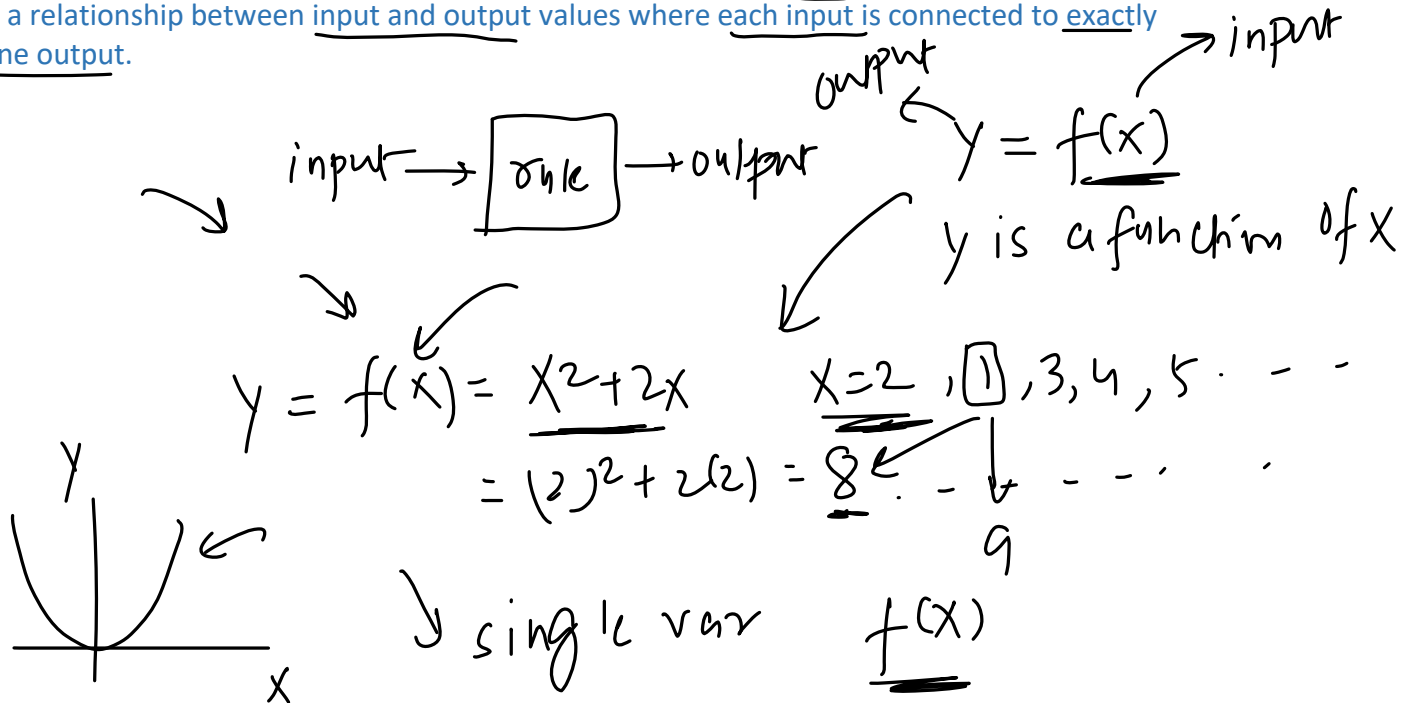


# Functions

21 April 2023 14:28

A function is a mathematical rule that takes an input value, processes it according to a specific formula or set of instructions, and produces a unique output value. In other words, a function is a relationship between input and output values where each input is connected to exactly one output.



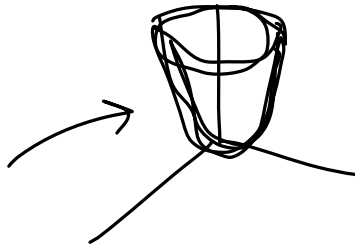
# Multivariable Functions

21 April 2023 16:35

$$z = f(x, y) = x^2 + y^2$$

↑      ↑ ↑  
output input

z+1



$$f(x) \quad 1+1 \rightarrow \underline{2}$$

n+1 dim

$$x=1 \quad y=1 \rightarrow z=2$$

$$\underline{x=1 \quad y=1} \rightarrow \underline{z=2}$$

$$\underline{f(x_1, x_2, x_3, \dots, x_n)}$$

# Parameters in a Function

21 April 2023 16:36

$$\underline{ax^2 + bx + c}$$

In mathematics, parameters of a function are the variables that are used to define the behaviour of the function. The parameters influence the function's output by determining how the input values are processed.

