$x_2$	$x_3$	$y$
10	50	20	1
11	31	22	0
11	12	15	0
20	55	20	1
23	41	27	1
31	12	35	0
13	18	12	0
21	55	16	0
32	56	27	1
12	22	11	?
?	?	?	?

$x_1$	$x_2$	$x_3$	$y$
10	50	20	10
11	31	22	12
11	12	15	4
20	55	20	22
23	41	27	1
31	12	35	9
13	18	12	23
21	55	16	16
32	56	27	22
12	22	11	?
?	?	?	?

What should be ?

## Regression

**Regression:** predicts the value of continuous target variable

- A simplest model for regression can be a linear combination of the input variables

$$Y(x, w) = w_0 + w_1 x_1 + \dots + w_n x_n$$

where input  $x$  is represented by a  $n$  dimensional feature vector  $(x_1, x_2, \dots, x_n)$

- It can be extended by considering linear combination of fixed nonlinear functions  $\phi$  (called basis functions)

$$Y(x, w) = w_0 + \sum_{i=1}^n w_i \phi_i(x)$$

- In short  $Y(x, w) = w^T \phi(x)$
- Objective is to choose  $w$ , such that it makes  $y(x^{(i)}, w)$  as close to  $y^{(i)}$  as possible

## Regression

- Determining  $w$ , is similar to solving a minimization problem. Let us define a **squared error cost function** as

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (y(x^{(i)}, w) - y^{(i)})^2$$

where  $m$  is number of training examples

- Then one have to minimize the value of  $J(w)$

$$\operatorname{argmin}_w J(w)$$

- Basic idea: Push  $w$ ; a bit against the direction of its gradient

## Example

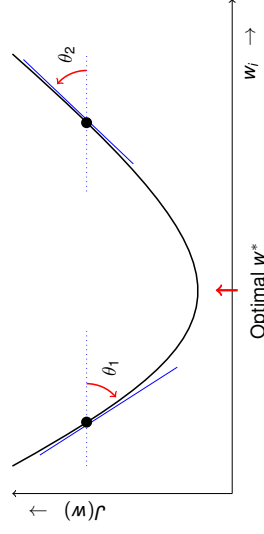
Assume for some chosen  $w$  we compute  $y(x, w)$

$x_1$	$x_2$	$x_3$	$y$	$\hat{y} = y(x, w)$	$(\hat{y} - y)^2$
10	50	20	10	8	4
11	31	22	12	9	9
11	12	15	4	3	1
20	55	20	22	26	16
23	41	27	1	1	0
31	12	35	9	4	25
13	18	12	23	30	49
21	55	16	16	13	9
32	56	27	22	21	1

Then

$$J(w) = \frac{1}{2 \times 9} \times 114 = 6.33$$

Consider  $w_i = w_i - \alpha \frac{\partial}{\partial w_i} J(w)$



- Slope  $\tan \theta_1$ , representing  $\frac{\partial}{\partial w_i} J(w)$  is  $-ve$  so the equation  $w_i = w_i - \alpha \frac{\partial}{\partial w_i} J(w)$  moves  $w_i$  towards  $w^*$
- $\tan \theta_2$ , being  $+ve$  the equation still moves  $w_i$  towards  $w^*$

## Gradient Descent

### Algorithm 1: Gradient Descent

- 1 Initialize  $w$  randomly
- 2 **repeat**
- 3 | Simultaneously update all  $w_j$  with  $w_j - \alpha \frac{\partial}{\partial w_j} J(w)$
- 4 **until** converge;
- 5 **return**  $w$

- Here  $\alpha$  is a learning rate. If  $\alpha$  is small enough then  $J(w)$  would decrease in every iteration (large  $\alpha$  can overshoot the minimum and may fail to converge)
- Susceptible to local minimum
- As it moves closer to local minimum, it automatically takes smaller steps as gradient decreases. **Time complexity**  $O(m)$

