

# Deep Learning with Keras and TensorFlow



# **Model Optimization and Performance Improvement**



# Learning Objectives

By the end of this lesson, you will be able to:

- Optimize a deep learning model to get the most accurate results
- Implement optimization algorithms like SGD, Momentum, AdaGrad, and Adadelata to improve model performance
- Use advanced optimizers like Adam, RMSprop, and NAG to enhance convergence and stability in deep learning models
- Implement batch normalization using Keras to improve training efficiency and model performance



# Learning Objectives

By the end of this lesson, you will be able to:

- 👁️ Apply regularization techniques to prevent overfitting and enhance model generalization using dropout and early stopping
- 👁️ Analyze the issues of vanishing and exploding gradients to understand their impact on model training
- 👁️ Distinguish between interpretability and explainability to better evaluate machine learning models



## Business Scenario

XYZ Corp. is an online retail company experiencing high cart abandonment rates on its website. It wants to optimize its deep learning-based product recommendation system to improve customer retention rates.

XYZ partners with a machine learning solutions provider that recommends utilizing the Adadelta optimization algorithm to improve the convergence speed of the neural network model. The provider also suggests implementing batch normalization to expedite the training process and mitigate overfitting. XYZ incorporates dropout and early stopping techniques to prevent overfitting and enhance the model's generalization capability.

Upon training and deploying the model, the company observes a substantial reduction in cart abandonment rates, resulting in increased sales revenue and improved customer satisfaction.

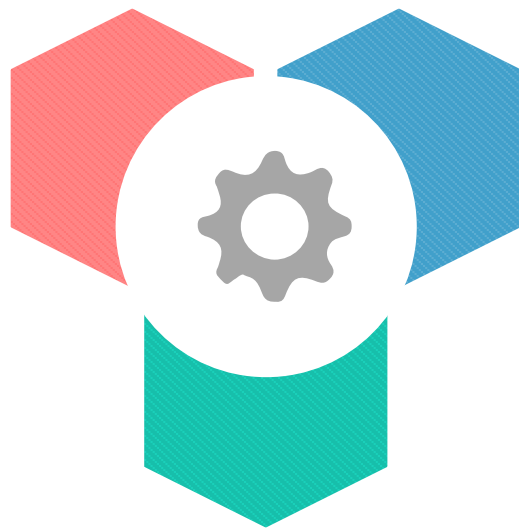




# **Introduction to Optimization Algorithms**

# What Is Optimization?

Minimizing or maximizing any mathematical expression involves finding the values that result in the lowest or highest outcome. This makes the predictions as accurate and optimized as possible.



The parameters of a model are adjusted and changed during the training phase to minimize the loss function.

## Example: Optimization

Building bridges involves balancing weight capacity, cost, safety, and material use. This equilibrium is referred to as optimization.



- Aim for maximum load capacity given the design and usage
- Minimize material usage without compromising integrity
- Prioritize cost-effectiveness
- Uphold stringent safety standards



## Example: Optimization

The optimization of warehouse location and placement in logistics aims to minimize shipping costs and improve efficiency.



- Choose locations near transport routes and customers
- Consider supplier proximity for efficient inbound logistics
- Plan the warehouse layout for optimal space and goods movement

# Optimization Algorithms

To minimize the loss function, optimization algorithms iteratively change the model's parameters during training. In deep learning, optimization algorithms optimize the cost function **J**.

The equation is represented as follows:

$$J(W, b) = \sum_{i=1}^m \frac{L}{m} (y'_i, y_i)$$

# Optimization Algorithms: Equation

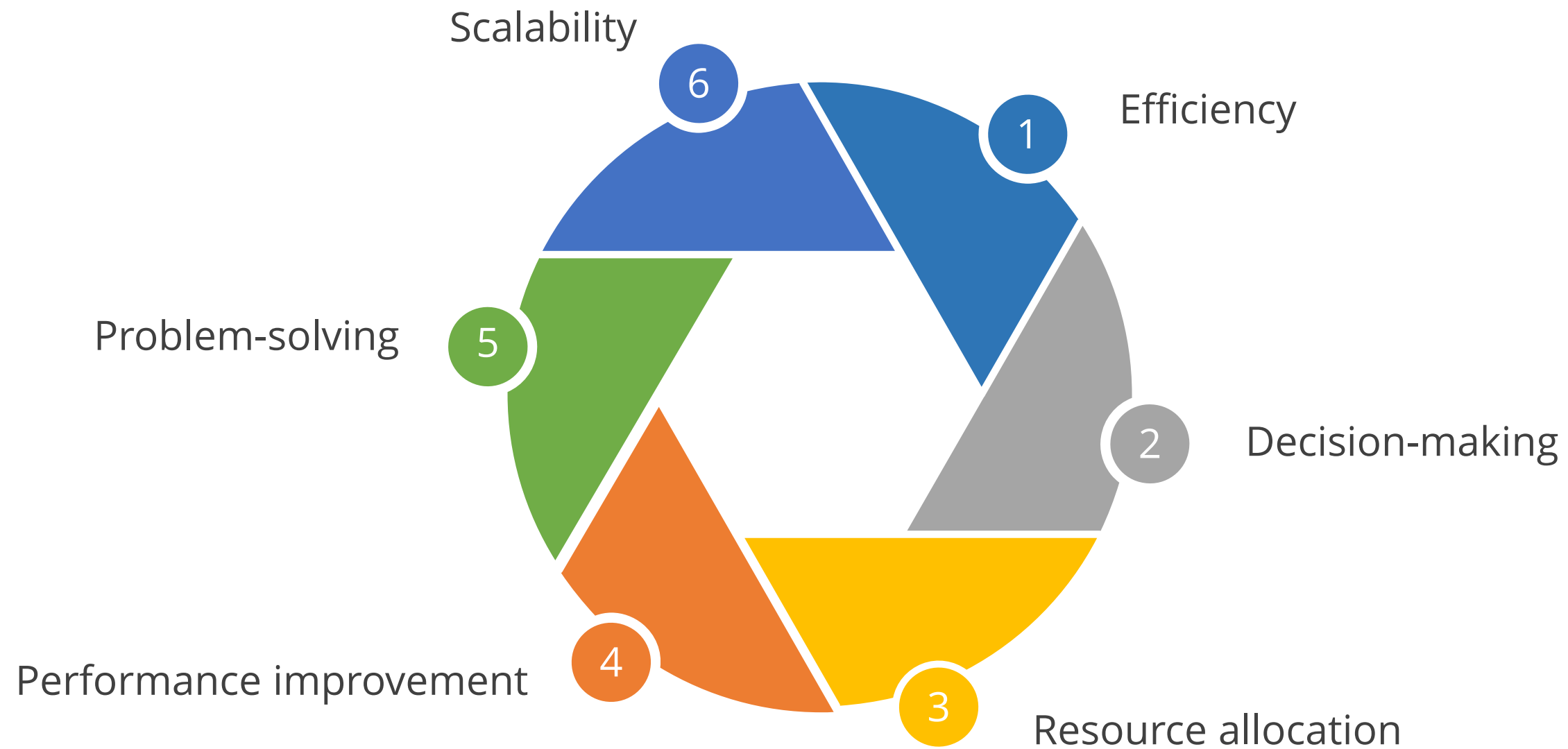
The equation is represented as follows:

$$J(W, b) = \sum_{i=1}^m \frac{L}{m}(y'_i, y_i)$$

- The value of the cost function  $J$  is the mean of the loss  $L$  between the predicted value  $y'$  and the actual value  $y$ .
- The value  $y'$  is obtained during the forward propagation step and uses the weights  $W$  and biases  $b$  of the network.
- With the help of optimization algorithms, you minimize the value of the cost function  $J$  by updating the values of the trainable parameters  $W$  and  $b$ .

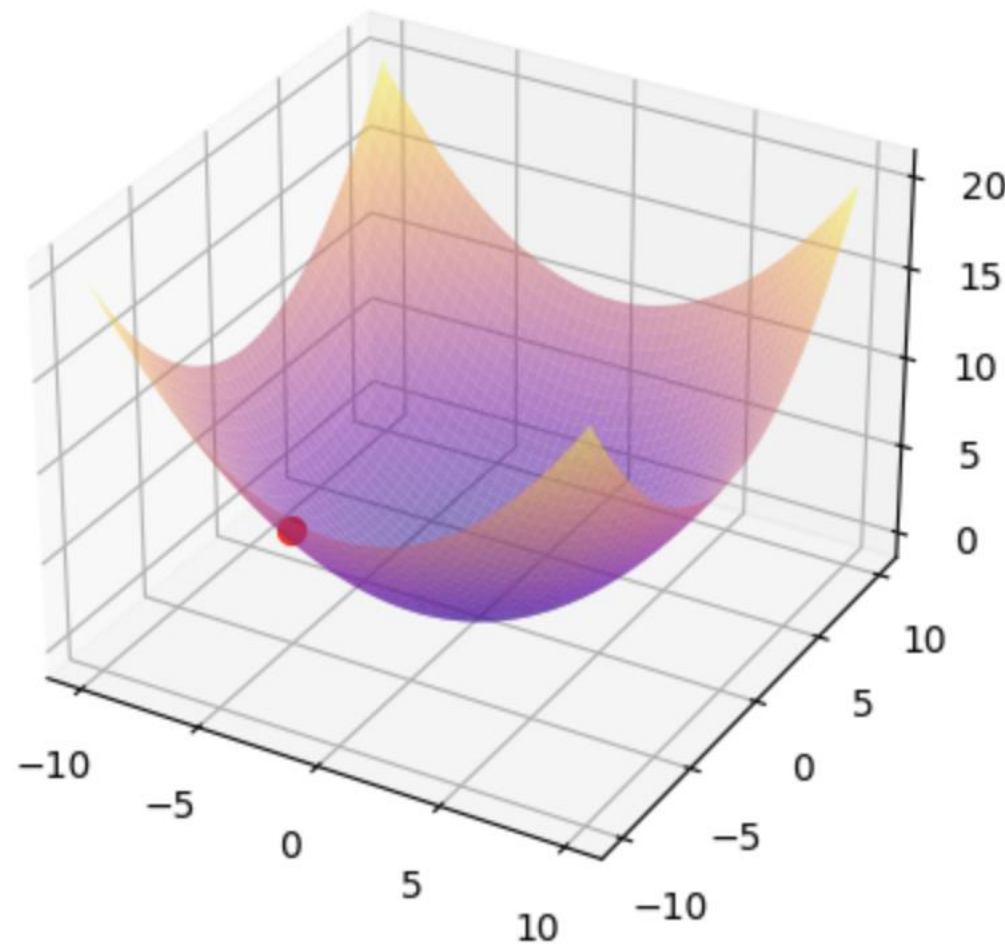
# Importance of Optimization Algorithms

The following list of factors illustrates the significance of optimization algorithms:



