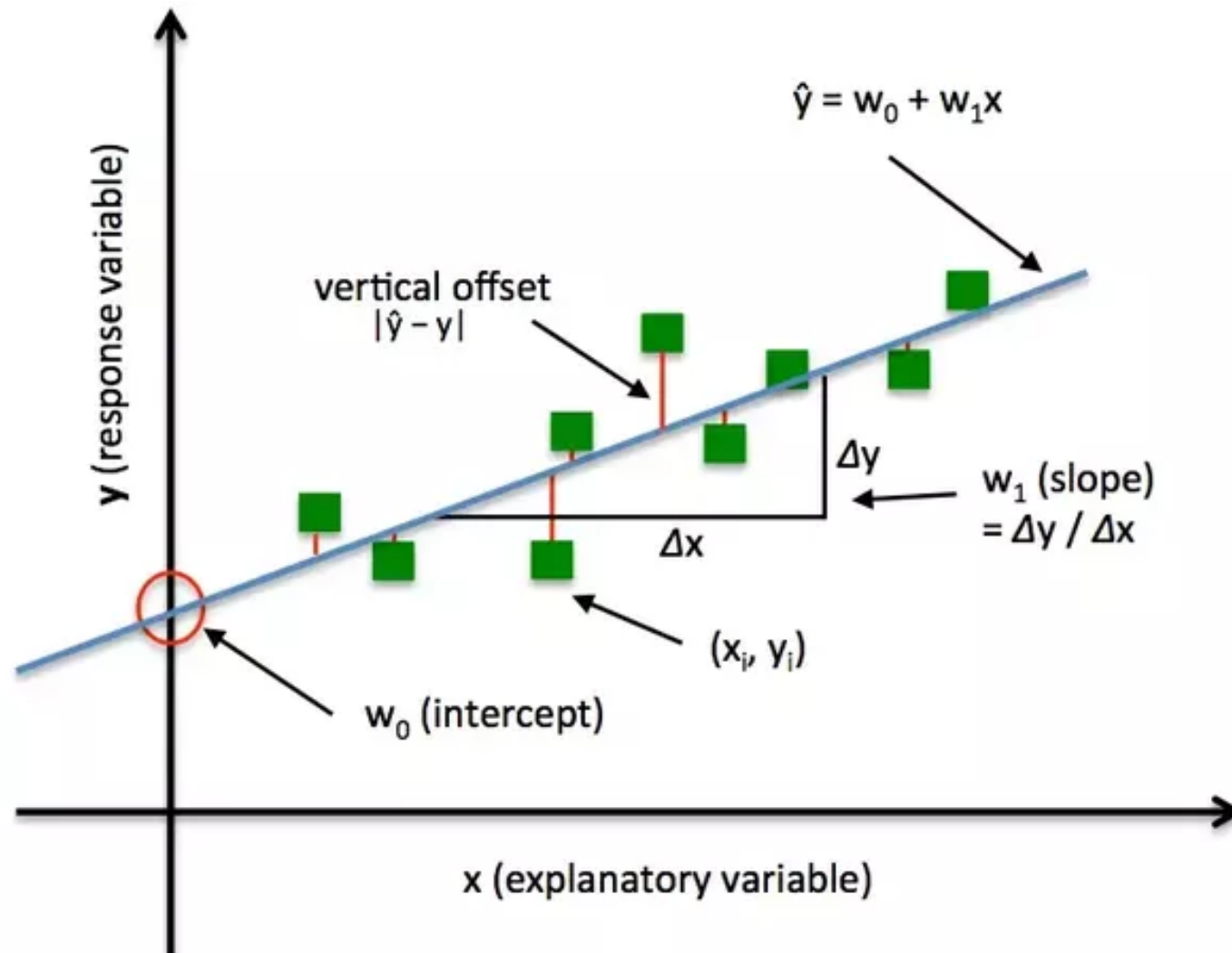# Linear regression

# Assumptions of linear regression

Linear regression has the following **assumptions**, failing which the linear regression model does not hold true:

1. The dependent variable should be a linear combination of independent variables
2. No autocorrelation in error terms
3. Errors should have zero mean and be normally distributed
4. **No or little multi-collinearity**
5. Error terms should be homoscedastic

