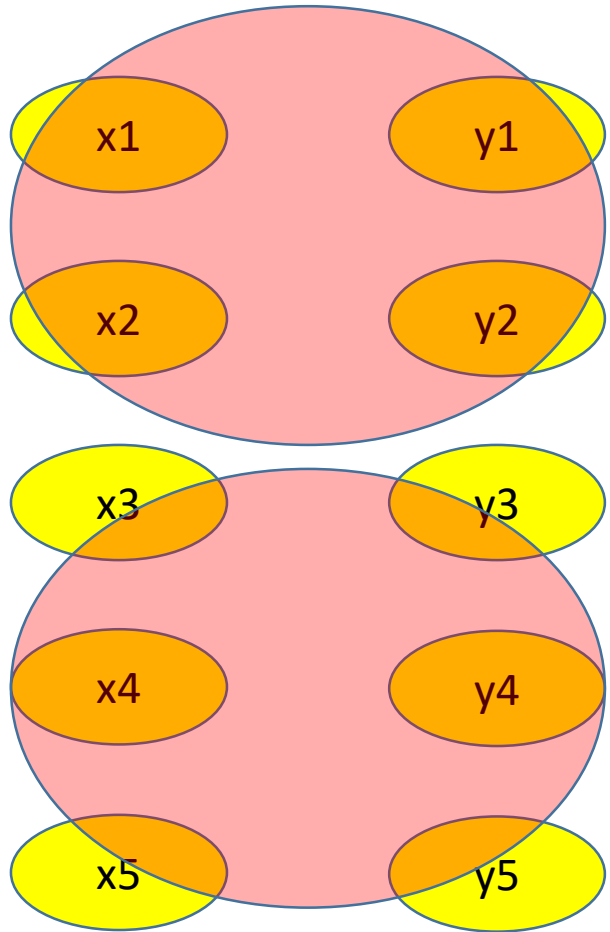


2D Steepest Descent

Dr. Kalidas Y.

In this lecture you will understand how to formulate 2D steepest descent

Left side (X) \rightarrow Machine Learning \rightarrow Right side (Y)



Equation of Line :- $y = m * x + c$

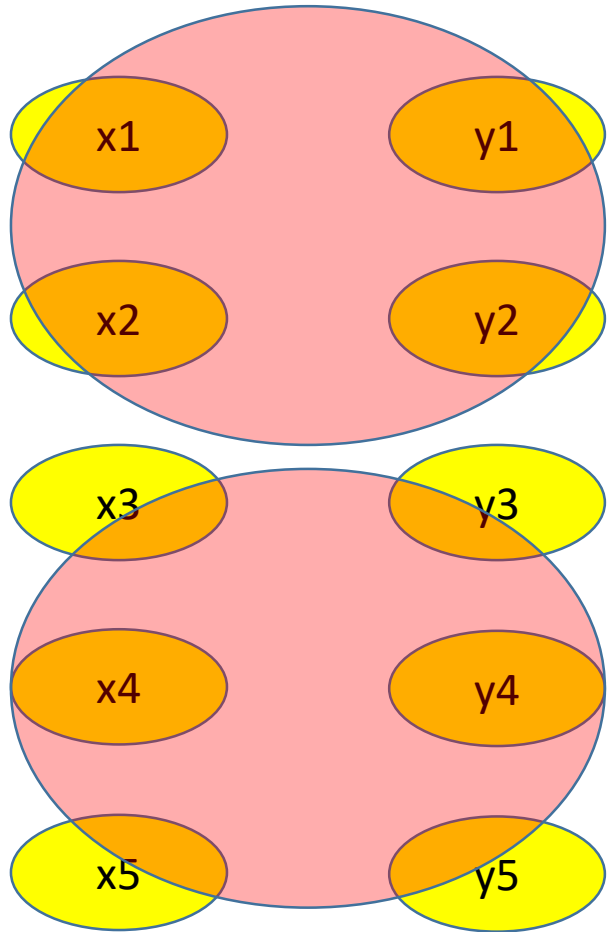
$$L(m, c) = \sum_{i=1}^{i=5} (y_i - (m * x_i + c))^2$$

$$\frac{\partial L}{\partial m}, \frac{\partial L}{\partial c}, \frac{\partial^2 L}{\partial m^2}, \frac{\partial^2 L}{\partial m \partial c}, \frac{\partial^2 L}{\partial c \partial m}$$

