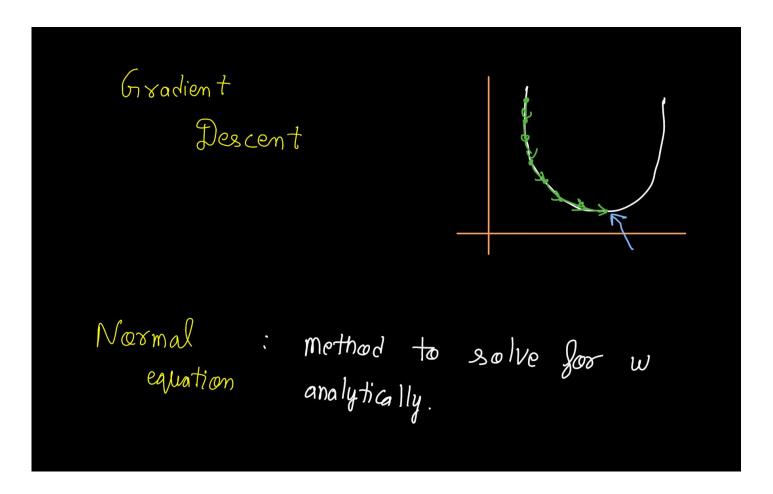## Outline

- Normal Equation/ Closed Form solution

## ⌄ Closed form solution:

Derivation: https://www.youtube.com/watch?v=g8qF61P741w

reference video: https://www.youtube.com/watch?v=B-Ks01zR4HY&ab_channel=ArtificialIntelligence-AllinOne

blog: https://towardsdatascience.com/normal-equation-in-python-the-closed-form-solution-for-linear-regression-13df33f9ad71
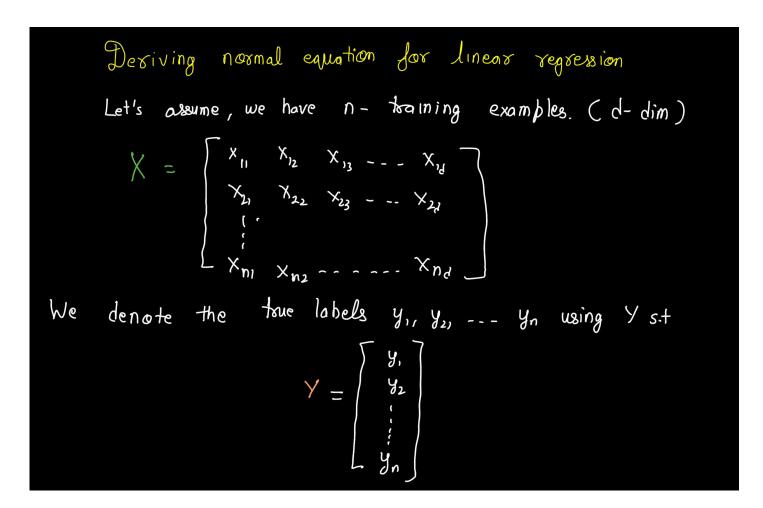


Gradient Descent is an iterative algorithm meaning that you need to take multiple steps to get to the Global optimum

## ⌄ Is there any other way to solve Linear Regression?

Yes. Using Normal Equation

It solves for the optimal values of the parameter $w$ in one step without needing to use an iterative algorithm and this algorithm is called the Normal Equation. It works only for Linear Regression and not any other algorithm.

- Normal Equation is the Closed-form solution
- It means that we can obtain the optimal parameters by just using a formula that includes a few matrix multiplications and inversions.

Objective: We want to find weights $w_1, w_2, \ldots, w_d$

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

For simplicity, assume $w_0 = 0$ ( passes through origin )

Error $(e_i) = y_i - \hat{y}_i$

In matrix form,

$$\text{Error} = \begin{bmatrix} y_1 - (w_1 x_{11} + w_2 x_{12} + \cdots + w_d x_{1d}) \\ y_2 - (w_1 x_{21} + w_2 x_{22} + \cdots + w_d x_{2d}) \\ \vdots \\ \vdots \\ y_n - (w_1 x_{n1} + w_2 x_{n2} + \cdots + w_d x_{nd}) \end{bmatrix} \begin{matrix} \rightarrow e_1 \\ \rightarrow e_2 \\ \\ \vdots \\ \\ e_n \end{matrix}$$

$$= \quad Y - Xw$$

$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} & - - - & x_{1d} \\ x_{21} & x_{22} & x_{23} & - - - & x_{2d} \\ \vdots & & & & \\ x_{n1} & x_{n2} & - - - - - - & x_{nd} \end{bmatrix} \qquad W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

$$XW = \begin{bmatrix} x_{11}w_1 + x_{12}w_2 + x_{13}w_3 + - - + x_{1d}w_d \\ x_{21}w_1 + x_{22}w_2 + - - - - - + x_{2d}w_d \\ \vdots \\ x_{n1}w_1 + x_{n2}w_2 + - - - + x_{nd}w_d \end{bmatrix}$$