## Batch-Gradient Descent

### Algorithm 2: Batch-Gradient Descent

- 1 Initialize  $w$  randomly
- 2 **repeat**
- 3 | Simultaneously update all  $w_j$  with  $w_j - \alpha \frac{1}{m} \sum_{i=1}^m (y(x^{(i)}, w) - y^{(i)}) x_j^{(i)}$
- 4 **until** converge;
- 5 **return**  $w$

- At every step it evaluate all training examples
- Sometime it is also called multi-variate linear regression

## Example: Gradient Descent (learning rate $\alpha$ )

### Consider following data

	$x_1$	$x_2$	$x_3$	$y$
1	10	50	20	10
2	11	31	22	12
3	11	12	15	4
4	20	55	20	22
5	23	41	27	1
6	31	12	35	9
7	13	18	12	23
8	21	55	16	16
9	32	56	27	22
10	8	22	35	11

```
J=396.662506
w=(0.5000,0.5000,0.5000,0.5000)
J=19454472.000000
w=(-2.0155,-51.0770,-1.00,970.-62.640)
J=1096526813184.000000
w=(590.236,11518.771,23902.906,13778.348)
J=5230041021218816.000000
w=(-13.8991,362.-26539762.50.-5525792.009,-3170425.000)
J=2942665354556228763648.000000
w=(31365578.000,612476928.000,1275686924.000,731686912.000)
J=156806927293873988314651668.000000
w=(-72.40111104,0.000,-141378551.808,0.000,-2944657358208.000)
J=8352966854652697102769827428352.000000
w=(1671254376448.000,32834791002112.000,67972376530304.000,38386479304704.000)
J=44520007922259187970766868873062340.000000
w=(-3857802,70759936.000,-7533178626784768.000)
-15690251045437440.000,-8995357718200320.000)
```

Learning rate  $\alpha = 0.1$

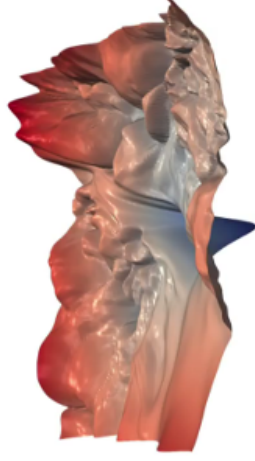
## The Partial Derivative term

$$\begin{aligned} \frac{\partial}{\partial w_j} J(w) &= \frac{\partial}{\partial w_j} \frac{1}{2m} \sum_{i=1}^m (y(x^{(i)}, w) - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m \frac{\partial}{\partial w_j} (y(x^{(i)}, w) - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m 2(y(x^{(i)}, w) - y^{(i)}) \frac{\partial}{\partial w_j} (y(x^{(i)}, w) - y^{(i)}) \\ &= \frac{1}{m} \sum_{i=1}^m (y(x^{(i)}, w) - y^{(i)}) \frac{\partial}{\partial w_j} y(x^{(i)}, w) \end{aligned}$$

For  $y(x^{(i)}, w) = w_0 + w_1 x_1^{(i)} + \dots + w_n x_n^{(i)}$  we have  $\frac{\partial}{\partial w_j} y(x^{(i)}, w) = x_j^{(i)}$

$$\frac{\partial}{\partial w_j} J(w) = \frac{1}{m} \sum_{i=1}^m (y(x^{(i)}, w) - y^{(i)}) x_j^{(i)}$$

## Real Loss Landscape



- Loss Landscape is rude <sup>1</sup>

<sup>1</sup> Visualizing the Loss Landscape of Neural Nets, Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, Tom Goldstein, 2018