You can treat it like any other **parameterized** function

Linear Regression

Left side (X) \rightarrow Machine Learning \rightarrow Right side (Y)



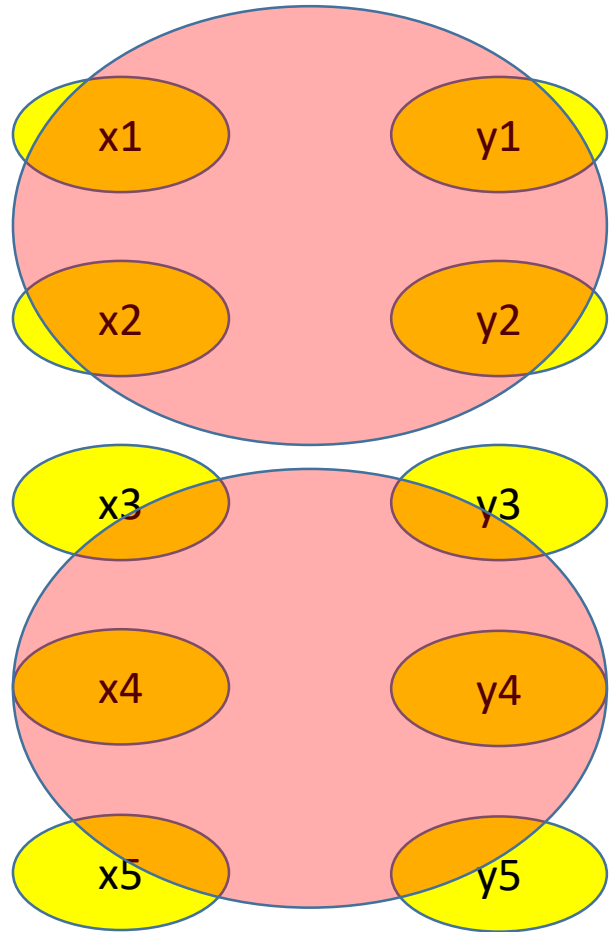
Equation of Line :- $y = m * x + c$

$$L(m, c) = \sum_{i=1}^{i=5} (y_i - (m * x_i + c))^2$$

$$\frac{\partial L}{\partial m} = \sum_{i=1}^{i=5} 2 * (y_i - (m * x_i + c))^1 * (-x_i)$$

Linear Regression

Left side (X) \rightarrow Machine Learning \rightarrow Right side (Y)



