# Linear Regression-2



THE BEST VIEW COMES AFTER THE HARDEST CLIMB

# Agenda

① Model interpretability

② feature importance

③ Mathematics
$\quad\quad$ ↳ gradient Descent

# Model Interpretability

Model
$$\hat{y} = w_1 x_1 + w_2 x_2 + \ldots + w_d x_d + w_0$$

$$\hat{y} = \ldots (-1000) \cdot age + (-10) \, odo. + (200) engine + \ldots$$

Case I:  $\quad w_j^- \rightarrow -ve$

$$x_j \uparrow \rightarrow \hat{y} \downarrow$$

$$odo \uparrow \rightarrow Price \downarrow$$

$$age + 1 \rightarrow \hat{y} : \hat{y} - 10000$$

Case II $\quad w_j \to$ +ve

Case III $\quad w_j \to 0$

No impact on $\hat{y}$ due to $f_j$

## feature importances

"weights"

magnitude of $w_j \uparrow \; = \;$ importance$(f_j) \uparrow$

$$\hat{y} = (0.8) \cdot x_2 + \cdots + (0.42) x_5$$

$x_2$ has higher impact

$$\hat{y} = (-1.9) x_2 + \cdots + (-4.4) \cdot x_5$$

$x_5$ has higher impact

$$\hat{y} = (2.3) x_2 + \cdots + (-5.1) \cdot x_5$$

$x_5$ has higher impact

**Feature importance in linear regression is determined by :**

85 users have participated

| | | | |
|---|---|---|---|
| ✓ | A | The magnitude of the regression coefficients. | 92% |
| | B | The number of observations in the dataset. | 0% |
| | C | The correlation between the independent variables. | 2% |
| | D | The average squared difference between the predicted and actual values. | 6% |

End Quiz Now

*feature Scaling*

age

Odometer

$$\hat{y} = \ - \ - \ . \quad -10000 \ age \ + \quad -10. \ odo$$

[1-15]

[10k - 150k]

age is Better  odo ✗

**When assessing model interpretability in Linear Regression, what is the impact of feature scaling?**

86 users have participated

| A | Feature scaling does not affect model interpretability | 0% |
|---|---|---|
| ✗ | Feature scaling improves model interpretability | 35% |
| ✓ C | Feature scaling can help compare the magnitudes of different coefficients | 65% |

End Quiz Now

**Consider the following Linear Regression model equation: y = 5.2x1 - 3.8x2 + 2.1x3 + 0.01x4 - 1.5 if we were to drop one feature, which one would be the best choice ?**

37 users have participated

| A | x1 | 0% |
|---|---|---|
| B | x2 | 5% |
| C | x3 | 3% |
| ✓ D | x4 | 92% |

End Quiz Now

$$y = 5.2\,x_1 - 3.8\,x_2 + 2.1\,x_3 + 0.01\,x_4 - 1.5$$

# Gradient Descent

$$L = y = f(x) = (x-5)^2$$
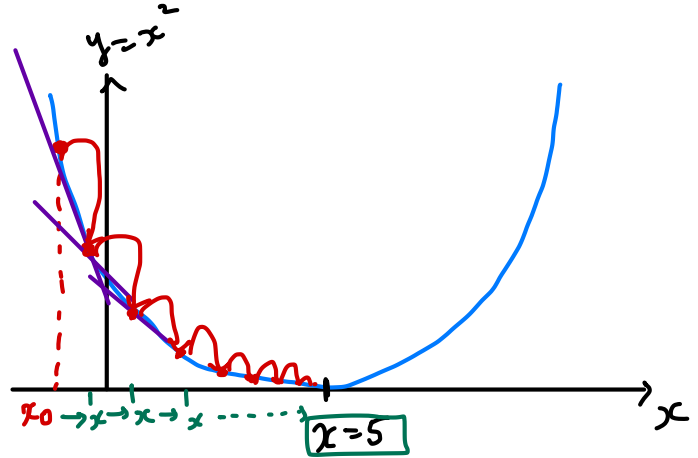

$$y = x^2$$
$$x = 5$$

G.D
=

1. Pick randomly $x_0$

2. $\frac{\partial L}{\partial x}\Big|_{x_0}$ $-ve$    move towards neg derivatives