# Optimization in Machine Learning

(A) Deterministic Problems

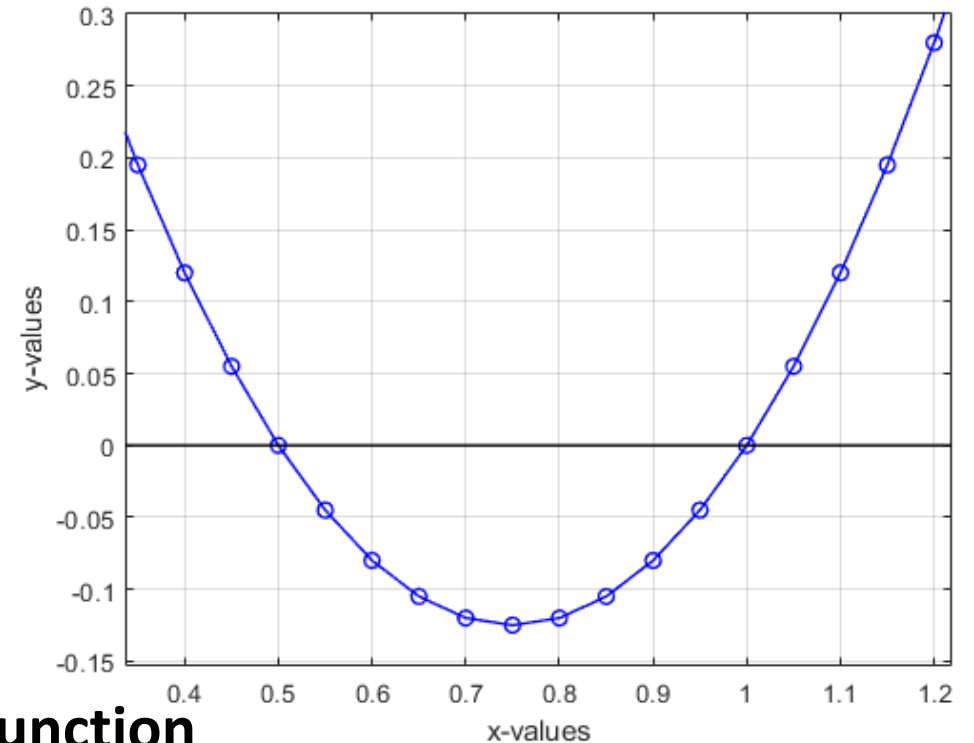Q. Find the roots of the $f(x)$, where $f(x) = 2x^2 - 3x + 1$

$$Solution: x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = 0.5 \ and \ 1.0$$

(B) Convex Optimization Problems

Q. $\min f(x)$, where $f(x) = 2x^2 - 3x + 1$

Q. $\underset{x}{\operatorname{argmin}} \ f(x)$, where $f(x) = 2x^2 - 3x + 1$

f(x) is called as objective function, **cost function / Loss function**

# Optimization in Machine Learning

**(A) Deterministic Problems:** Q. Find the roots of the $f(x)$, where $f(x) = 2x^2 - 3x + 1$

$$Solution: x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = 0.5 \; and \; 1.0$$
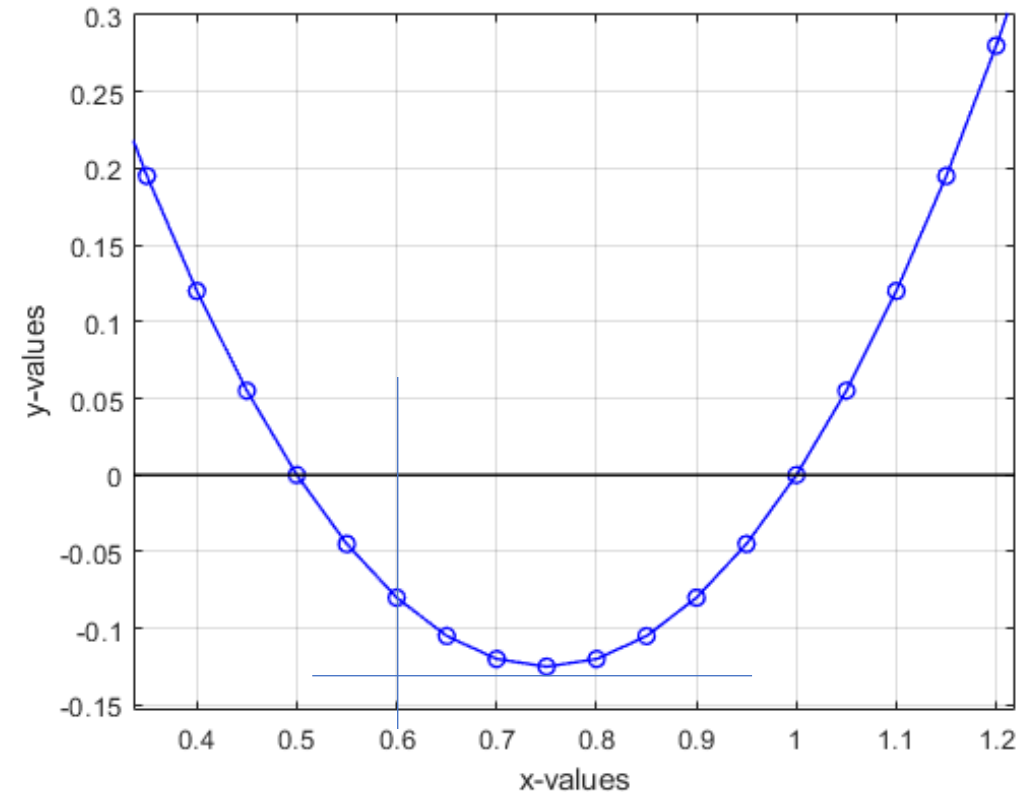
**(B) Optimization Problems:**

minimization/maximization of an objective function, subject to some constraints.