# What Is Optimizer?

It's the algorithms or methods that reduce losses by changing neural network attributes such as weights and learning rates.



Optimizers establish a connection between the loss function and model parameters, modifying the model based on the loss function's output.

They manipulate the neural network weights, molding and refining the model to achieve the highest possible accuracy.

# Example of an Optimizer

The loss function serves as a roadmap for the optimizer to indicate if it has taken the correct or wrong path.

Consider hikers attempting to descend a mountain with a blindfold.



The hikers can't determine which way to travel but can tell if they are going down (making progress) or up (losing progress).

If the hikers continue downhill, they will eventually reach the bottom.

# Types of Optimizers

**GD**  
**(Gradient Descent)**

**SGD (Stochastic  
Gradient Descent )**

**Momentum**

**AdaDelta**

**RMSprop**

**Adam**

**AdaGrad**

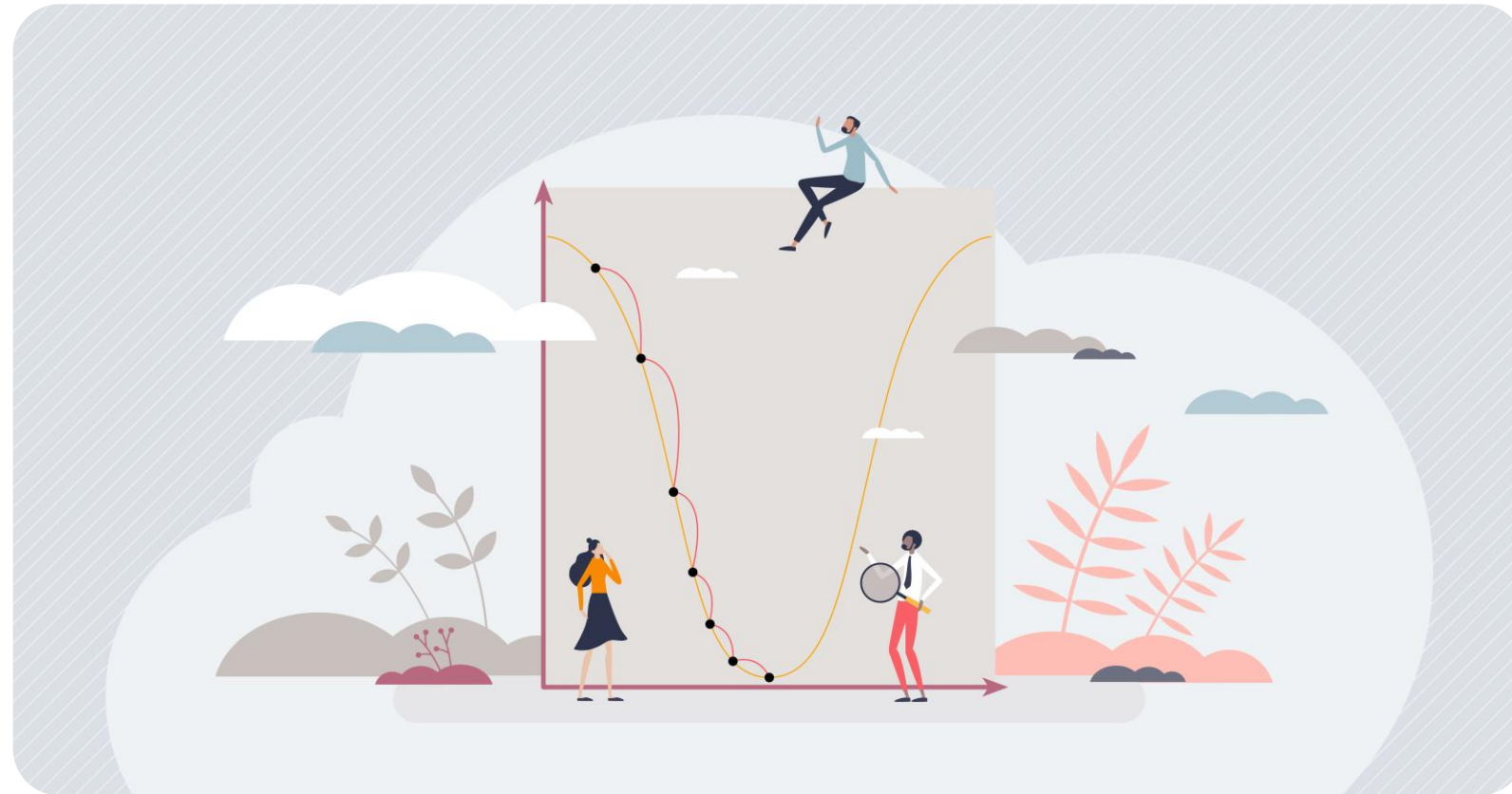


# **Introduction to Stochastic Gradient Descent**



# What Is Gradient Descent?

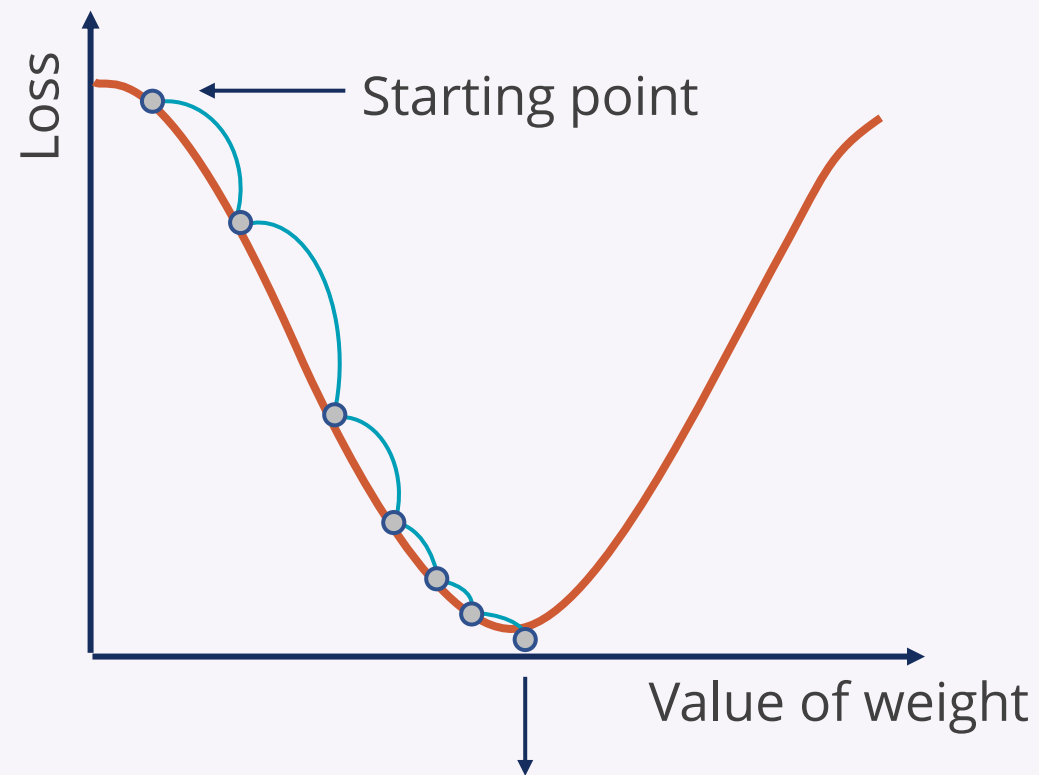
Gradient descent is a typical optimization algorithm and the foundation for training deep learning models.



Gradient descent iteratively adjusts the parameters in the direction that minimizes the loss function, gradually improving the model's fit to the data.

# Gradient Descent

This technique involves calculating how much the model's error changes with its changes in parameters.



The point of convergence is where the cost function is at its minimum.

It modifies those parameters in a way that will drastically decrease the error.

The optimal weights for a model are not known before training but can be found through trial and error using the loss function.

# Gradient Descent

The goal of gradient descent is to reduce the difference between what the model predicts and the actual outcomes. To achieve this, it uses learning rate and direction.