3. $x = x - \eta \cdot \frac{\partial L}{\partial x}$

   $\hookrightarrow$ learning rate : 0.1

# In gradient descent, what does the gradient represent ?

90 users have participated

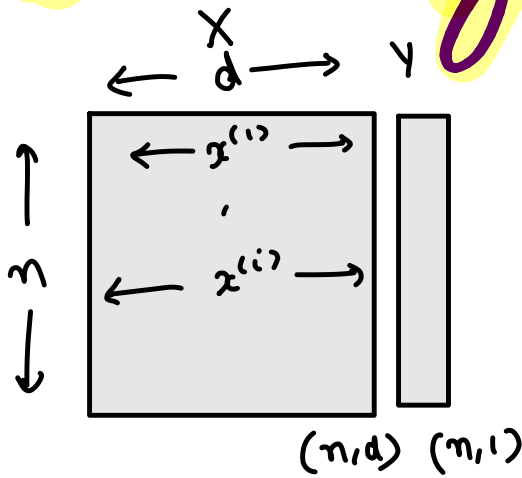| | | | |
|---|---|---|---|
| ✓ | A | The direction of steepest increase of the cost function | 34% |
| | B | The direction of steepest decrease of the cost function | 62% (-gradient) |
| | C | The number of training examples in the dataset | 1% |
| | D | The number of layers in the neural network | 2% |

**End Quiz Now**

# Linear Regression

$$\hat{y}^{(i)} = w^T \cdot x^{(i)} + w_0$$

$$\hat{y} = X \cdot W + w_0$$

$(n,d) \quad (d,1)$

$(n,1)$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \quad (d,1)$$

$X \quad (n,d) \qquad Y \quad (n,1)$

$x^{(1)}$, $x^{(i)}$

$$\hat{Y} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(n)} \end{bmatrix} \quad (n,1)$$

## Loss

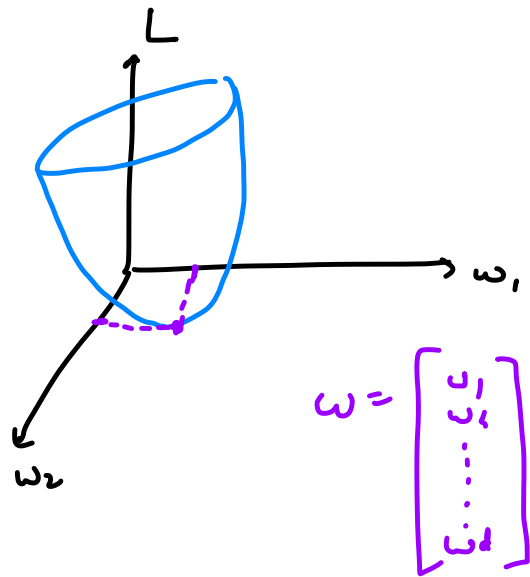$$MSE = \frac{1}{n} \sum_{i=1}^{n} \left( y^{(i)} - \hat{y}^{(i)} \right)^2$$

# Gradient Descent

$$\text{Loss}(w) = \frac{1}{n} \sum_{i=1}^{n} \left( y^{(i)} - \hat{y}^{(i)} \right)^2$$

minimise ( MSE )

$$\underset{w}{\text{argmin}} \quad \frac{1}{n} \sum_{i=1}^{n} \left( y^{(i)} - \hat{y}^{(i)} \right)^2$$



$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

1. randomly 'init' 'w'

2. $\dfrac{\partial L}{\partial w} = \nabla_w L$

3. repeat n-ith times $\left\{ \quad w_j' = w_j - \eta \cdot \dfrac{\partial L}{\partial w_j} \right\}$

# Gradients

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

$$L = \frac{1}{n} \sum_{i=1}^{n} \left( y^{(i)} - \boxed{\hat{y}^{(i)}} \right)^2$$

$$\rightarrow = w_1 x_1^{(i)} + w_2 x_2^{(i)} + \cdots$$
$$\cdots + w_d x_d^{(i)}$$

$$\frac{\partial L}{\partial w} = \begin{bmatrix} \dfrac{\partial L}{\partial w_1} \\[2mm] \dfrac{\partial L}{\partial w_2} \\[2mm] \vdots \\[2mm] \dfrac{\partial L}{\partial w_j} \\[2mm] \vdots \\[2mm] \dfrac{\partial L}{\partial w_d} \end{bmatrix}$$

away

$$L = \frac{1}{n} \sum_{i=1}^{n} \left[ y^{(i)} - \hat{y}^{(i)} \right]^2$$

$$\frac{\partial L}{\partial w_j}$$

for 1 datapoint

$$\frac{\partial L}{\partial w_j} = \frac{\partial (y - \hat{y})^2}{\partial w_j}$$

$$\frac{d}{dx} f(g(x)) = \frac{d\,f(x)}{d\,g(x)} \cdot \frac{d\,g(x)}{dx}$$

$$= \frac{\partial (y-\hat{y})^2}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_j}$$

$$= -2(y-\hat{y}) \cdot \frac{\partial}{\partial w_j} \; w_1 x_1 + w_2 x_2 + \ldots w_j x_j + \ldots w_d x_d$$

$$= -2(y-\hat{y}) \cdot x_j \;\Rightarrow\; 2(\hat{y}-y)\cdot x_j$$

derivative for all Points

$$\frac{\partial L}{\partial w_j} = \frac{2}{n} \sum_{i=1}^{n} (\hat{y}^{(i)} - y^{(i)}) \cdot x_j^{(i)}$$