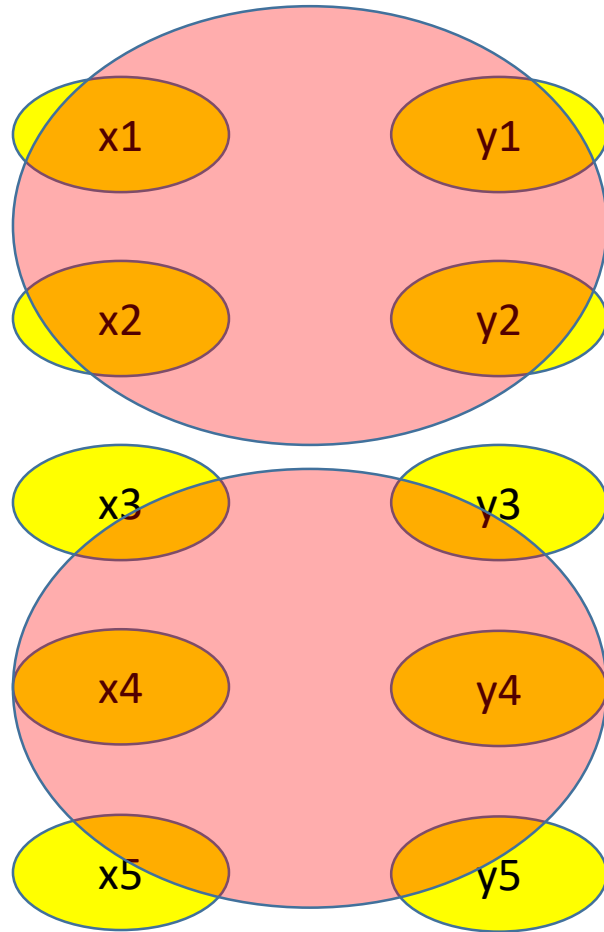
Equation of Line :- $y = m * x + c$

$$L(m, c) = \sum_{i=1}^{i=5} (y_i - (m * x_i + c))^2$$

$$\frac{\partial L}{\partial m} = \sum_{i=1}^{i=5} 2 * (y_i - (m * x_i + c))^1 * (-x_i)$$

$$\frac{\partial L}{\partial c} = \sum_{i=1}^{i=5} 2 * (y_i - (m * x_i + c))^1 * (-1)$$

Left side (X) \rightarrow Machine Learning \rightarrow Right side (Y)



Equation of Line :- $y = m * x + c$

$$L(m, c) = \sum_{i=1}^{i=5} (y_i - (m * x_i + c))^2$$

$$\frac{\partial L}{\partial m} = \sum_{i=1}^{i=5} 2 * (y_i - (m * x_i + c))^1 * (-x_i)$$

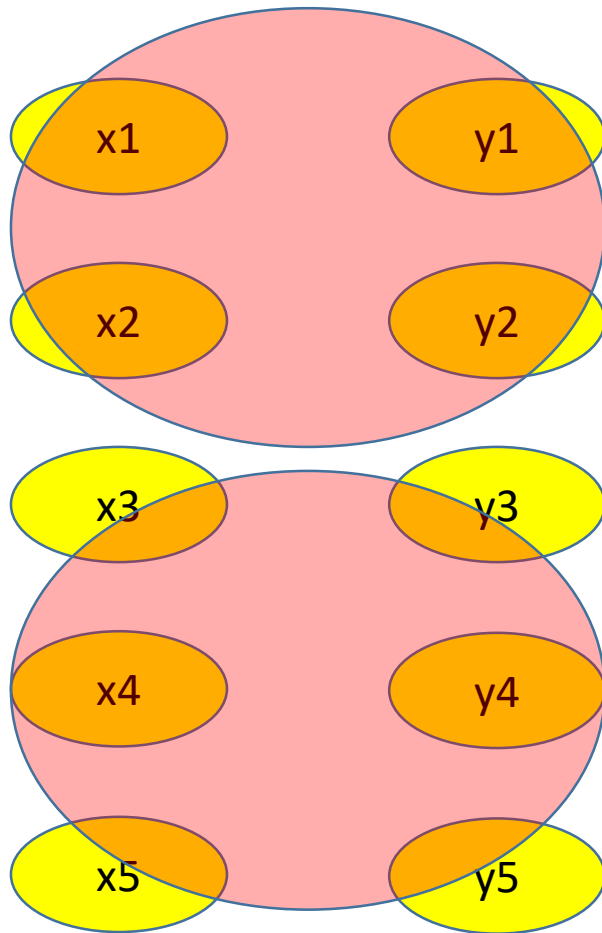
$$\frac{\partial L}{\partial c} = \sum_{i=1}^{i=5} 2 * (y_i - (m * x_i + c))^1 * (-1)$$

$$\frac{\partial^2 L}{\partial m^2} = \sum_{i=1}^{i=5} 2 * (-x_i) * (-x_i) > 0$$

$$\frac{\partial^2 L}{\partial c^2} = \sum_{i=1}^{i=5} 2 * (-1) * (-1) > 0$$

$$\frac{\partial^2 L}{\partial m \partial c} = \sum_{i=1}^{i=5} 2 * (-1) * (-x_i) = \frac{\partial^2 L}{\partial c \partial m}$$

Left side (X) \rightarrow Machine Learning \rightarrow Right side (Y)



Equation of Line :- $y = m * x + c$

$$L(m, c) = \sum_{i=1}^{i=5} (y_i - (m * x_i + c))^2$$

$$= \sum_{i=1}^{i=5} (y_i^2 + m^2 x_i^2 + 2 m x_i c + c^2 - 2 y_i m x_i - 2 y_i c)$$

$L(m1, c1)$

$L(m2, c2)$

...