The learning rate (also referred to as step size or alpha) is the number of steps taken to reach the minimum.

Gradient descent uses direction to gradually arrive at the local or global minimum (that is, the point of convergence).

# Gradient Descent

- GD minimizes the cost function **J** and obtains the optimal weight **W** and bias **b**.

- This equation represents a change in matrix **W**.

$$W = W - \alpha \frac{\partial}{\partial W} J(W)$$

Learning  
rate

- This equation represents a change in bias **b**.

$$b = b - \alpha \frac{\partial}{\partial b} J(b)$$

- The learning rate and derivatives of J for W and b determine the change in value. This process is repeated until **J** has been minimized.

# Gradient Descent

In gradient descent, the parameter  $W$  is updated iteratively to minimize the loss function  $J(W)$ .

Gradient Descent update rule:

$$W = W - \alpha \frac{\partial}{\partial W} J(W)$$

Here,  $\alpha$  represents the learning rate, a positive scalar controlling the step size.

Here,  $\alpha$  (the learning rate) scales the derivative, and  $W$  is adjusted in the opposite direction of the gradient. In both cases, the goal is to adjust  $W$  in a direction that reduces  $J(W)$  and approaches the minimum loss function during training.

# Gradient Descent

The dataset is shuffled for each iteration to avoid biasing the descent and improve the randomness of the mini-batch selection process.

Steps to be followed:

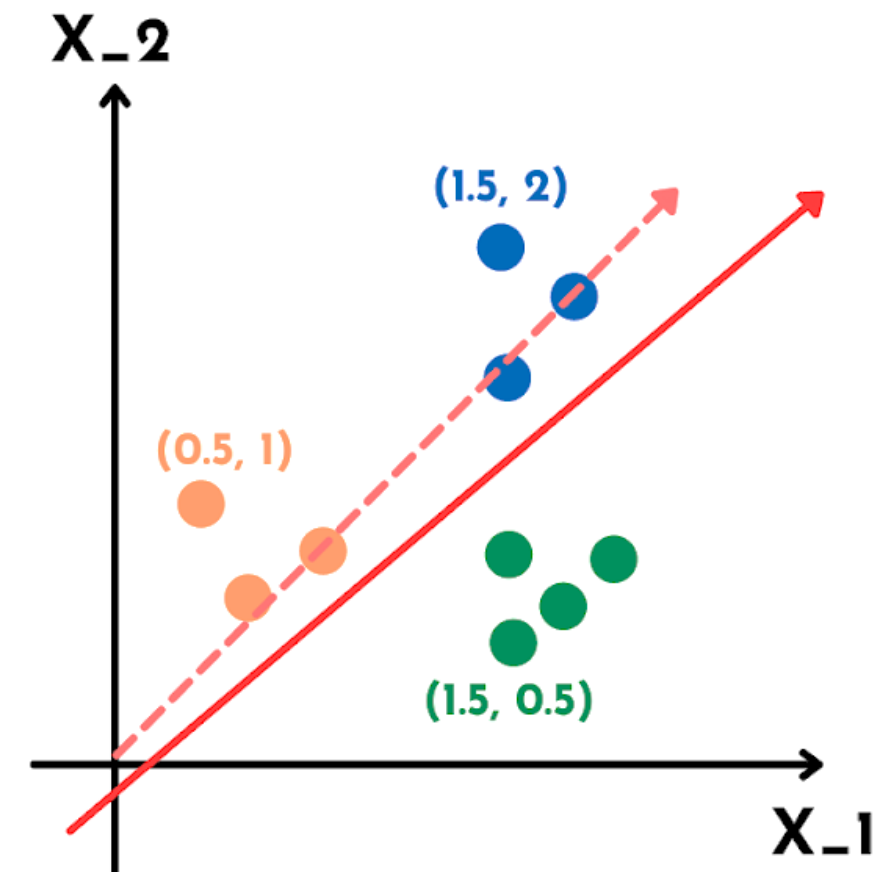
Given equation:  $X_2 = m(X_1) + b$   
let,  $m = 1$  and  $b = 0$

$$\text{Loss} = \frac{1}{N} \sum_{i=1}^N (y_i - (mx_i + b))^2$$

1. Compute the gradients:

$$d(\text{loss})/d(m) = -2*[X_2 - m(X_1) - b]*X_1 = M$$

$$d(\text{loss})/d(b) = -2*[X_2 - m(X_1) - b] = B$$



# Gradient Descent

2. Update the parameters:

$$m'' = m - L(M) \text{ and}$$

$$b'' = b - L(B)$$

where  $L$  = learning rate

3. Repeat the above steps for several iterations, and adjust the parameters based on the gradients and the learning rate.

After several iterations, the equation can be represented as:

$$X_2 = mf(X_1) + bf$$

Where  $mf$  and  $bf$  are the optimal slopes and intercepts obtained using GD

# Stochastic Gradient Descent

Stochastic gradient descent (SGD) updates parameters by evaluating the loss and gradient on mini-batches of data. It enables efficient iterative optimization in deep learning.



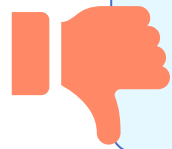
SGD, a variant of the gradient descent algorithm, accelerates model learning in deep learning.



# Stochastic Gradient Descent

The term **stochastic** alludes to the random nature of the algorithm.

SGD randomly picks one sample from a dataset to approximate the loss and the gradient and updates the parameters.



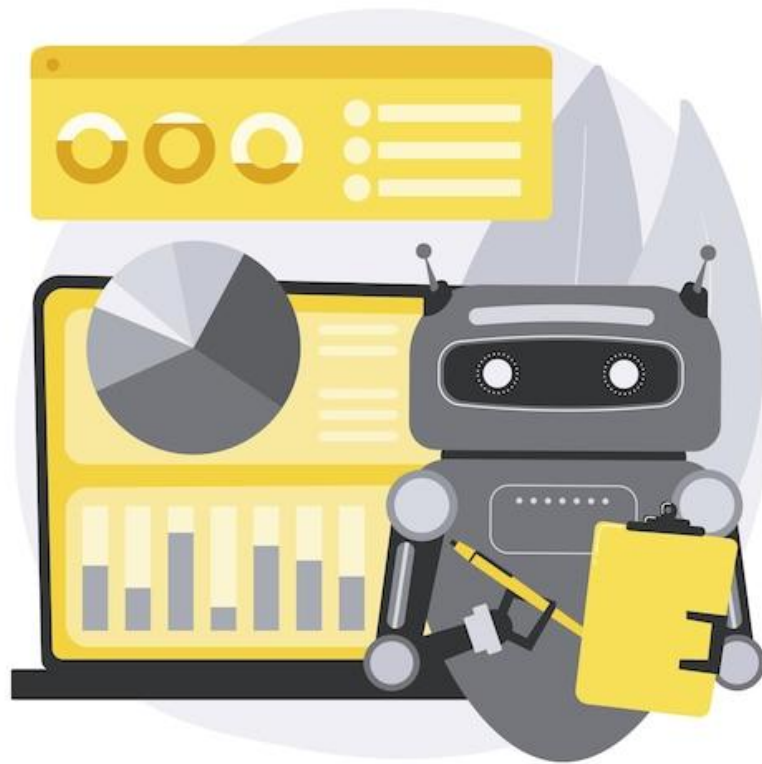
The resulting **descent** can be noisy and might take more iterations.



It is computationally less expensive than vanilla gradient descent.

# Stochastic Gradient Descent

Optimal convergence in SGD can be achieved by systematically decreasing the learning rate over time, allowing the algorithm to settle closer to the global minimum.



- The global minimum is the point in the parameter space where the cost function attains its lowest possible value.
- At this point, the model parameters (weights and biases) are optimized so that the cost function is minimized globally across the entire parameter space.

A common practice is to set the learning rate value to 0.01.

# Key Features of Stochastic Gradient Descent

Deep learning optimization uses mini-batches of data to find optimized weights.

SGD shuffles the data points within each mini-batch for improved generalization.

The optimizer aims to iteratively update the weights to find the optimal solution, considering factors like mini-batch randomness and noise in the data.

# Advantages of Stochastic Gradient Descent

The descent need not be redone when a new datapoint is introduced.

The algorithm can continue from the latest equation and update the parameters to suit the new data.

While SGD is highly random and liable to noise, it is a highly efficient algorithm and can provide optimal solutions under the right conditions.

# Stochastic Gradient Descent-Mini Batch (SGD-Mini Batch)

It is a combination of vanilla GD and SGD, which distribute the training data in its entirety in mini-batches.