# G.D

1. Randomly init $'w'$

2. $\dfrac{\partial L}{\partial w_j}$

3. Repeat $\{$

$$w_j = w_j - \eta \cdot \left[ \frac{1}{n} \sum_{i=1}^{n} \left( \hat{y}^{(i)} - y^{(i)} \right) \cdot x_j^{(i)} \right]$$

$\}$

$\dfrac{\partial L}{\partial w} \rightarrow \begin{bmatrix} \ \\ \ \end{bmatrix}_{(d,1)}$

$$\frac{\partial L}{\partial w} = \frac{2}{n} \sum_{i=1}^{n} \left( \underbrace{\hat{y}^{(i)} - y^{(i)}}_{A} \right) \cdot \underbrace{x_j^{(i)}}_{B}$$

$$\underbrace{\left( \hat{Y} - Y \right)}_{(n,1)} \qquad X \longrightarrow (n,d)$$

$$\begin{bmatrix} \\ \\ \\ \\ \end{bmatrix}$$
$(d,1)$

$$\underset{\downarrow}{X^T} \cdot \underset{\downarrow}{\left( \hat{Y} - Y \right)}$$

$$(d,n) \cdot (n,1)$$

$$(d,1)$$

## What is the objective of Gradient Descent in linear regression?

64 users have participated

| | | |
|---|---|---|
| A | Minimize the absolute error | 14% |
| ✓ B | Minimize the squared error | 73% |
| C | Maximize the R-squared score | 5% |
| D | Maximize the accuracy | 8% |

**End Quiz Now**

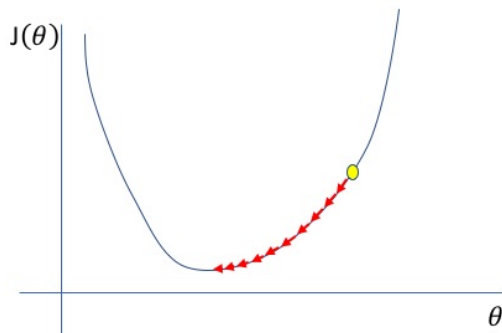## What happens if the learning rate in gradient descent for linear regression is set too large?

21 users have participated

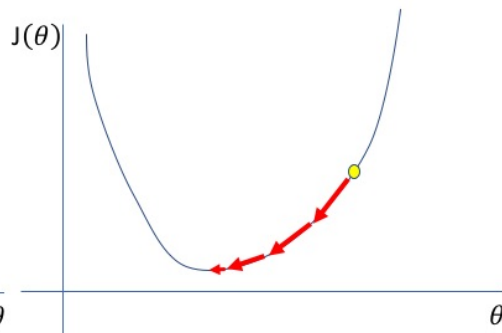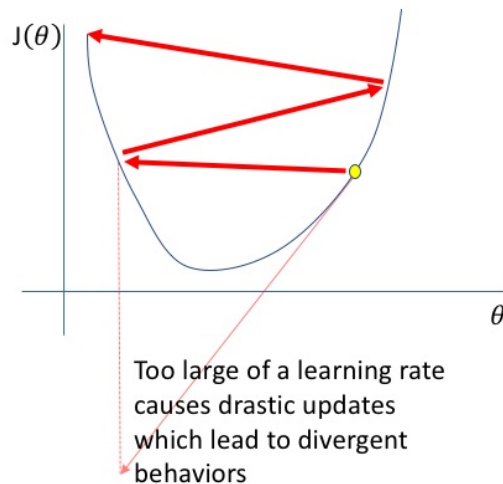| | | |
|---|---|---|
| A | The algorithm will converge faster to the optimal solution. | 10% |
| B | The model will overfit the training data, leading to poor generalization. | 10% |
| ✓ C | The algorithm may fail to converge, and the coefficients may oscillate or diverge. | 81% |
| D | The cost function will be overestimated, resulting in an inflated R2 score. | 0% |

**End Quiz Now**

**Too low**

$J(\theta)$

$\theta$

A small learning rate requires many updates before reaching the minimum point

**Just right**

$J(\theta)$

$\theta$

The optimal learning rate swiftly reaches the minimum point

**Too high**

$J(\theta)$

$\theta$

Too large of a learning rate causes drastic updates which lead to divergent behaviors

Vanilla/ Batch G.D

Mini-Batch G.D

Stochastic G.D

°C    °F

$$\left[ \begin{array}{c} \text{و}\text{ دا wo} \\ \text{imp} \end{array} \right]$$

$w_0$