The parameters are the constants or coefficients that appear in the function's formula. For example, in the quadratic function  $f(x) = ax^2 + bx + c$ , 'a', 'b', and 'c' are the parameters of the function. By changing the values of these parameters, you can modify the shape and position of the parabola represented by the function.

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

↑      ↑  
parameter input

$$6x^2$$

$$y = f(x) = 9x^2$$

↑  
parameter

# [ML models] as Mathematical Function

21 April 2023

16:36

cgpa	iq	placement
9	90	9
10	100	0
8	80	8

ml model

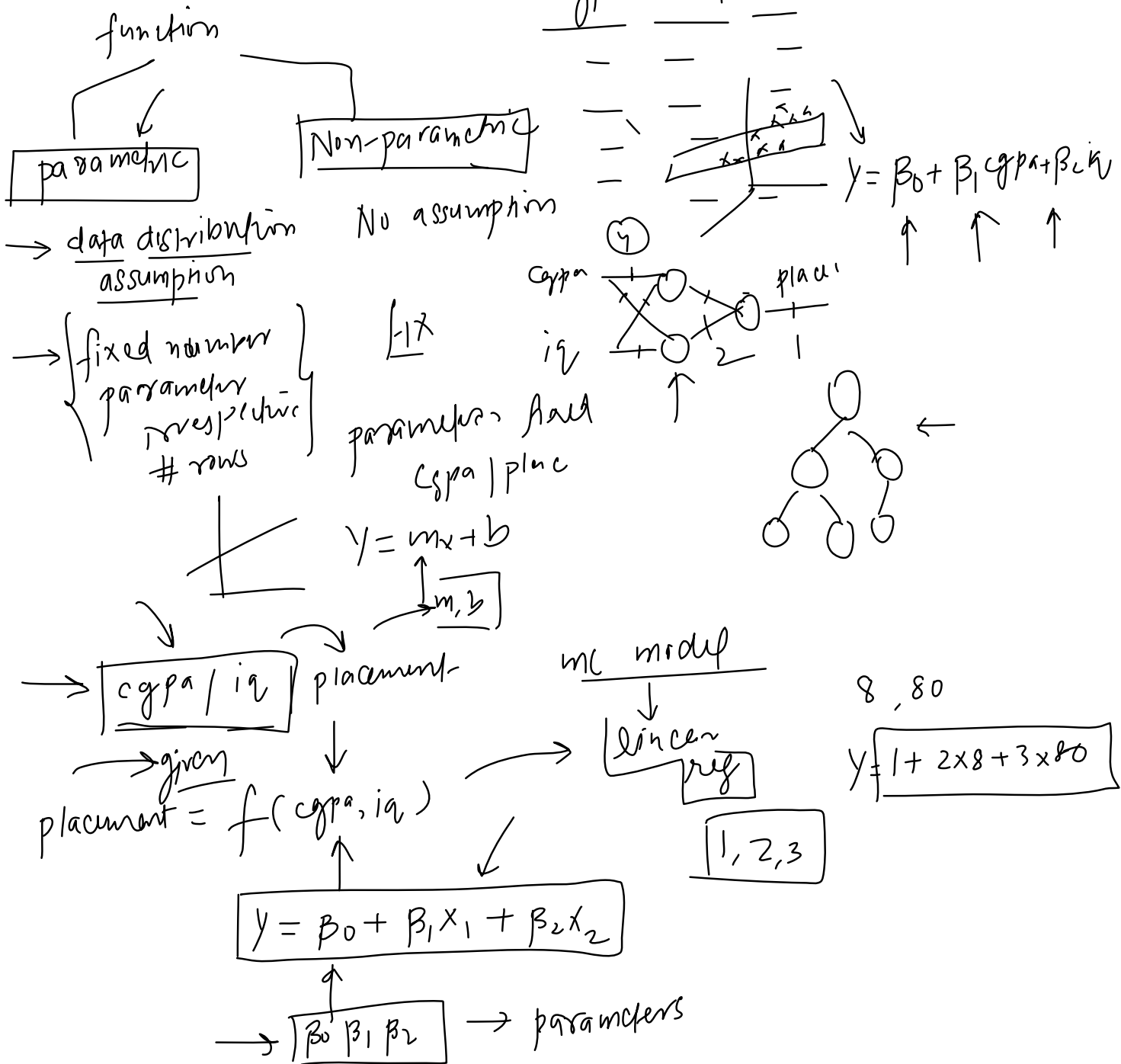
$y \rightarrow \text{placement}$

$x \rightarrow \text{cgpa, iq}$

$$\text{placement} = f(\text{cgpa}, \text{iq})$$

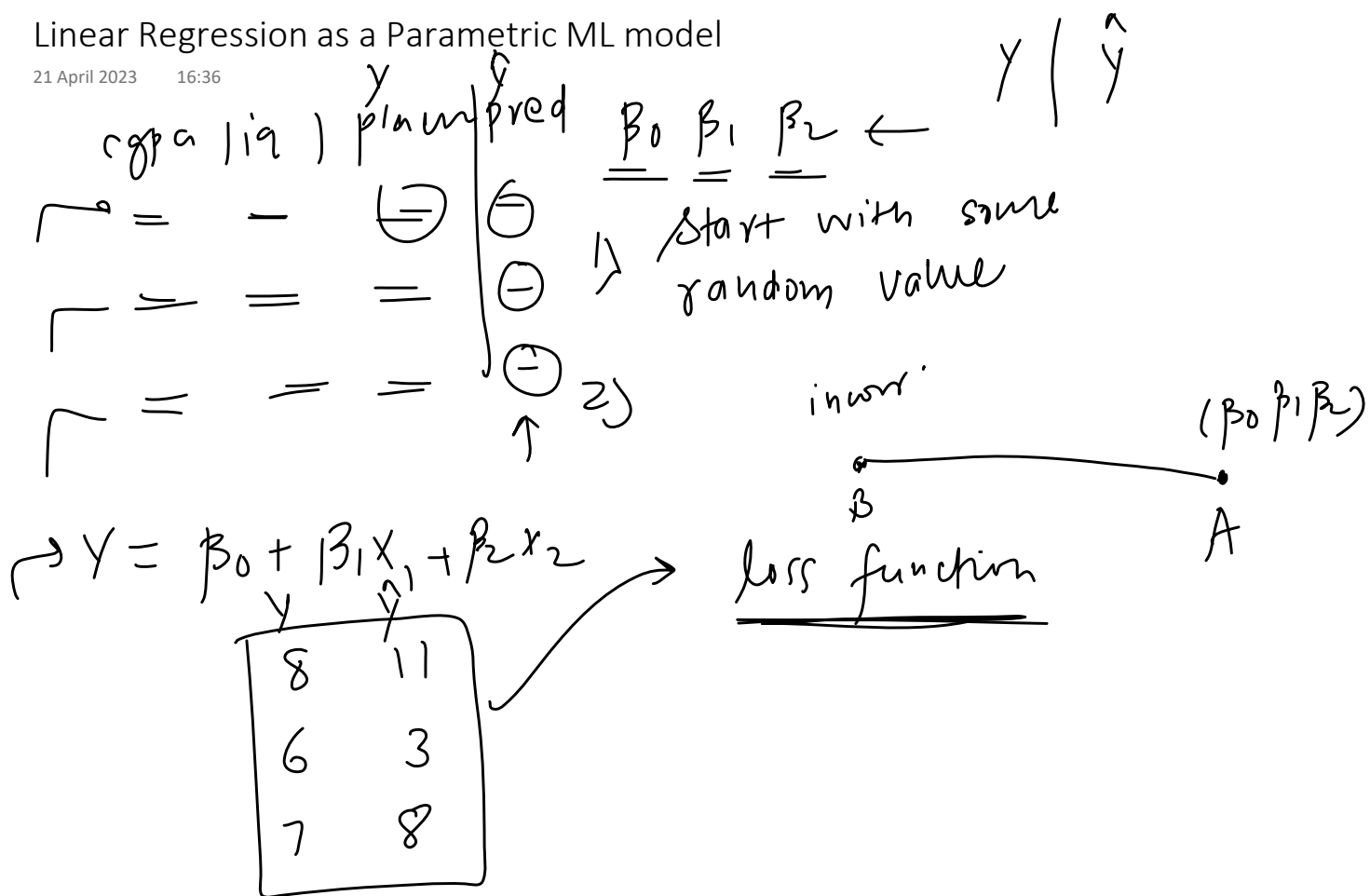
# Parametric Vs Non-Parametric ML models