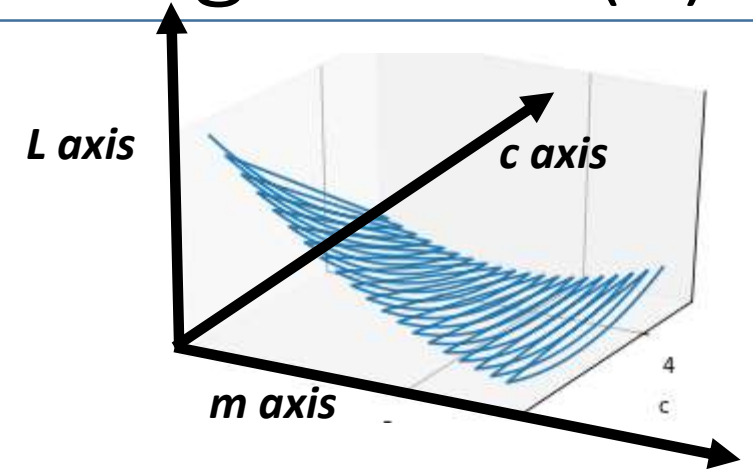
...

...

$L(mN, cN)$

$$= \alpha_1 m^2 + \beta_1 m + \alpha_2 c^2 + \beta_2 c + \gamma$$

For each given, c , the above curve is a **parabola** in m and L axes
For each given m , the above curve is a **parabola** in c and L axes