$\min f(x)$, where $f(x) = 2x^2 - 3x + 1$

$\underset{x}{\operatorname{argmin}} \; f(x)$, where $f(x) = 2x^2 - 3x + 1$

$\underset{x}{\operatorname{argmin}} \; f(x)$, where $f(x) = 20x^2 - 30x + 10$

Both lead to the same "optimal values" of x

# Optimization in Machine Learning

Types of Optimization Problems, in general:

(A) Unconstrained Optimization

$$\min f(x), \text{ where } f(x) = 2x^2 - 3x + 1$$

$$\operatorname*{argmin}_{x} f(x), \text{ where } f(x) = 2x^2 - 3x + 1$$

(B) Constrained Optimization:
   a.   Linear or non-linear constraints

   b. Equality or inequality constraints

$$\min (x_1 - 2)^2 + (x_2 - 1)^2 \quad \text{subject to} \begin{cases} x_1^2 - x_2 & \leq 0, \\ x_1 + x_2 & \leq 2. \end{cases}$$
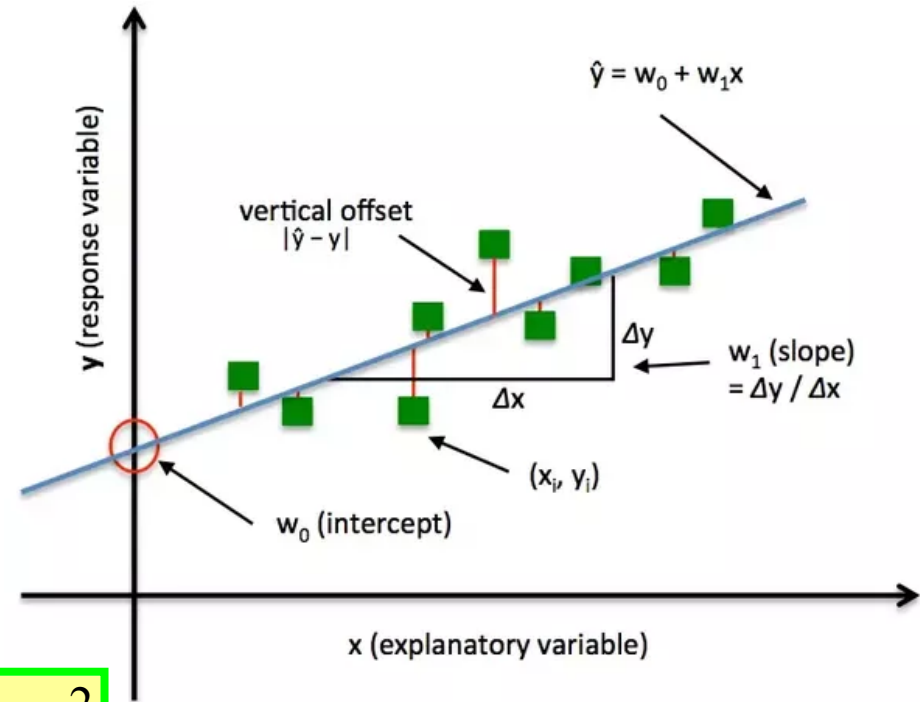
# Possible Loss (Cost) functions in ML (Regression)