21 April 2023 16:36



# Linear Regression as a Parametric ML model

21 April 2023 16:36

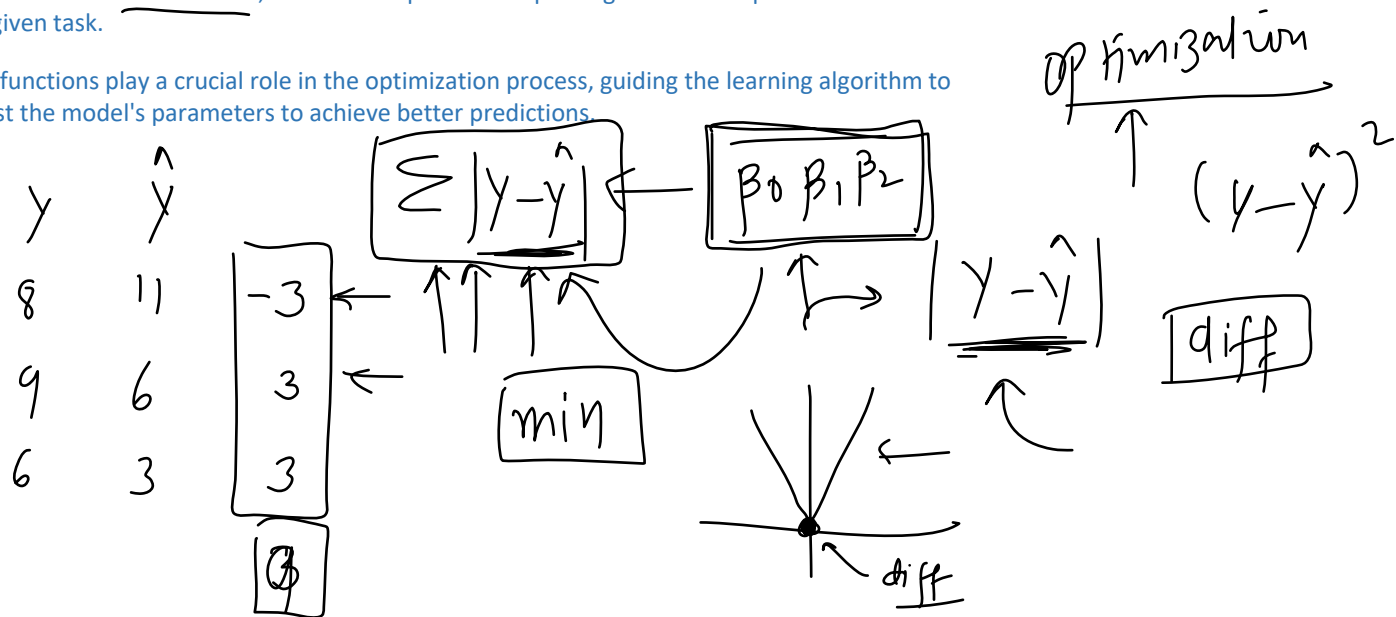


# Loss Function

21 April 2023 16:37

A loss function, also known as a cost function or objective function, is a mathematical function that measures the difference between the predicted output and the actual target values in a machine learning model. The primary goal of training a machine learning model is to minimize the value of the loss function, which corresponds to improving the model's performance on the given task.

Loss functions play a crucial role in the optimization process, guiding the learning algorithm to adjust the model's parameters to achieve better predictions.