...sum of two parabolas **NEED NOT BE a parabola**

...some kind of a bowl shaped surface

Consider parabola 1: $y = 3x^2 + 4$

Consider parabola 2: $y = -3x^2 + 7x$

parabola 1 + parabola 2: $y = 7x + 4$

This is a line! Not a parabola

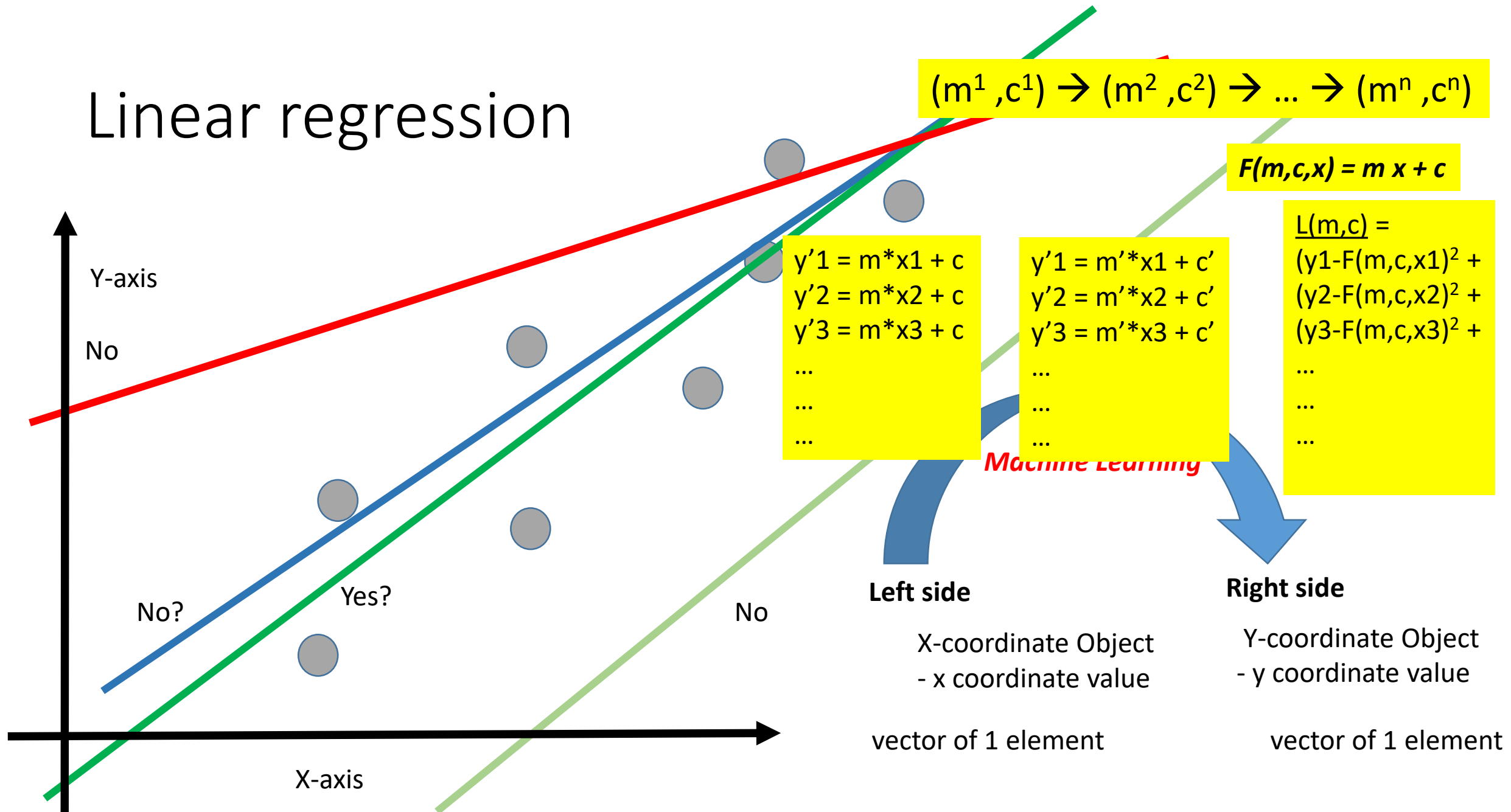
Sum of degree-K polynomials \rightarrow A polynomial with degree $\leq K$

Consider parabola 1: $y = 3x^2 + 4$

Consider parabola 2: $y = \sqrt{7}x^2 + 8x - 13$

parabola 1 + parabola 2: $y = (3 + \sqrt{7})x^2 + 8x - 9$ ***This is a parabola again!***

Linear regression



Matrix Formulation of L(m,c)

$$\bullet X = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \\ x_4 & 1 \\ x_5 & 1 \end{bmatrix}, W = \begin{bmatrix} m \\ c \end{bmatrix}, Y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{bmatrix}$$

$$\bullet L(W) = (XW - Y)^T (XW - Y)$$

Gradient Descent Formulation

Loss function of 2D vector

- $L(W) = (XW - Y)^T (XW - Y)$
- $L(W) = \left(\begin{bmatrix} x_1 & 1 \\ \dots & \dots \\ x_n & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix} \right)^T \left(\begin{bmatrix} x_1 & 1 \\ \dots & \dots \\ x_n & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix} \right)$
- $L(W) = \sum_{i=1}^n (x_{i,1} w_1 + 1 * w_2 - y_i)^2$

Gradient of Loss function (2D scenario)

- $L(W) = \sum_{i=1}^n (x_{i,1}w_1 + w_2 - y_i)^2$

- $\frac{\partial L(W)}{\partial w_1} = \sum_{i=1}^n 2 * (x_{i,1}w_1 + w_2 - y_i)^1 x_{i,1} = 2 * X[:, \text{column 1}]^T (XW - Y)$

- $\frac{\partial L(W)}{\partial w_2} = \sum_{i=1}^n 2 * (x_{i,1}w_1 + w_2 - y_i)^1 = 2 * X[:, \text{column 2}]^T (XW - Y)$

- $\nabla L(W) = \begin{bmatrix} \frac{\partial L(W)}{\partial w_1} \\ \frac{\partial L(W)}{\partial w_2} \end{bmatrix} = \begin{bmatrix} 2 * X[:, 1]^T (XW - Y) \\ 2 * X[:, 2]^T (XW - Y) \end{bmatrix} = 2 * \begin{bmatrix} X[:, 1]^T \\ X[:, 2]^T \end{bmatrix} * (XW - Y)$

- $= 2 X^T (XW - Y)$

Alternatively...