(1) Sum of errors (SE): $L = \sum_{i=1}^{N} \left( \hat{Y}_i - Y_i \right)$

(2) Sum of Absolute Errors (SAE): $L = \sum_{i=1}^{N} \left| \hat{Y}_i - Y_i \right|$
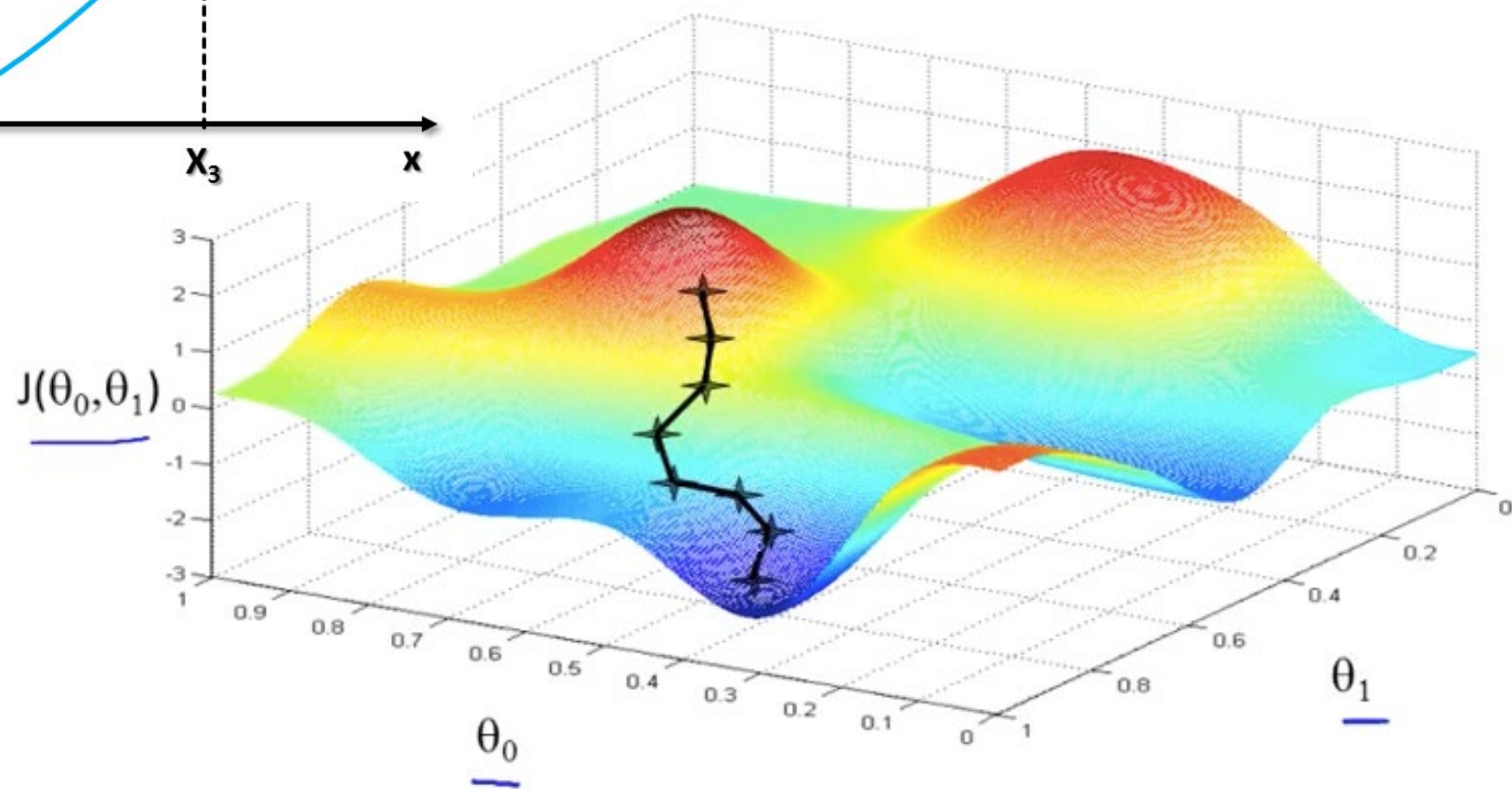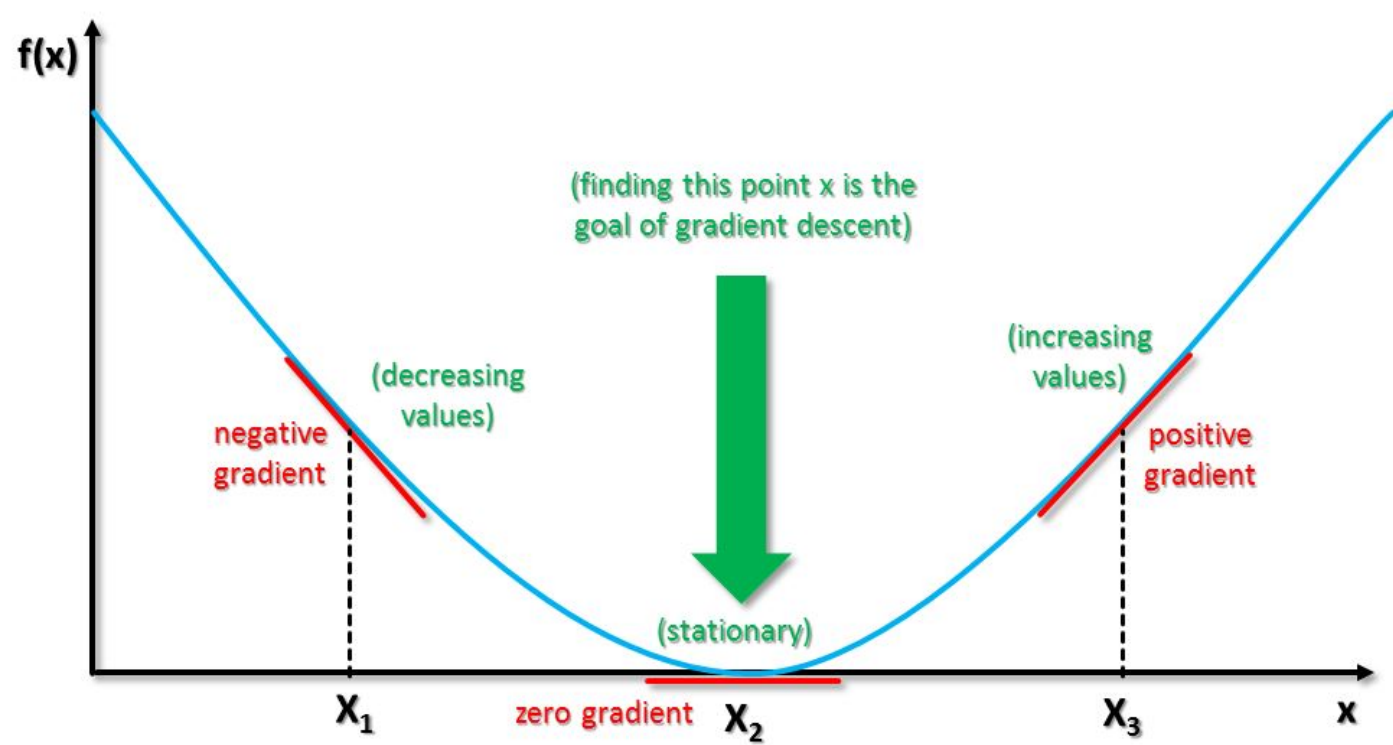
(3) Sum of Squares of Errors (SSE): $L = \sum_{i=1}^{N} \left( \hat{Y}_i - Y_i \right)^2$

(4) Mean of Squares of Errors (MSE): $L = \dfrac{1}{N} \sum_{i=1}^{N} \left( \hat{Y}_i - Y_i \right)^2$