# How to select a good Loss Function

21 April 2023 16:37

1. Problem type: The choice of a loss function depends on the type of problem you are solving. For example, in regression tasks, mean squared error (MSE) or mean absolute error (MAE) are commonly used. For binary classification, cross-entropy loss or hinge loss can be employed. For multi-class classification, categorical cross-entropy or multi-class hinge loss can be used. Choose a loss function that aligns with the objectives of the specific problem you are addressing.
2. Robustness to outliers: Some loss functions, like mean squared error, are more sensitive to outliers, which can lead to a model that is overly influenced by extreme values. If your dataset contains outliers or is prone to noise, consider using a loss function that is more robust to outliers, such as mean absolute error (MAE) or Huber loss.
3. Interpretability and ease of use: A good loss function should be interpretable and easy to implement. Simple loss functions like mean squared error or cross-entropy loss are widely used because they are easy to understand, compute, and differentiate. When possible, opt for a loss function that is easy to work with and can be easily incorporated into your optimization process.
4. Differentiability: Most optimization algorithms, like gradient descent, require the loss function to be differentiable. Choose a loss function that has continuous first-order derivatives, which makes it easier to compute the gradients needed for optimization.
5. Compatibility with the model: Ensure that the chosen loss function is compatible with the model architecture you are using. Some models have specific requirements or assumptions about the loss function. For example, linear regression assumes a Gaussian noise distribution, which is why mean squared error is a suitable loss function in that case.

# Calculating Parameters From a Loss Function (the easy way and problem)

21 April 2023 16:37

cgpa	iq	placement
8	80	8
7	70	7
6	60	6

$\hat{y}$
11
6
3

$\beta_0$	$\beta_1$	$\beta_2$
1	2	3

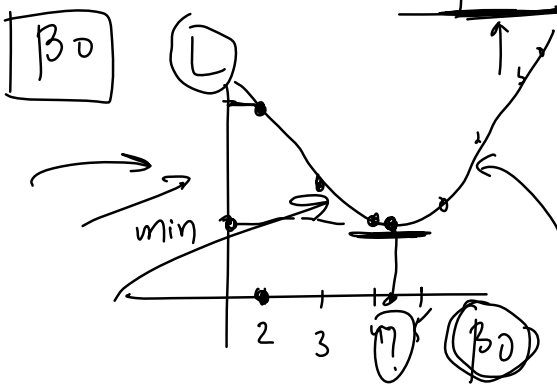
$$\sum (y - \hat{y})^2 \quad (3)^2 + (1)^2 + (3)^2$$

$$\frac{1}{3} 19 = 6.3$$

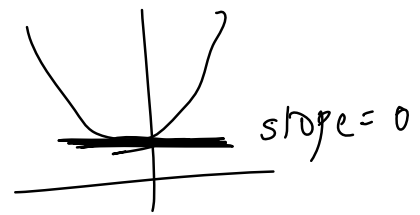
$$L(\beta_0, \beta_1, \beta_2) = \arg \min_{\beta_0, \beta_1, \beta_2} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$1.5 \quad 2.5 \quad 3.5$$

parabolic



optimization  
activation function



$$\frac{dL}{d\beta_0} = 0$$

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

log(1+ex) linear  
mse -> convex  
OLS

$$\hat{y}_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

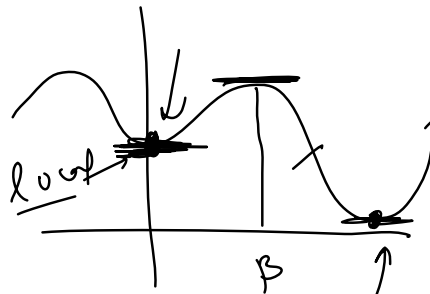
$$y_i = \beta_0 + x_1 + 2x_2$$

$$\frac{dL}{d\beta_0} = \sum_{i=1}^n (y_i - \beta_0 - x_1 - 2x_2)$$

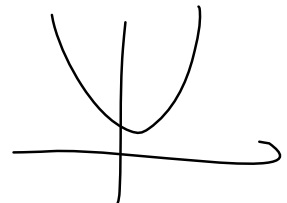
$$-2(y_i - \beta_0 - x_1 - 2x_2) = 0$$

$$(\sum y_i - x_1 - 2x_2) = \beta_0$$

$$\text{slope} = 0$$



$$\frac{dL}{d\beta_0} = 0$$



global  
minimized summation

$$\boxed{\frac{\partial L}{\partial \beta_0} = 0} \rightarrow \begin{matrix} \text{global} \\ \text{limited summation} \end{matrix} \quad \swarrow$$