- Start, $m = c = 0$ //or any random
- Iterate 1000 steps (*or whatever*)
 - $m_{(new)} = m_{(old)} - \eta \frac{\partial L}{\partial m} \big|_{m=m_{(old)}, c=c_{(old)}}$
 - $c_{(new)} = c_{(old)} - \eta \frac{\partial L}{\partial c} \big|_{c=c_{(old)}, m=m_{(new)}}$
- Output m, c

Steepest Descent for 2D Loss function

$$W_{(new)} = W_{(old)} - \eta \nabla L(W) \Big|_{W=W_{(old)}}$$

This is a **step size**
a.k.a **learning rate**

This is a **function**

Function computed **at**

specified point of interest

The diagram illustrates the steepest descent formula for a 2D loss function. The formula is written in red: $W_{(new)} = W_{(old)} - \eta \nabla L(W) \Big|_{W=W_{(old)}}$. Annotations in green provide context for the terms: a blue arrow points from the text 'This is a step size a.k.a learning rate' to the symbol η ; another blue arrow points from 'This is a function' to the $\nabla L(W)$ term; a third blue arrow points from 'Function computed at' to the vertical bar notation; and a final blue arrow points from 'specified point of interest' to the $W_{(old)}$ term within the vertical bar.

What are the different ways?

2nd argument is
the *list of Parameters*

BruteForceSolver(*L*, [*m1*, *c1*])

minval = +10000.0

FOR *m* ∈ [−100, ..., +100]

FOR *c* ∈ [−100, ..., +100]

Compute *v* = *L*(*m*, *c*)

IF *v* < *minval*:
 v = *minval*
 m1 = *m*,
 c1 = *c*

1st argument is the *Loss function*

Questions...

- -100 to +100, who gave the range?
- What is the step size?
- What if the solution is highly fine, (3.451, -89.1123)?
- Can we speed up?

What are the different ways?

2nd argument is
the *list of Parameters*

RandomSolver(*L*, [*m1*, *c1*])

minval = +10000.0

FOR ITER= 1:100000

FOR m = RAND(-100, 100)

FOR c = RAND(-100, 100)

Compute $v = L(m, c)$

IF $v < minval$:
 $v = minval$
 $m1 = m$,
 $c1 = c$

1st argument is the *Loss function*

Questions...

- -100 to +100, who gave the range?
- What is the step size?
- What if the solution is highly fine, (3.451, -89.1123)?
- Can we speed up?

What are the different ways?

GradientSolver(*L*, [*m1*, *c1*], *gradL*)

2nd argument is
the *list of Parameters*

1st argument is the *Loss function*

3rd argument is
the *gradient function*

???

What are the different ways?

2nd argument is
the *list of Parameters*

GradientSolver(*L*, [*m*, *c*])

- INIT $m, c = 0$

- Iterate 1000 steps (or whatever)

- $m_{(new)} = m_{(old)} - \frac{\partial L}{\partial m} \big|_{m=m_{(old)}}$

- $c_{(new)} = c_{(old)} - \frac{\partial L}{\partial c} \big|_{c=c_{(old)}}$

1st argument is the *Loss function*

What are the different ways?

SteepestDescentSolver(*L*, [*m1*, *c1*], *gradL*)

//*m1*,*c1* = some initial value is already input to the function

- Iterate 1000 steps (or whatever)

- IF $L(m1, c1) < \text{Threshold}$

- RETURN

- $[m1, c1] = [m1, c1] - \text{gradL}(m1, c1)$

2nd argument is
the *list of Parameters*

1st argument is the *Loss function*

3rd argument is
the *gradient function*

This function outputs a list or vector

Result being used: A symmetric matrix is positive semidefinite if and only if all its eigenvalues are nonnegative.