(5) Root Mean of Squares of Errors (RMSE): $L = \sqrt{\dfrac{1}{N} \sum_{i=1}^{N} \left( \hat{Y}_i - Y_i \right)^2}$



| X1 | Y |
|---|---|
| 1 | 4.8 |
| 3 | 11.4 |
| 5 | 17.5 |
| .. | ... |

# Linear Regression Problem Formulation:

$$Objective\ Fn : MSE : \ L = \frac{1}{N}\sum_{i=1}^{N}(\hat{Y}_i - Y_i)^2 = \frac{1}{N}\sum_{i=1}^{N}(error)^2$$
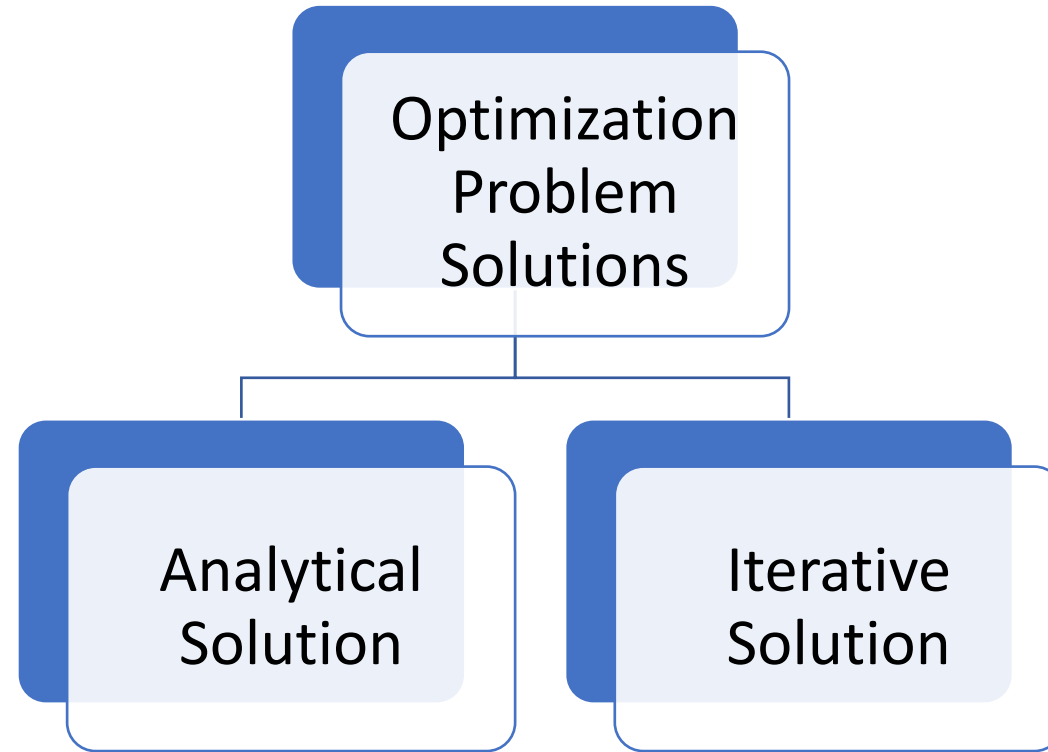
$$\arg\min_{w_j} L(w_j \mid X, Y)$$

$$where,\ \ \hat{Y} = w_0 + w_1 X_1$$

Says, Find the **optimal weights ($w_j$)** for which the **MSE**

Loss function has **min value**, for a **GIVEN X,Y data.**

| X1 | Y |
|----|------|
| 1 | 4.8 |
| 3 | 11.4 |
| 5 | 17.5 |
| .. | ... |

# Optimization Problem Solution Methods:

Optimization Problem Solutions

Analytical Solution

Iterative Solution

**- Theoretical Solution**, which gives the "exact" solution to the problem, provided the optim. problem has a "closed-form solution"

**-Approximate Solution** to the optim. problem, based on some iterative algorithm
- Can solve all types of optim. problems.

# Analytical Solution for Unconstrained Optimization:

**FOC: Necessary Condition:** States that the first (odd) derivative (gradient) of the objective function must vanish at the optimal points (points of maxima/minima)

$$e.g. \ f(x) = 2x^2 - 3x + 1 \qquad \frac{df}{dx} = 4x - 3 = 0 \qquad \Rightarrow x^* = 0.75$$

**SOC: Sufficiency Condition:** States that the second (even) derivative (gradient) of the objective function **evaluated at the optimal points**, must be:

- **Positive >> for minima**

- Negative >> for maxima

$$\left. \left\| \frac{d^2 f}{dx^2} \right\| \right|_{x^*=0.75} = 4 \ (positive)$$

which means minima occurs at $x^* = 0.75$

c) worst-case scenario: d2f/dx2 evaluated at some optimal x is ZERO...,
THEN, that optimal point is neither a point of maxima/minima. It could just a point of inflexion or a saddle point.

if you want to maximize a fcn:
1) minimize the negative of the obj function
2) minimize the inverse of the obj fcn. >> used on SVM

# Linear Regression – Analytical Solution:

$$MSE : L = \frac{1}{N}\sum\left(\hat{Y} - Y\right)^2 \quad where, \quad \hat{Y} = w_0 + w_1 X_1$$

| X1 | Y |
|----|------|
| 1 | 4.8 |
| 3 | 11.4 |
| 5 | 17.5 |

$$\arg\min_{w_j} L(w_j \mid X, Y)$$

Find the **optimal weights (wj)** for which the MSE Loss function has min value, for **GIVEN X,Y data**.