Now,

$$\text{Error} = Y - XW$$

We know, Loss is sum of squared errors

$$\boxed{\mathcal{L} = \sum_{i=1}^{n} e_i^2}$$

$$= e_1^2 + e_2^2 + - - - + e_n^2$$

Now,

$$\mathcal{L} = (Y - Xw)^T (Y - Xw)$$

$$\because [e_1 \ e_2 \ e_3 \ -- \ e_n] \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \vdots \\ e_n \end{bmatrix}$$

$$= e_1^2 + e_2^2 + \cdots e_n^2$$

Multiplying the terms,

$$= Y^T Y - Y^T Xw - w^T X^T Y + w^T X^T Xw$$

$$(A+B)^T = A^T + B^T$$
$$(AB)^T = B^T A^T$$

We want to minimize loss $\mathcal{L}$ in order to minimize loss, we take its derivative & equate to $0$.

$$\nabla_w \mathcal{L} = \begin{pmatrix} \dfrac{\partial \mathcal{L}}{\partial w_1} \\ \dfrac{\partial \mathcal{L}}{\partial w_2} \\ \vdots \\ \dfrac{\partial \mathcal{L}}{\partial w_d} \end{pmatrix} = 0 \implies \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix}$$

We know,

$$\nabla_w (w^T a) = \nabla_w (a^T w) = a$$

where $\quad w^T a = a^T w = w_1 a_1 + w_2 a_2 + \cdots w_d a_d$

$$\begin{pmatrix} \dfrac{\partial (w^T a)}{\partial w_1} \\ \vdots \\ \dfrac{\partial (w^T a)}{\partial w_d} \end{pmatrix} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_d \end{pmatrix} = a$$

---

$$\mathcal{L} = y^T y - y^T x w - w^T x^T y + w^T x^T x w$$

① As $y^T y$ is constant term, its derivative will be zero.

② $y^T x w$ can be written as $(x y^T)^T w \qquad (\because (AB)^T = B^T A^T$

Now, $\nabla_w (w^T A) = \nabla_w (A^T w) = A$

Using this $\quad \nabla_w (y^T x w) = \nabla_w (x^T y)^T w = x^T y$

③ Similarly, $\quad \nabla_w (w^T x^T y) = x^T x$

$$\left( \because \nabla_w (w^T A) = \nabla_w (A^T w) = A \right)$$

④ To calculate the derivative of $w^T x^T x w$, we'll use product rule.

$$\left[ d(u \cdot v) = v\, d(u) + u \cdot d(v) \right]$$

$$\nabla_w w^T \underbrace{x^T x w}_{\text{assume constant}} = x^T x w \qquad (\because \nabla_w w^T a = a)$$

$$\nabla_w \underbrace{w^T x^T x}_{\text{assume constant}} w = \nabla_w (x^T x w)^T w \qquad (\because \nabla_w a^T w = a)$$
$$= x^T x w$$

$$\nabla_w \mathcal{L} = 0 - x^T y - x^T y + x^T x w + x^T x w$$

$$= -2 x^T y + 2 x^T x w.$$

---

$$\nabla_w \mathcal{L} = -2 x^T y + 2 x^T x w.$$

Equating to zero

$$\Rightarrow -2 x^T y + 2 x^T x w = 0$$

$$\Rightarrow 2 x^T x w = 2 x^T y$$

$$\Rightarrow x^T x w = x^T y$$

$$\boxed{w = (x^T x)^{-1} x^T y}$$

$$\left( \because \begin{array}{l} \text{if} \\ A B = C \\ B = A^{-1} C \end{array} \right)$$

By equating the values of X and Y, we can compute the values of W

## Comparing Normal equation to GD

- There is no need to for feature scaling in Normal Equation
- GD works well for any number of dimensions.
- But, We need to calculate $(X^T X)^{-1}$ in normal equation