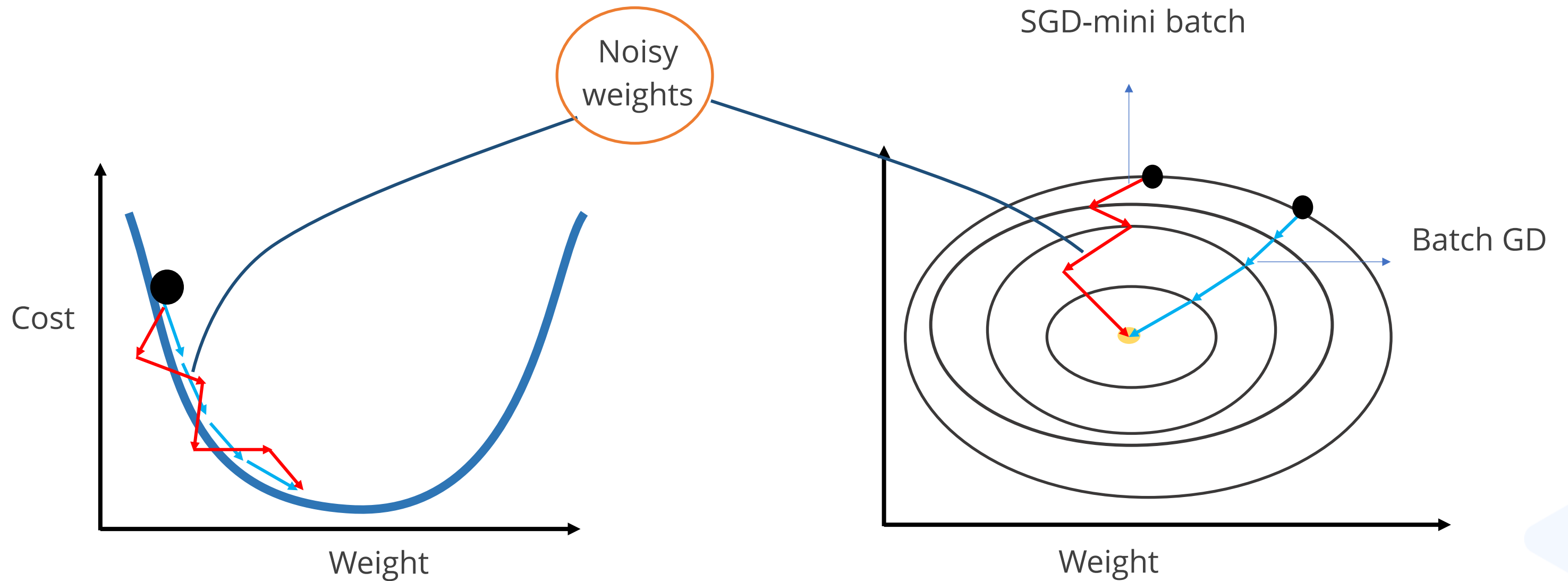
It divides the training data into small batches so that the network can easily be trained on the data.

The mathematical formulation is the same as vanilla GD, but the training occurs batch-wise.

For example, a training set of 400 examples can be divided into 10 batches of 40 training examples each. In this example, the weight evaluation equation will be iterated 10 times the number of batches.

## GD vs. SGD-Mini Batch

GD is computationally expensive, but it converges to the global minimum smoothly. In contrast, SGD creates more noisy weight, which, in turn, takes more time to reach the global minimum.



# Learning Rate of GD, SGD, and SGD-Mini Batch

The learning rate in weight initialization using the optimizers GD, SGD, and SGD-mini batch remains constant. This is a drawback, as it hinders model optimization and performance.

$$W_{new} = W_{old} - \alpha * J(W_{old})$$

Constant  
learning rate

## Assisted Practices



Let's understand the concept of SGD using Jupyter Notebooks.

- 7.03\_Implementation of SGD

**Note:** Please refer to the **Reference Material** section to download the notebook files corresponding to each mentioned topic.

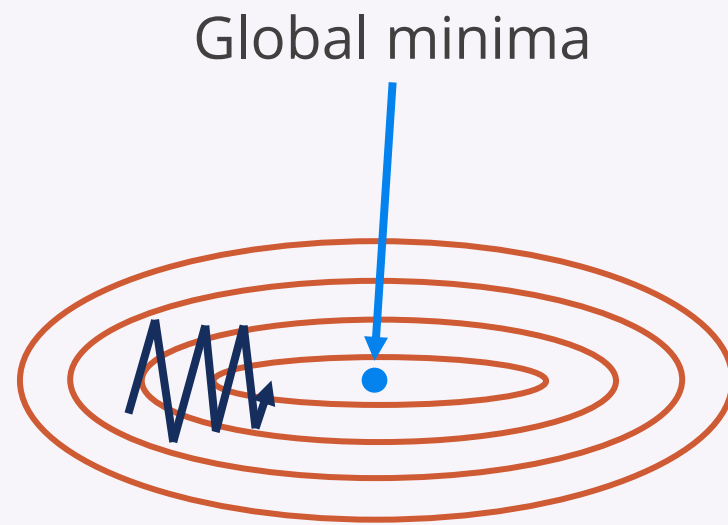




# Introduction to Momentum

# Momentum

It improves convergence and stability. It helps overcome local optima by introducing inertia and consistent direction during weight updates.



At iteration  $t$ ,

$$W_t = W_{t-1} - L * \frac{d(loss)}{d(w_{t-1})}$$

Where,

$W_t$  is the updated value of the parameter.

$W_{t-1}$  is the previous value.

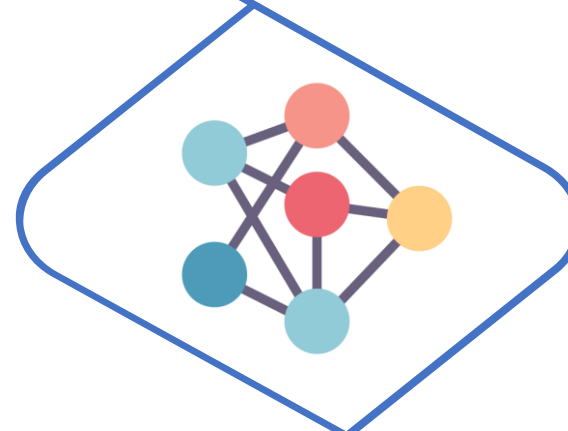
$L$  is the learning rate (typically 0.01).

During model training, algorithm parameters are initialized randomly and updated iteratively, progressively getting closer to the optimal value of the function.

# Momentum

The momentum algorithm uses an assigned learning rate to accelerate convergence and reduce divergence.

Since the learning rate value cannot be large, the process slows down.

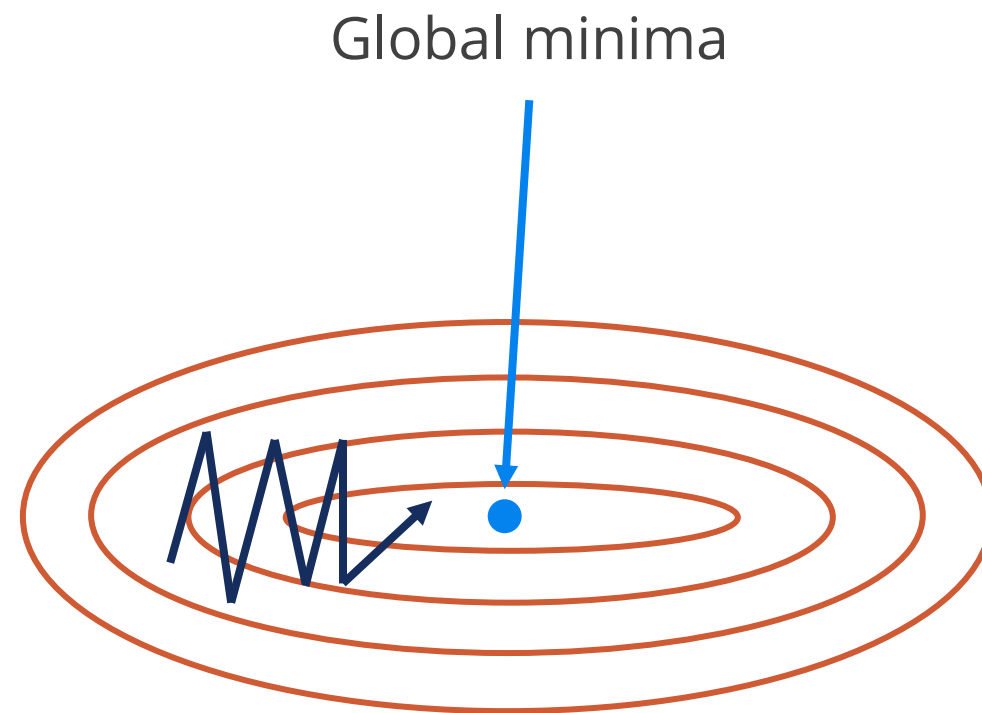


If the algorithm encounters a plateau, it may be deceived into thinking it has reached the minimum.

To resolve these issues and correct the learning rate effect, the concept of momentum is used.

# Momentum

Although it can easily handle smaller datasets, momentum is typically used in the vast, noisy datasets of neural networks.



The only disadvantage to using momentum is that it adds further complexity to the algorithm.

# Momentum

The objective is to reach the global minimum swiftly.