## STEP – 1: Get the Gradients:

**FINAL GRADIENTS**

$$\frac{\partial L}{\partial w_j} = \frac{\partial L}{\partial \hat{Y}} \times \frac{\partial \hat{Y}}{\partial w_j}$$

$$\frac{\partial L}{\partial w_j} = \frac{2}{N}\sum\left(\hat{Y} - Y\right) \times \frac{\partial \hat{Y}}{\partial w_j}$$

$$\frac{\partial L}{\partial w_0} = \frac{2}{N}\sum\left(\hat{Y} - Y\right) \times 1$$

$$\frac{\partial L}{\partial w_1} = \frac{2}{N}\sum\left(\hat{Y} - Y\right) \times X_1$$

# Linear Regression – Analytical Solution:

## STEP – 2: Equate the Gradients to zero and solve:

$$\frac{2}{N} \sum (\hat{Y} - Y) \times 1 = 0 \qquad -- (1)$$

$$\frac{2}{N} \sum (\hat{Y} - Y) \times X_1 = 0 \qquad -- (2)$$

**FINAL SOLUTION** ⬅ **Ordinary Least Squares (OLS) !**

$$w_1^* = \frac{\left[ \overline{XY} - \overline{X}\,\overline{Y} \right]}{\left[ \left( \overline{X^2} \right) - \left( \overline{X} \right)^2 \right]}$$

$$\hat{Y} = w_0 + w_1 X_1$$

$$\overline{Y} = w_0 + w_1 \overline{X_1}$$

Intercept / Bias Term:

$$w_0^* = \overline{Y} - \left( w_1^* \cdot \overline{X} \right)$$

# Linear Regression – Analytical Solution:

## STEP – 3: Calculate the optimal weights

$$w_1^* = \frac{\left[\overline{XY} - \overline{X}\,\overline{Y}\right]}{\left[\left(\overline{X^2}\right) - \left(\overline{X}\right)^2\right]}$$

$$w_0^* = \overline{Y} - \left(w_1^* \cdot \overline{X}\right)$$

| X | Y |
|---|---|
| 1 | 4.8 |
| 3 | 11.4 |
| 5 | 17.5 |

Substitute:

$$w_1^* = \frac{\left[547.87 - (3 \times 11.67)\right]}{\left[(11.67) - (3)^2\right]}$$

| | X | Y | X² | XY |
|---|---|---|---|---|
| | 1 | 4.8 | 1 | 1x4.8 |
| | 3 | 11.4 | 9 | 3x11.4 |
| | 5 | 17.5 | 25 | 5x17.5 |
| **Mean** | 9/3 = 3 | 35/3 =11.67 | 35/3 = 11.67 | (1643.6)/3 = 547.87 |

# Linear Regression – Analytical Solution:

$$w_1^* = \frac{\left[\overline{XY} - \overline{X}\,\overline{Y}\right]}{\left[\left(\overline{X^2}\right) - \left(\overline{X}\right)^2\right]}$$

**is also equivalent to:**

$$w_1^* = \frac{\sum_i \left(X_i - \overline{X}\right)\left(Y_i - \overline{Y}\right)}{\sum_i \left(X_i - \overline{X}\right)^2} = \frac{\text{cov}(X, Y)}{\text{var}(X)}$$

## OLS in Vector Form:

$$L = \frac{1}{N}\sum\left(\hat{Y}_i - Y_i\right)^2 = \frac{1}{N}\left[\left(\hat{Y} - Y\right)^T\left(\hat{Y} - Y\right)\right]$$ **Loss Function in Vector Notation**

$$W = (X^T X)^{-1} X^T Y$$ ⬅ **OLS Solution (Analytical Solution) !**

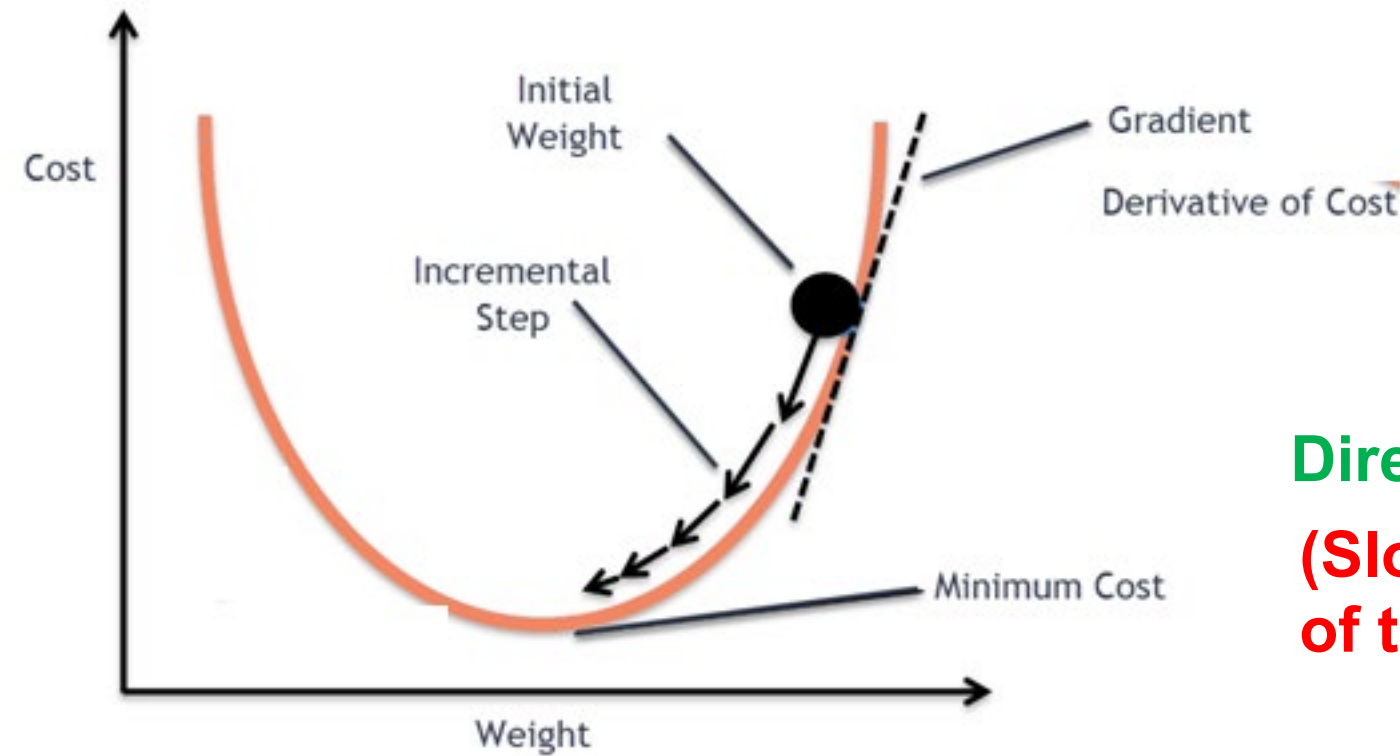$$\frac{\partial^2 L}{\partial W^2} = \frac{1}{N}\left[(X^T X)(2)\right] = +ve \ \ (\text{minima})$$