## Problem with the easy way

21 April 2023 16:44

1. **Non-convexity:** The loss function may not always be convex, meaning that it might have multiple local minima and maxima. In such cases, setting the gradient to zero might lead to a local minimum or maximum, which is not necessarily the global minimum (the optimal solution).
2. **Complexity:** For some models, the loss function can be highly complex, and finding the analytical solution by setting the gradient to zero might be computationally expensive or even impossible. This is particularly true for deep learning models, where the loss functions involve a large number of parameters and complex relationships between them.
3. **Scalability:** In large-scale machine learning problems with massive amounts of data or high-dimensional feature spaces, computing the analytical solution by setting the gradient to zero can be computationally prohibitive due to the high cost of processing and storing the data.
4. **Online learning and streaming data:** In some applications, the data is not available all at once but arrives in a continuous stream. In these scenarios, models need to be updated incrementally as new data arrives, and an analytical solution would not be practical. Gradient descent and its variants, such as stochastic gradient descent, are well-suited for online learning and handling streaming data.

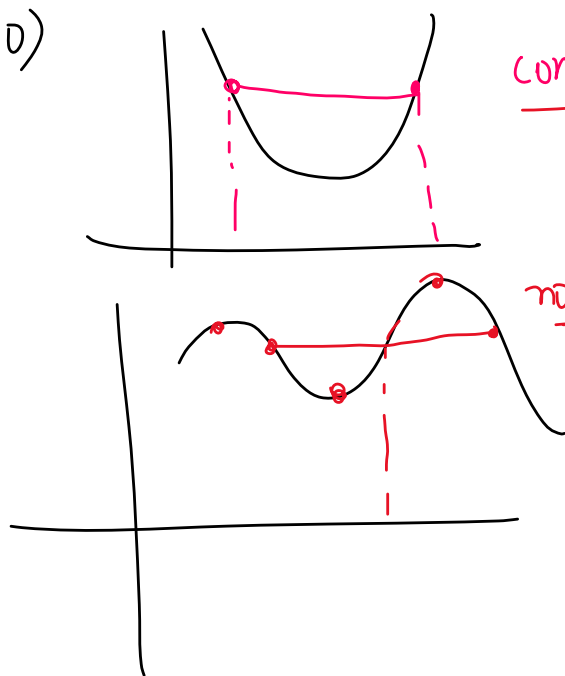