To prove: The matrix $G = \begin{bmatrix} z & \bar{x} \\ \bar{x} & 1 \end{bmatrix}$ is positive semidefinite.

Proof. Since G is a symmetric matrix, by the above result, it is enough to prove that all the eigenvalues of the matrix G are nonnegative. The eigenvalues of the matrix G are given by the roots of the equation $(z - \lambda)(1 - \lambda) - \bar{x}^2 = 0$.

That is, the roots of the equation

$$\lambda^2 - (z + 1)\lambda + z - \bar{x}^2 = 0. \quad (1)$$

This gives

$$\lambda = \frac{(z + 1) \pm \sqrt{(z + 1)^2 - 4(z - \bar{x}^2)}}{2}.$$

Claim 1. All the roots of equation (1) are real numbers.

Proof of Claim 1. The discriminant of equation (1) is given by $(z + 1)^2 - 4(z - \bar{x}^2)$.

Now,

$$\begin{aligned} (z + 1)^2 - 4(z - \bar{x}^2) &= z^2 + 2z + 1 - 4z + 4\bar{x}^2 \\ &= z^2 - 2z + 1 + 4\bar{x}^2 \\ &= (z - 1)^2 + 4\bar{x}^2, \end{aligned}$$

which is nonnegative, as it is the sum of squares of real numbers.

This proves Claim 1.

Claim 2. All the roots of equation (1) are nonnegative.

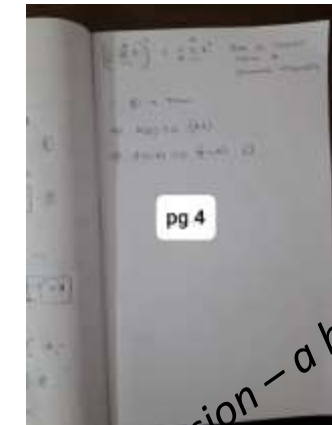
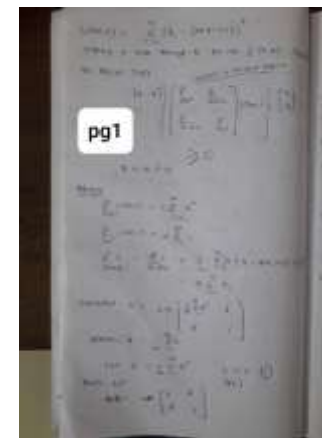
Proof of Claim 2. By Jensen's inequality, we have $\bar{x}^2 \leq z$. This yields

$$\begin{aligned} \bar{x}^2 &\leq z \\ \Rightarrow -4(z - \bar{x}^2) &\leq 0 \\ \Rightarrow (z + 1)^2 - 4(z - \bar{x}^2) &\leq (z + 1)^2 \\ \Rightarrow \sqrt{(z + 1)^2 - 4(z - \bar{x}^2)} &\leq z + 1 \quad (\text{Since } (z + 1) > 0) \\ \Rightarrow \frac{(z + 1) \pm \sqrt{(z + 1)^2 - 4(z - \bar{x}^2)}}{2} &\geq 0 \end{aligned}$$

That is, the roots of the equation (1), and hence the eigenvalues of the matrix G are nonnegative. Hence the matrix G is positive semidefinite.

Remarks.

1. Here, Claim 1 (nonnegativity of the discriminant) is redundant, as the proof of Claim 2 does not assume that the roots are real.
2. The quadratic loss function is infinitely differentiable. So the Hessian matrix is symmetric and the above technique can be used to check if the Hessian is positive semidefinite in the multivariate scenario, with even more than two variables.



Dr. Kalidas's version – a bit lengthier proof though

Thanks to Dr. Lakshmi, IISER Mathematics Faculty for the proof on left (simplified version)

Data, Label, Parameter – Simple supervised case