# Optimization Problem Solution Methods:

```
           ┌─────────────────┐
           │  Optimization   │
           │    Problem      │
           │   Solutions     │
           └────────┬────────┘
             ┌──────┴──────┐
        ┌────┴────┐   ┌─────┴──────┐
        │   OLS   │   │  Gradient  │
        │         │   │  Descent   │
        └─────────┘   └────────────┘
```

**- Analytical Solution**, which gives the "exact" solution to the problem, provided the optim. problem has a "closed-form solution

**-Approximate Solution (iterative)** to the optim. problem, based on some iterative algorithm
- Can solve all types of optim. problems.

# Gradient Descent Algorithm:



**Gradient Descent Update Rule:**

$$w_j^{k+1} = w_j^k - \Delta w_j$$

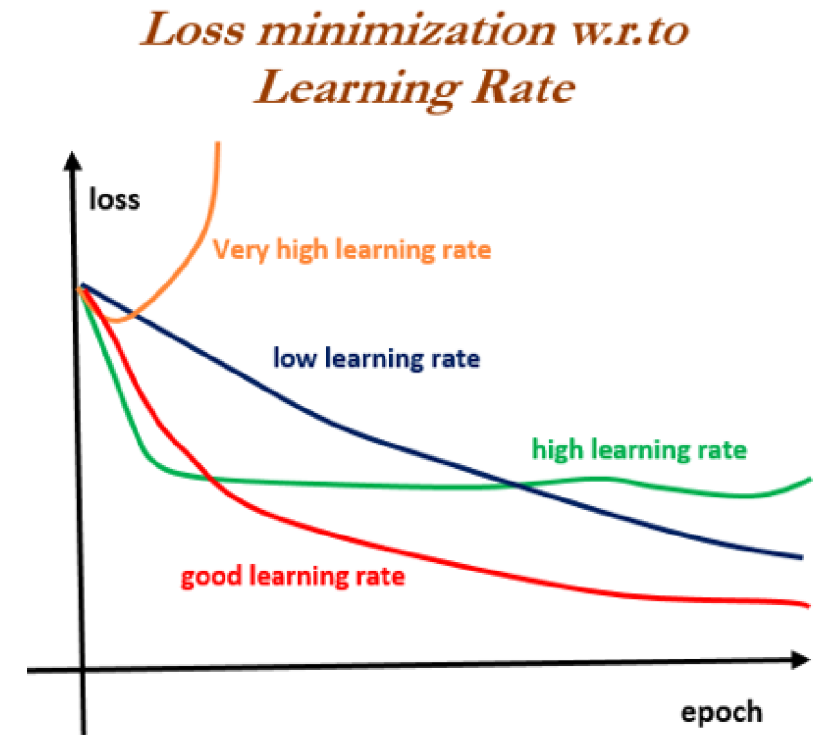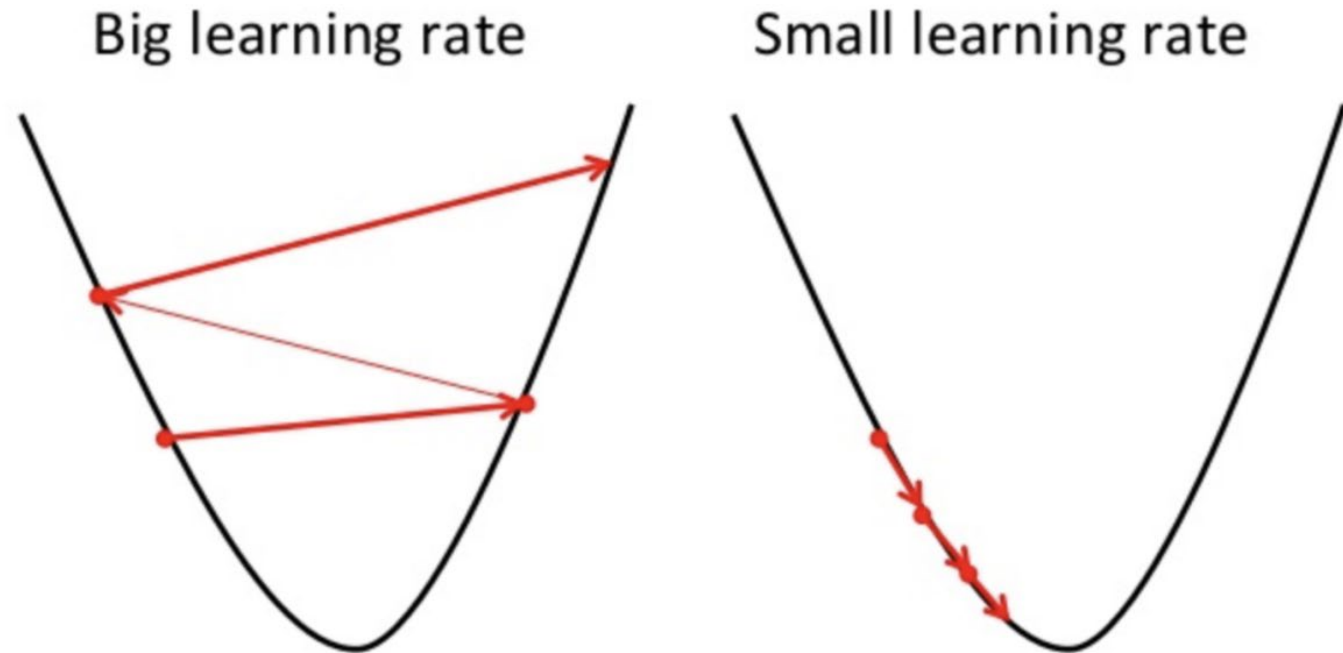**Direction of Update**