## Example: Gradient Descent (learning rate $\alpha$ )

### Consider following data

	$x_1$	$x_2$	$x_3$	$y$
1	10	50	20	10
2	11	31	22	12
3	11	12	15	4
4	20	55	20	22
5	23	41	27	1
6	31	12	35	9
7	13	18	12	23
8	21	55	16	16
9	32	56	27	22
10	8	22	35	11

```
J=396.663
w=(0.500,0.500,0.500,0.500)
684.137
(0.474,-0.016,-0.515,-0.131)
1317.021
(0.508,0.631,0.881,0.828)
1943.882
(0.464,-0.249,-0.910,-0.435)
3557.655
(0.523,0.888,1.482,0.914)
5815.401
(0.448,-0.630,-1.641,-0.908)
10067.481
(0.549,1.396,2.918,1.466)
17124.684
(0.492,-1.163,-2.264,1.055)
30124.684
(0.592,2.737,4.926,2.058)
52847.000
(0.359,-2.383,-5.221,-3.028)
91828.805
(0.688,3.630,7.314,4.151)
159276.761
(0.263,-4.302,-9.200,-5.330)
277927.562
(0.789,6.152,12.560,7.155)
481973.594
(0.083,-7.633,-16.125,-9.316)
837646.250
(1.025,10.537,21.725,12.387)
1455901.375
(-0.201,-13.418,26.168,-16.234)
2330747.500
(1.418,16.923,37.811,21.487)
4000000.000
(1.418,16.923,37.811,21.487)
7642525.500
(2.697,31.415,65.218,37.319)
13262503.000
(-3.278,54.449,113.196,64.832)
23064938.000
(3.162,-71.310,148.738,-85.405)
40121436.000
(-5.329,94.483,196.678,112.653)
6970584.000
(5.329,94.483,196.678,112.653)
121190936.000
(8.884,164.069,341.494,195.769)
210628446.000
(15.089,261.382,526.385,320.035)
382929515.000
(15.089,261.382,526.385,320.035)
682929515.000
(-18.729,375.224,781.738,-448.478)
1327783808.000
(25.852,485.147,1031.086,591.231)
3340036608.000
(32.908,652.287,1358.811,-779.468)
```

Learning rate  $\alpha = 0.001$

## Example: Gradient Descent (Learning rate $\alpha$ )

### Consider following data

	$x_1$	$x_2$	$x_3$	$y$
1	10	50	20	10
2	11	31	22	12
3	11	12	15	4
4	20	55	20	22
5	23	41	27	1
6	31	12	35	9
7	13	18	12	23
8	21	55	16	16
9	32	56	27	22
10	8	22	35	11

Learning rate  $\alpha = 0.0001$

## Example: Gradient Descent (Feature scaling)

### Feature scaling

	$x_1$	$x_2$	$x_3$	$y$
1	0.08	0.86	0.35	10
2	0.12	0.43	0.43	12
3	0.12	0.00	0.13	4
4	0.50	0.98	0.35	22
5	0.62	0.66	0.65	1
6	0.96	0.00	1.00	9
7	0.21	0.14	0.00	23
8	0.54	0.98	0.17	16
9	1.00	1.00	0.65	22
10	0.00	0.23	1.00	11

Learning rate  $\alpha = 0.1$

## Example: Stochastic Gradient Descent

### Data

	$x_1$	$x_2$	$x_3$	$y$
1	10	50	20	10
2	11	31	22	12
3	11	12	15	4
4	20	55	20	22
5	23	41	27	1
6	31	12	35	9
7	13	18	12	23
8	21	55	16	16
9	32	56	27	22
10	8	22	35	11

Learning rate  $\alpha = 0.0008$

## Some Tricks

- Apply feature scaling (like min-max normalization)
- Sometime mean normalization is also good idea
- Declare convergence when  $J(w)$  decreases less then a small constant  $\epsilon$  (say  $\epsilon = 0.001$ )
- Choose  $\alpha$  as 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, 3, ...
- Having with some insight of the problem, you may define new features as  $f(x_1, x_2, \dots, x_k)$
- In case one wish to use  $y(x; w) = w_0 + w_1 x_1 + w_2 x_2^2 + w_3 x_3 + \dots$  Replace the features  $x_2$  with  $x_2^2$  and  $x_3$  with  $x_3^4$  and so on. This would enable the use of gradient descent possible once again