At iteration  $t$ ,

$$M_t = PM_{t-1} + (1 - P) \frac{d(loss)}{d(w_{t-1})}$$
$$w_t = w_{t-1} - LM_t$$

Where,

$M_t$ : Updated value of momentum

$M_{t-1}$ : Previous value

$P$ : Momentum hyperparameter (typically between 0.5 and 0.9)

$W_t$ : Updated value of the parameter

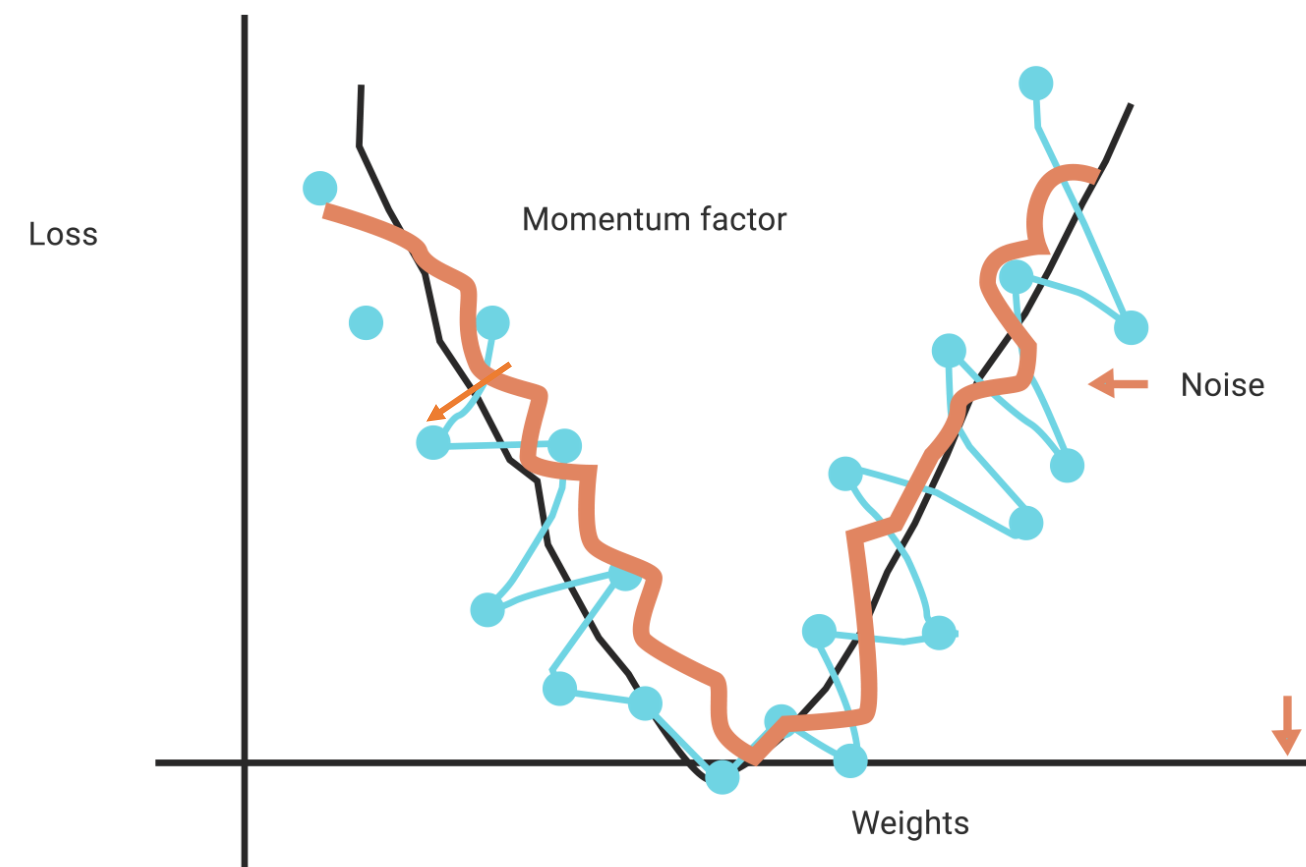
$W_{t-1}$ : Previous value

$L$ : Learning rate (typically 0.01)

The descent must be steep initially and slow down after a certain point to avoid overshooting the minima.

# SGD with Momentum

It is an advanced optimization algorithm that uses a moving average to update the trainable parameters.



SGD with momentum is well-suited to overcoming the noisy weights of SGD.

# SGD with Momentum

- A moving average is a calculation to analyze data points by creating a series of averages of different subsets of the full dataset.
- For example, for times  $t_1$ ,  $t_2$ , and  $t_3$ , the corresponding data points are  $b_1$ ,  $b_2$ , and  $b_3$ .

According to the moving average:

Variable  
at time  $t_1$

$$V_1 = b_1$$

$$V_2 = \gamma * V_1 + b_2$$

$$V_3 = \gamma * V_2 + b_3$$

$$V_3 = \gamma 2b_1 + \gamma b_2 + b_3$$

for  $\gamma = 0.5$ ,

$$V_2 = 0.5b_1 + b_2$$

$$V_3 = 0.25b_1 + 0.5b_2 + b_3$$

Here,  
 $\gamma$ : Decay factor

# SGD with Momentum

- Apply a moving average to the weighted evaluation formula.
- In most cases, gamma = 0.9 works well.

$$J(W_{old}) = \frac{\partial L}{\partial W_{old}}$$

$$W_{new} = W_{old} - \alpha * J(W_{old})$$
$$= W_{old} - [\gamma V_{t-1} + \alpha * J(W)]$$

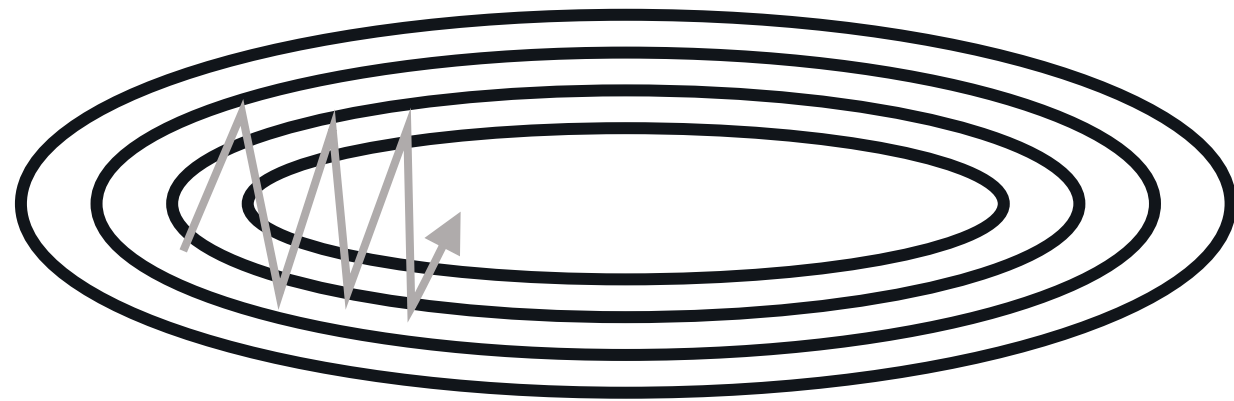
$$V_{t-1} = \gamma * J(W)_t + \gamma^2 * J(W)_{t-1} + \gamma^3 * J(W)_{t-2} \dots$$

Momentum

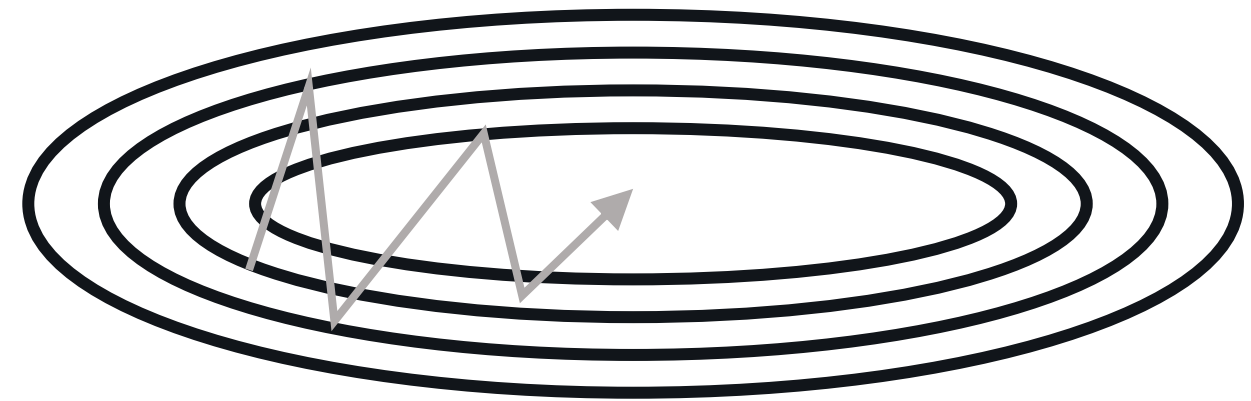


# SGD With and Without Momentum

SGD with momentum clearly shows that momentum makes the steps smooth and less noisy.