$$\beta = (X^T X)^{-1} X^T Y$$

↑  
real time

# Convex And Non Convex Loss Functions

21 April 2023 16:38

$L(\beta_0)$



convex function

only 1

global  
minima

non-convex function

# Gradient Descent

21 April 2023 16:38

optimization techniques

linear reg

logistic

- perceptron

deep learning

first order der

$$\frac{dL}{d\beta_0}$$

$$f'(x)$$

$$f''(x)$$

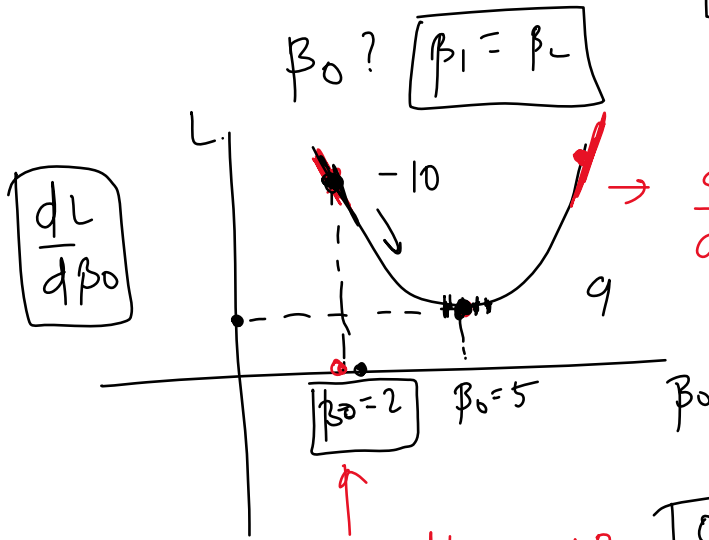
$$\frac{d^2L}{d\beta^2}$$

SVM → quadratic prog  
PCA

convex

maxima

learning



$$\beta_{0, \text{new}} = \beta_{0, \text{old}} - \eta \frac{dL}{d\beta}$$

$$2 - (-10) = 12$$

$$12 - 10 = 2$$

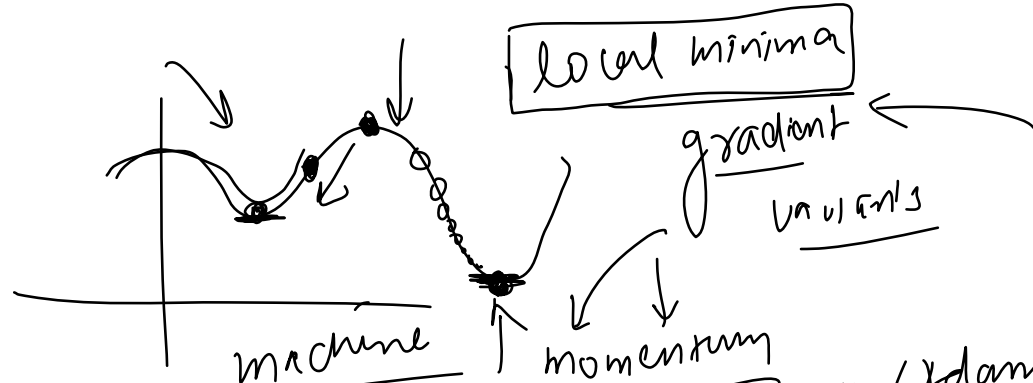
$$\beta_0 = 2 \quad \frac{dL}{d\beta} = -10$$

$$0.01$$

$$\beta_0 = \beta_{0, \text{old}} - \eta \frac{dL}{d\beta}$$

$$\beta_0 = 2 + 0.01 \times 10 = 2.1$$

$$2.1 + 0.01 \times 9 = 2.19$$



deep learning

12

machine learning

$$\beta_{n+1} - \text{slope} = 0$$

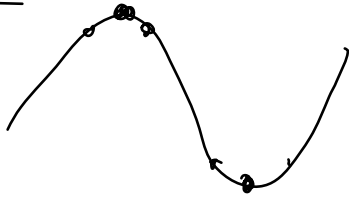
$$\beta_{\text{new}} = \beta_{\text{old}}$$

1

$$\beta_{0_{new}} = \beta_{old} - \text{slope} = 0$$

$$\beta_{new} = \beta_{old}$$

epochs

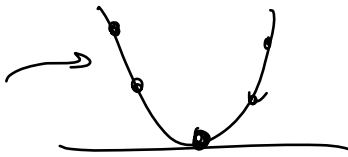


linear reg

$$\boxed{\beta = 5} \quad \begin{matrix} \text{min} & \text{max} \end{matrix}$$

$$f(\beta) = \beta^2 + 2 = 27$$

$$\beta = 4.9$$



convex  
minima

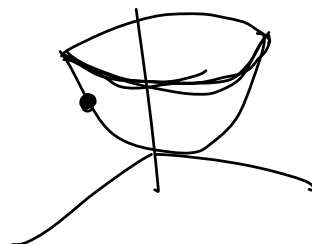
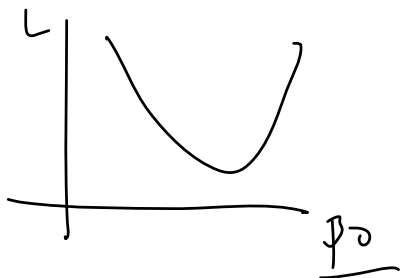
$$\text{epochs } (4.9)^2 + 2 = \boxed{26}$$