## Stochastic Gradient Descent

Gradient descent is computationally expensive for large training set

- Don't wait! consider single training example during iteration. <sup>2</sup>

### Algorithm 3: Stochastic-Gradient Descent

- 1 Initialize  $w$  randomly
- 2 Randomly shuffle the dataset
- 3 repeat
- 4     for  $i = 1$  to  $m$  do
- 5         Simultaneously update all  $w_j$  with  $w_j - \alpha(y(x^{(i)}, w) - y^{(i)})x_j^{(i)}$
- 6     until converge;
- 7     return  $w$

- The idea can be applied to other learning method as well

<sup>2</sup> Takes a random looking path and wander around minimum (may not converge)

## Mini-Batch Gradient Descent

- Doing somewhat in between the batch-Gradient Descent (all  $m$  training sample) and Stochastic-Gradient Descent (only one training sample)
- Let  $b$  be the size of mini batch

### Algorithm 4: Mini-Batch Gradient Descent

- 1 Initialize  $w$  randomly
- 2 Randomly shuffle the dataset
- 3 repeat
- 4     for  $i = 1$  to  $m - b + 1$  in steps of  $b$  do
- 5         Simultaneously update all  $w_j$  with  $w_j - \alpha \frac{1}{b} \sum_{j=i}^{i+b} (y(x^{(i)}, w) - y^{(i)})x_j^{(i)}$
- 6     until converge;
- 7     return  $w$

## Direct method of MSE optimization

If we have two data points  $(x_1, y_1)$  and  $(x_2, y_2)$  having two attributes  $m = 2, n = 1 + 1$

- Finding parameters  $w_0$  and  $w_1$  such that  $y = w_0 + w_1x$  is easy
- Let
$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \end{bmatrix} \quad y = [y_1, y_2] \quad w = [w_0, w_1]$$
- Write  $y = w \cdot X^T$
- That leads to  $w = y \cdot (X^T)^{-1}$

But most of the time we have lot many data points as compared to number of attributes ( $m \gg n$ ).

**In such over constrained systems, inverse  $(X^T)^{-1}$  does not exist**

## Normal Equation

- Consider  $X$  as *design matrix*, constructed using training data of size  $m \times (n + 1)$
- And  $Y$  be the target values (matrix of size  $m \times 1$ )
- Weights could be computed directly using the calculus
$$w = y \cdot X \cdot (X^T X)^{-1}$$
- Computation of matrix inverse could need **pseudo-inverse** computation if  $X^T X$  is non invertible (singular)
- This happens if **features are linear combination of others**, or the **matrix is not square**. Delete some features or apply regularization.
- No need to chose  $\alpha$
- Matrix computation is hard takes roughly  $O(m^3)$  time. One can avoid normal equation if  $m > 1000$

Thank You!

Thank you very much for your attention! (Reference<sup>3</sup>)

Queries ?

## Direct method of MSE optimization

Gradient is **zero** at the point where loss is minimum.

Loss function being  $J_w = (y - w \cdot X^T)^2$  gradient is given by  $\nabla J_w = -(y - w \cdot X^T)X$ . Let us find where it is **zero**

$$\begin{aligned} -(y - w \cdot X^T)X &= 0 \\ y \cdot X - w \cdot X^T \cdot X &= 0 \\ w \cdot X^T \cdot X &= y \cdot X \\ w &= y \cdot X \cdot (X^T \cdot X)^{-1} \end{aligned}$$

- here  $X \cdot (X^T \cdot X)^{-1}$  is called as **pseudo-inverse**

With large data size  $m$ , it takes  $O(m^3)$  time to compute. However, gradient descent would be much faster, taking only  $O(m)$  time.

## MSE

$$J_\theta = \sum (y - \hat{y})^2$$

- Use of MSE is computationally convenient.
- It measures the variance of residue
- Corresponds to the likelihood under Gaussian model of noise
- It is very sensitive to outliers

### Issues with gradient descent

- Very general
- Gets stuck at local minima
- Difficult to find appropriate step size (learning rate  $\alpha$ )

<sup>3</sup> [1] Book - *Pattern Recognition And Machine Learning*, Bishop, Springer-2006 (CH-3), [2] Book - *Machine Learning*, ch6, Tom M. Mitchell.