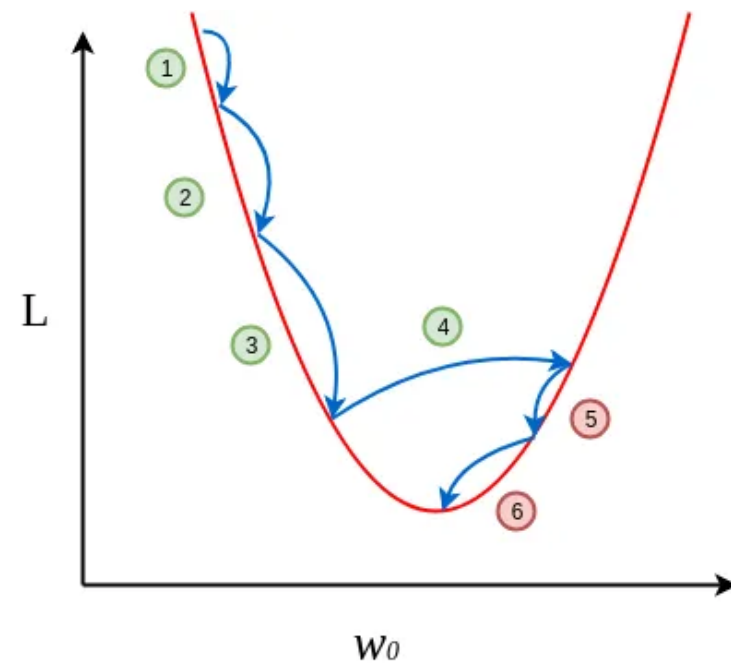
SGD with momentum



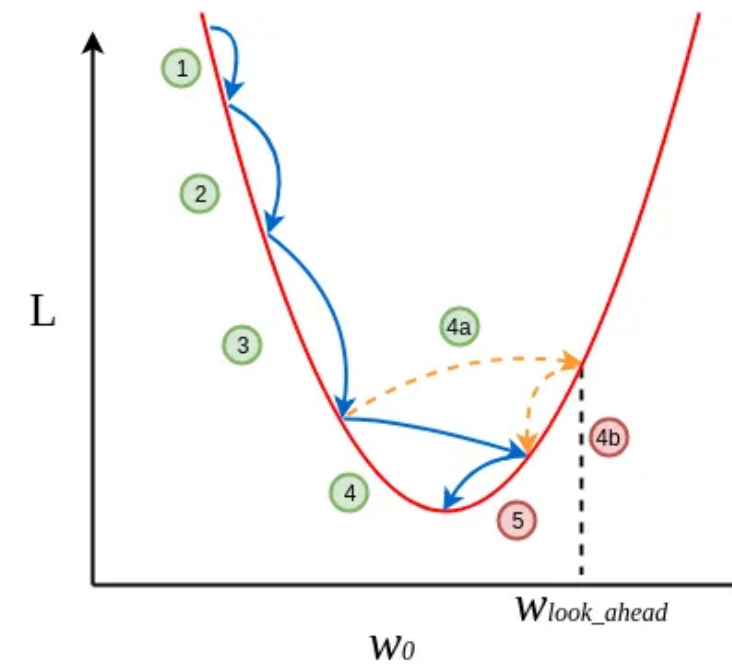
SGD without momentum

# Nesterov Accelerated Gradient (NAG)

NAG combines momentum-based updates with lookahead adjustments for faster convergence compared to traditional SGD with or without momentum.



(a) Momentum-Based Gradient Descent



(b) Nesterov Accelerated Gradient Descent

$$\text{Green Circle} \Rightarrow \frac{\partial L}{\partial w_0} = \frac{\text{Negative}(-)}{\text{Positive}(+)}$$

$$\text{Red Circle} \Rightarrow \frac{\partial L}{\partial w_0} = \frac{\text{Negative}(-)}{\text{Negative}(-)}$$

## NAG vs. SGD with Momentum

The characteristics of NAG and SGD with the momentum method may resemble each other or differ significantly under various scenarios.

- Both methods give a distinct output when the learning rate  $\eta$  is reasonably large.
- In such a case, NAG allows a larger decay rate  $\alpha$  than SGD with the momentum method while preventing oscillations.
- The theorem also shows that SGD with momentum and NAG become equal when  $\eta$  is small.

## Assisted Practices



Let's understand the concept of momentum using Jupyter Notebooks.

- 7.05\_Implementation of Momentum

**Note:** Please refer to the **Reference Material** section to download the notebook files corresponding to each mentioned topic.



# Introduction to AdaGrad

# Adaptive Gradient (AdaGrad)

Adagrad is an optimization algorithm that adjusts the learning rate for each parameter based on historical gradients, improving convergence and performance.

- Accumulates gradients over time from the entire training dataset
- Scales the learning rate individually for each parameter based on their historical gradients
- Is effective for handling sparse data with varying parameter importance

# AdaGrad

Adaptive gradient (AdaGrad) optimization iteratively updates different learning rates for each parameter without manual tuning.

At iteration  $t$ ,

$$G_t = \sum_{i=1}^t \left[ \frac{d(\text{loss})}{d(w_i)} \right]^2$$

$$L_t = \frac{L_{t-1}}{\sqrt{G_t + E}}$$

$$w_t = w_{t-1} - L_t \left( \frac{d(\text{loss})}{d(w_{t-1})} \right)$$

Where,

$G_t$ : Sum of the squares of gradients over time

$L_t$ : Updated value of the learning rate

$E$ : Small positive number

$w_t$ : Updated value of the parameter

# AdaGrad

AdaGrad updates the learning rate in each iteration to reach the global minimum swiftly.

It calculates the custom learning rate (step size) for one parameter as follows:

$$\frac{\text{Step size}}{1 \times 10^{-8} + \sqrt{s(t)}} = \text{cust step size}(t + 1)$$

Where,

cust step size(t+1) is the estimated step size for an input variable at a certain point.

Step size is the beginning step size.

$1 \times 10^{-8}$  is the small constant value.

$s(t)$  is the sum of the squared partial derivatives for the input variable.



# AdaGrad

The AdaGrad optimizer is commonly used in machine learning and deep learning algorithms to adaptively adjust the learning rate during the training process.

The mathematical formulation for the AdaGrad optimizer is:

$$W_t = W_{t-1} - \alpha * t * J(W_{t-1})$$

$$\alpha * t = \frac{\alpha}{\sqrt{\beta t + \epsilon}}$$

$$\beta t = \sum_{i=1}^t J(W_i)^2$$

Initial  
learning rate

# AdaGrad

As the epochs increase, the learning rates decrease, which results in a gradual decrease in weight.

$\epsilon$  is considered in the formula because if  $\beta t$  becomes zero, then the weight will be constant.

$$\alpha * t = \frac{\alpha}{\sqrt{\beta t + \epsilon}}$$

Small nonzero  
value

The only drawback with AdaGrad is that the  $\beta t$  value can sometimes turn out to be large.

$$\beta t = \sum_{i=1}^t J(W_i)^2$$

# Advantages of AdaGrad

01

It eliminates the need to manually tune the learning rate.

02

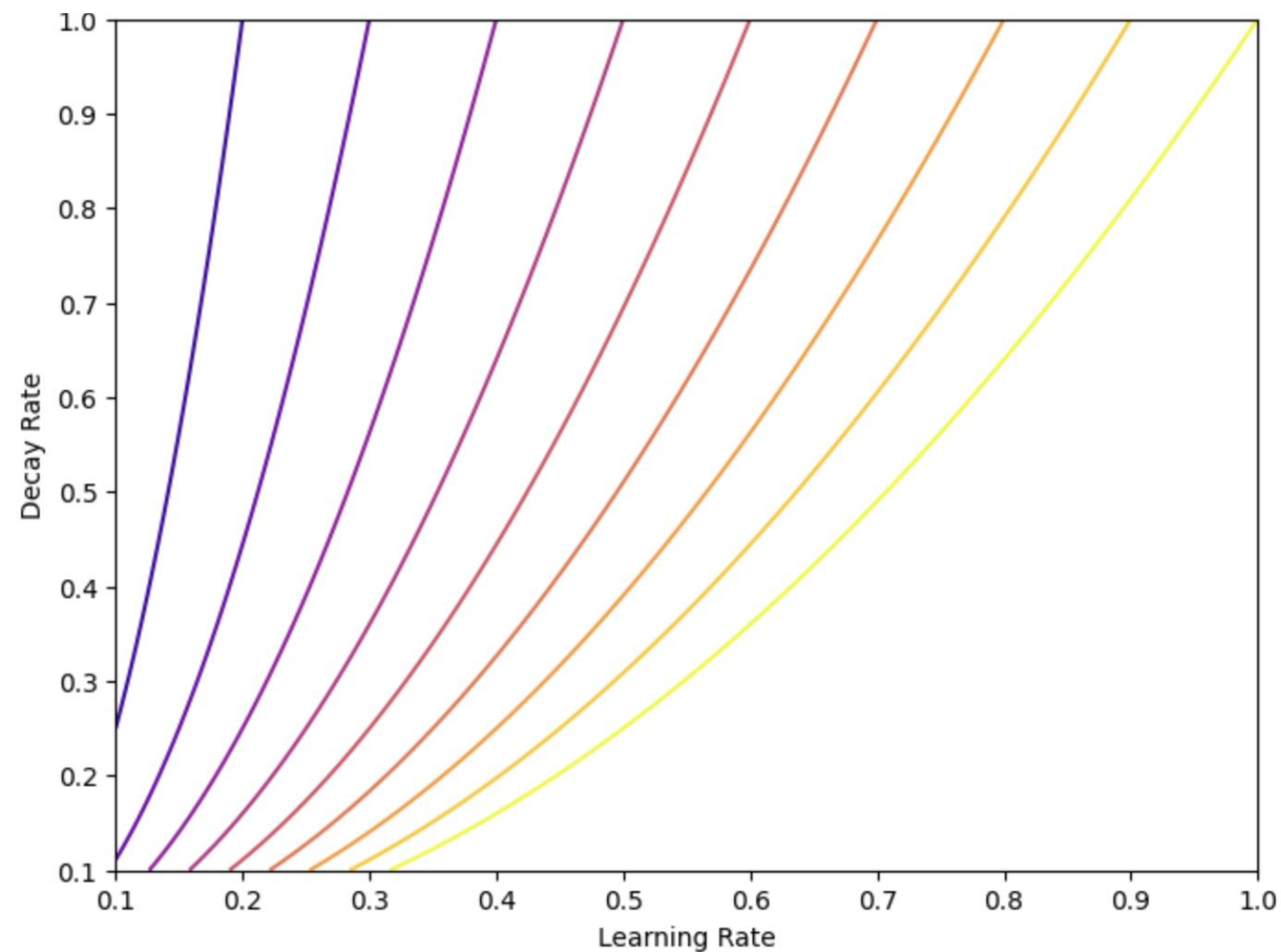
It facilitates faster convergence and is more reliable.