**(Slope / Gradient of the Loss Fn)**

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

$$\frac{\partial L}{\partial w_0} = \sum (\hat{Y} - Y) \times 1$$

$$\frac{\partial L}{\partial w_1} = \sum (\hat{Y} - Y) \times X_1$$

**Amount of Update**

**Step Size:** $\alpha$
**or Learning Rate**

**Final Gradient Descent Update Rule:**

$$w_j^{k+1} = w_j^k - \left( \alpha \frac{\partial L}{\partial w_j} \right)$$

# Effect of Learning Rate:



Big learning rate

Small learning rate

Loss minimization w.r.to Learning Rate

**Alpha** is a **Hyper-parameter** that <u>YOU</u> have to decide (based on your exp. & domain knowledge)... trail & error. HP are to be specified/decided before the start of the iterations start. [0.1, 0.02, 0.04, 0.05, 0.06, 0.08, 0.01]

**Model coefficients/weights (W)** >> **Model Parameters >> these are "learnt" by the optimizer algo from the DATA.** You don't specify this.

## Regularization in Machine Learning

$$MSE: \ L = \frac{1}{N}\sum_{i=1}^{N}(\hat{Y}_i - Y_i)^2 \qquad\qquad p-Norm \ (L_p) = \left\| w_j \right\|_p = \left( \sum \left| w_j \right|^p \right)^{1/p}$$

$$Ridge: \ L = \left\{ MSE \right\} + \lambda \left\| w_j \right\|_2^2 = \left\{ \frac{1}{N}\sum_{i=1}^{N}(\hat{Y}_i - Y_i)^2 \right\} + \lambda(w_1^2 + w_2^2 + ... + w_p^2)$$

$$LASSO: \ L = \left\{ MSE \right\} + \lambda \left\| w_j \right\|_1 = \left\{ \frac{1}{N}\sum_{i=1}^{N}(\hat{Y}_i - Y_i)^2 \right\} + \lambda \left( \left| w_1 \right| + \left| w_2 \right| + ... + \left| w_p \right| \right)$$

$$ElasticNet: \ L = \left\{ MSE \right\} + \lambda_1 \left\| w_j \right\|_1 + \lambda_2 \left\| w_j \right\|_2^2$$

where $j = 1, 2 ....... p$ number of features

# L1 Regularization (also called as LASSO penalisation)

Involves penalising sum of absolute values (1-norms) of regression coefficients

$$LASSO: \ L = \{MSE\} + \lambda \|w_j\|_1 = \left\{ \frac{1}{N} \sum_{i=1}^{N} (\hat{Y}_i - Y_i)^2 \right\} + \lambda \left( |w_1| + |w_2| + ... + |w_p| \right)$$

- Here we are familiar with the First half of the Cost Function.
- By adding all weights to the cost function, which we want to minimize, we're adding further restrictions on these parameters
- **Typically intercepts are not penalised.**
- The lambda parameter in Lasso tunes the strength of the penalty, and should be determined via cross-validation.

# L2 Regularization (also called as Ridge Penalisation)

This proceeds by penalising the sum of squares (2-norms) of the model coefficients

$$Ridge: \ L = \left\{ \frac{1}{N} \sum_{i=1}^{N} (\hat{Y}_i - Y_i)^2 \right\} + \lambda \left\| w_j \right\|_2^2 = \left\{ \frac{1}{N} \sum_{i=1}^{N} (\hat{Y}_i - Y_i)^2 \right\} + \lambda (w_1^2 + w_2^2 ... + w_j^2)$$

- The L2 regularization will force the parameters to be relatively small, the bigger the penalization, the smaller (and the more robust to overfitting) the coefficients are.

- Here we are considering every feature, but we are penalizing the coefficients based on how significant the feature is.