$$\boxed{\beta_n = \beta_0 - \eta \frac{\partial L}{\partial \beta}}$$

# Gradient Descent with multiple Parameters

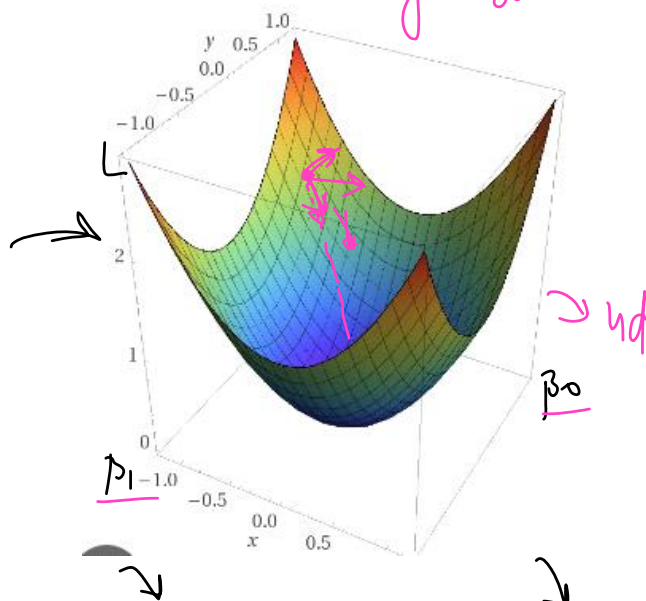
21 April 2023 16:39

$$L(\beta_0) \rightarrow \begin{matrix} \text{cols} & x_1 & x_2 & \dots & x_n \\ \beta_0 & \beta_1 & \beta_2 & \dots & \beta_n \end{matrix} \quad \begin{matrix} \\ (n+1) \end{matrix} \quad \begin{matrix} L & \beta_0 & \beta_1 \end{matrix}$$



gradient descent

$$\frac{\partial L}{\partial \beta_1} \quad \frac{\partial L}{\partial \beta_0} \quad \beta_0 \quad \beta_1 \quad \beta_2$$



$$\begin{aligned} \beta_0 &= \beta_0 - \eta \frac{\partial L}{\partial \beta_0} \\ \beta_1 &= \beta_1 - \eta \frac{\partial L}{\partial \beta_1} \end{aligned}$$

$$\beta_2 = \beta_2 - \eta \frac{\partial L}{\partial \beta_2}$$



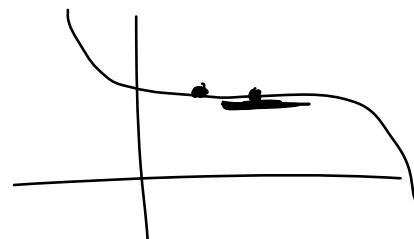
$$\nabla L(\beta_0 \dots \beta_n) = \left( \frac{\partial L}{\partial \beta_0}, \frac{\partial L}{\partial \beta_1}, \dots, \frac{\partial L}{\partial \beta_n} \right)$$



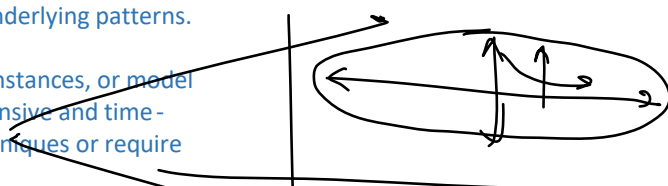
# [Problems faced in Optimization] → 9D

21 April 2023 16:39

1. **Non-convexity**: For many machine learning models, such as artificial neural networks, the loss function is non-convex, which means it has a complex landscape with multiple local minima, maxima, and saddle points. This makes it difficult for optimization algorithms to find the global minimum and can result in suboptimal solutions.
2. **Ill-conditioning**: The loss function may be ill-conditioned, meaning the gradients in some dimensions are much larger than in others. This can cause gradient-based optimization algorithms, such as gradient descent, to oscillate and converge slowly.
3. **Vanishing and exploding gradients**: In deep neural networks, the gradients can become very small (vanish) or very large (explode) as they propagate through the layers. This can lead to slow convergence or unstable training dynamics, making it difficult to optimize the loss function.
4. **Overfitting**: When optimizing the loss function, the algorithm may overfit the training data, resulting in a model that performs poorly on unseen data. This occurs when the model is too complex and learns the noise in the training data instead of the underlying patterns.
5. **Scalability**: For large-scale problems with a high number of features, instances, or model parameters, optimizing the loss function can be computationally expensive and time-consuming. This can limit the applicability of certain optimization techniques or require significant computational resources.



$$\beta_0 = \beta_0 - \eta \text{slope}$$

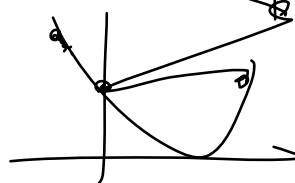


✓  
neuro

$$\frac{dL}{d\beta} \rightarrow \text{small big}$$

$$(m, b) \rightarrow \boxed{10 \text{ thousand}} \rightarrow \text{dist}$$

↓  
60 billion



## Other optimization techniques

21 April 2023 16:39

gradient  $\rightarrow$  d/s

constrained  $\leftarrow$   
SVM / PCA

argmin  $\beta_0$   
 $(y - \hat{y})^2$   
given  $\hat{y} = 2$   $\leftarrow$  SVM / PCA  
constraint