03

It can effectively adapt the learning process regardless of the initial learning rate.

# Disadvantages of AdaGrad

The main disadvantage is that the learning rate value might decrease rapidly, leading to slower convergence.



To resolve this, a decay factor is added to AdaGrad using root mean square propagation to maintain the sum of squares in the denominator at a manageable size.

## Assisted Practices



Let's understand the concept of AdaGrad using Jupyter Notebooks.

- 7.07\_Implementation of AdaGrad

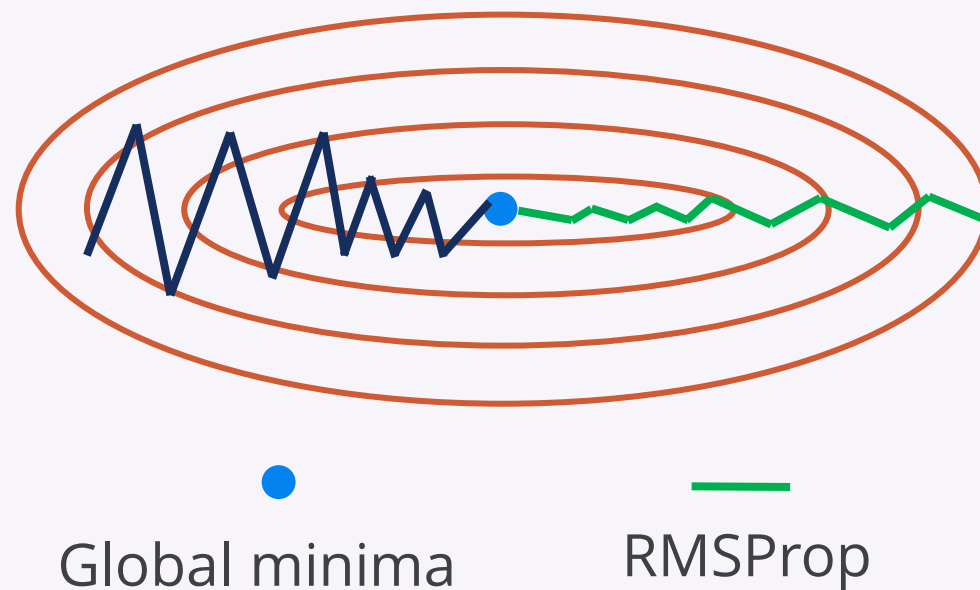
**Note:** Please refer to the **Reference Material** section to download the notebook files corresponding to each mentioned topic.



# Introduction to RMSProp

# What Is RMSProp?

The Root Mean Square Propagation (RMSProp) optimizer is a momentum-based version of the gradient descent technique.



It limits the oscillations in the vertical plane.

It boosts the learning rate, allowing the algorithm to take greater horizontal steps to converge faster.

RMSProp and gradient descent differ in the way the gradients are calculated.

## RMSProp: Equation

The RMSProp and gradient descent with momentum are determined in the following way:

$$v_{dw} = \beta * v_{dw} + (1 - \beta) * dw$$

$$v_{db} = \beta * v_{db} + (1 - \beta) * db$$

$$W = W - \alpha * v_{dw}$$

$$b = b - \alpha * v_{db}$$

Here,

$v_{dw}$ : Velocity for weights

$v_{db}$ : Velocity for biases

$\beta$ : Momentum coefficient

$dw$ : Current gradient of weights

$db$ : Current gradient of biases

$W$ : Weights

$b$ : Biases



## RMSProp: Equation

Beta is the measure of momentum and is commonly set to 0.9.

This normalization equalizes the step size or momentum.

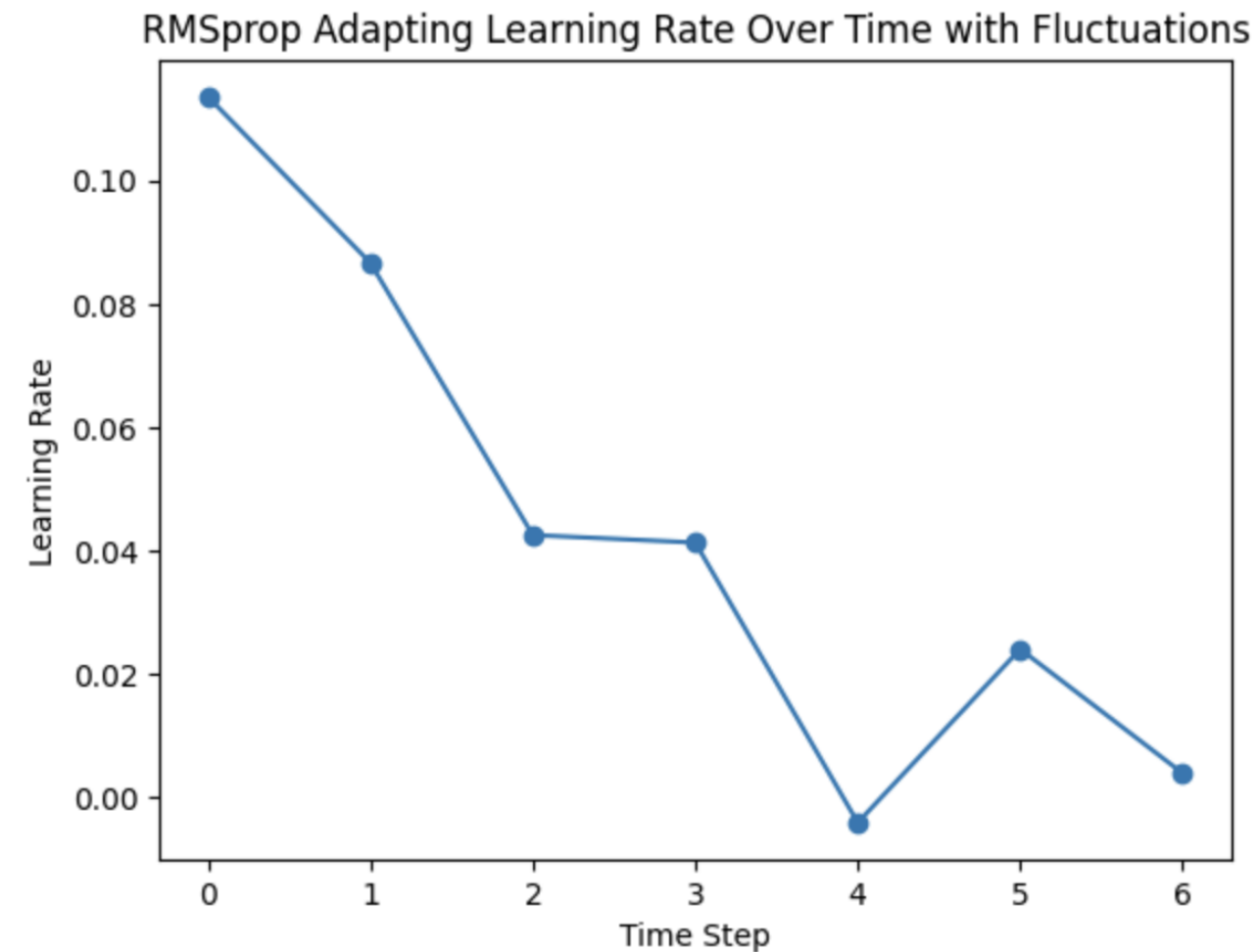


To normalize the gradient, RMSProp uses a moving average of squared gradients.

Normalization lowers the step size for high gradients to avoid exploding and raises it for minor gradients to avoid vanishing.

# Root Mean Square Propagation (RMSProp)

RMSProp treats the learning rate as an adaptive parameter rather than a hyperparameter.



This indicates that the rate of learning fluctuates with time.

# Root Mean Square Propagation (RMSProp)

RMSprop was developed to mitigate the drawbacks of AdaGrad.

The formula for RMSprop is as follows:

$$W_t = W_{t-1} - \alpha * t * J(W_{t-1})$$

$$\alpha * t = \frac{\alpha}{\sqrt{W_{avg(t)} + \varepsilon}}$$

$$W_{avg(t)} = \gamma * W_{avg(t-1)} + (1 - \gamma)J(W)^2$$

Where,

$W_t$ : Updated weights at time step  $t$

$W_{t-1}$ : Previous weights from the previous time step

$\alpha * t$ : Adaptive learning rate at time step  $t$

$\alpha$ : Initial learning rate (constant)

$t$ : Current time step

$W_{avg}(t)$ : Accumulated squared gradients at time step  $t$

$\gamma$ : Decay factor

$J(W)$ : Cost function

$\varepsilon$ : Constant

# RMSProp: Syntax

The Python code for implementing RMSprop is as follows:

Syntax:-

```
def RMSprop(index, beta, db, dw, vdw, vdb, alpha):  
    vdw = beta*vdw + (1 - beta)*dw**2  
    Vdb =(beta*vdb + (1-beta)*db**2  
  
    Model.layers[index].weight-=(alpha/(np.sqrt(vdw)+1e-8))*dw  
    Model.layers[index].bias-=(alpha/(np.sqrt(vdb)+1e-8))*db
```

Here, epsilon  $\epsilon = e-8$  is added for weight and bias to maintain numerical stability.

# Assisted Practices



Let's understand the concept of RMSProp using Jupyter Notebooks.

- 7.09\_Implementation of RMSProp

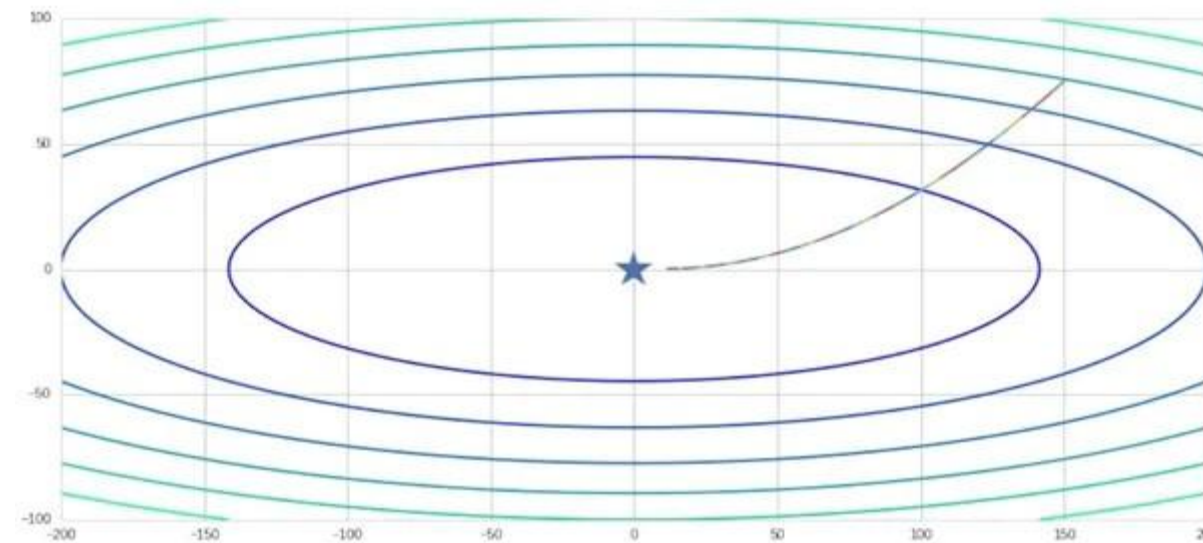
**Note:** Please refer to the **Reference Material** section to download the notebook files corresponding to each mentioned topic.



# Introduction to Adadelta

# Adadelta

Adadelta is an optimization algorithm used to train neural networks. It builds on AdaGrad and RMSProp, altering the custom step size calculation.



This approach removes the need for an initial learning rate hyperparameter.

# Adadelta: Selecting the Right Algorithm

To achieve optimal performance, choose the algorithm that aligns with the behavior of the data variables in the specific application.

## Understand the data

- Analyze the data variables' behavior and characteristics
- Identify challenges such as sparsity or noisy gradients

## Choose the appropriate algorithm

- Recognize the strengths and weaknesses of various algorithms
- Use adaptive learning rate methods without manual tuning



# Adadelata: Selecting the Right Algorithm

## Align with application

- Assess the alignment of Adadelata's features with data behavior
- Ensure the compatibility of the algorithm with the specific task

## Optimize performance

- Fine-tune hyperparameters, including the decay factor
- Evaluate and adjust continuously to align with the data variables

## Adadelta: Equation

Consider the parameter  $\rho$ , where the leaky updates are stored in the state variable as below:

The given equation represents how Adadelta updates the Exponentially Weighted Moving Average (EWMA) of the squared gradients, also referred to as **leaky updates**.

$$S_t = \rho S_{t-1} + (1 - \rho)g_t^2$$

Here,

$S_t$ : Current EWMA (Exponentially Weighted Moving Average) of squared gradients

$\rho$ : Decay factor for the contribution of the previous value

$S_{t-1}$ : Previous EWMA of squared gradients

$(1-\rho)$ : Weight for the current squared gradient

$g_t^2$ : Squared gradient at the current step

## Adadelta: Equation

The algorithm performs the update with the rescaled gradient:

$$X_t = X_{t-1} - g'_t$$

where

$$g'_t = \frac{\sqrt{\Delta X_{t-1} + \epsilon}}{\sqrt{s_t + \epsilon}} g_t$$

Here,  $\Delta X_{t-1}$  = Leaky average of the squared rescaled gradients

$g'_t$  = The rescaled gradient at time  $t$

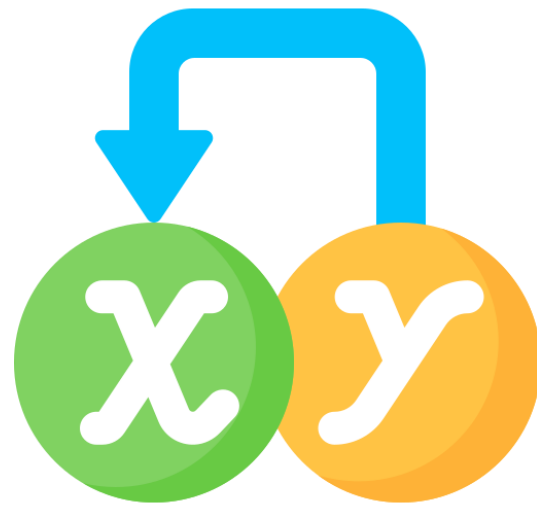
$\Delta X_0 = 0$  is updated at each step with the rescaled gradient as shown:

$$\Delta X_t = \rho \Delta X_{t-1} + (1 - \rho) g'^2_t$$

$\epsilon$  is in the range of  $10^{-5}$  to maintain numerical stability.

# Adadelta: Equation

In the algorithm, Adadelta uses two state variables to adaptively adjust the learning rate for each parameter:



$\mathbf{s}_t$  to store a leaky average of the second moment of the gradient

$\Delta \mathbf{x}_t$  to store a leaky average of the second moment of the change of parameters in the model

The second moment of a set of values is a statistical measure that represents the average of the squared deviations from the mean.

# Adadelta: Syntax

It can be used with the Keras package in the following way:

Syntax:-

```
tf.keras.optimizers.Adadelta(  
    Learning_rate=0.001, rho=0.95, epsilon=1e-07, name="Adadelta", **kwargs  
)
```

The learning rate can be set to higher values for better performance.

$\rho$  is the decay rate and can be a tensor or a floating-point value.

$\epsilon$  maintains numerical stability to offset cases like dividing by zero and rounding off issues.

# Summary of Adadelta

Adadelta requires two state variables to store the second moment of the gradient and the change.

It has no learning rate parameter.



It uses leaky averages to estimate the appropriate statistics.

## Assisted Practices



Let's understand the concept of Adadelta using Jupyter Notebooks.

- 7.11\_Implementation of Adadelta

**Note:** Please refer to the **Reference Material** section to download the notebook files corresponding to each mentioned topic.



# Adam Optimizer



# What Is Adam Optimizer?

The Adam optimizer is a popular algorithm in deep learning that optimizes stochastic objectives using adaptive estimates. It efficiently handles sparse gradients and noisy problems with low resources.

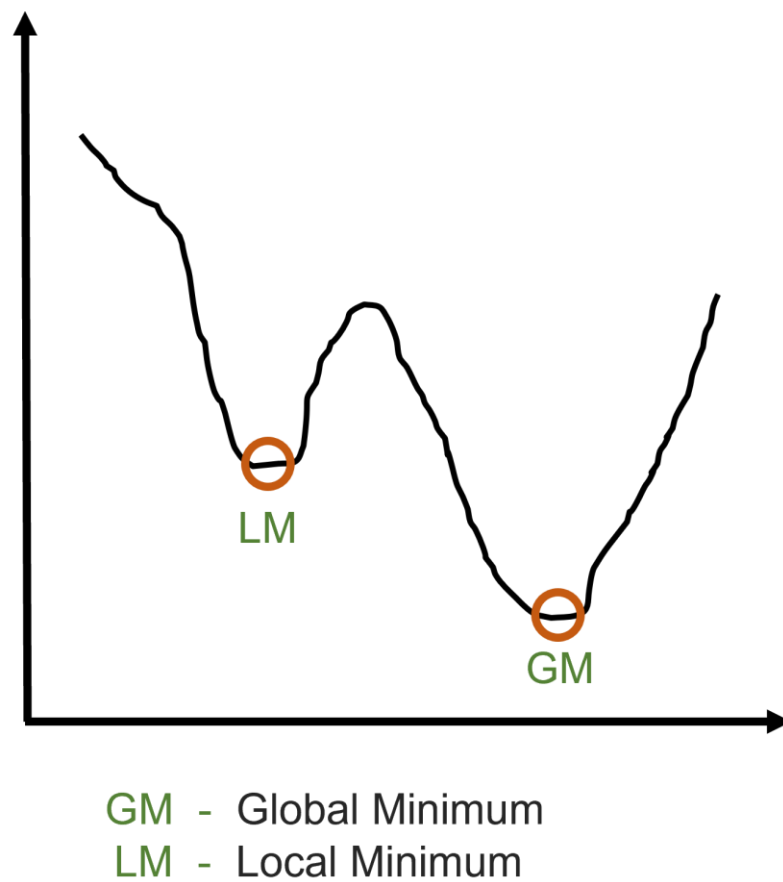
It effectively handles large-scale problems with extensive data or numerous parameters.

It doesn't require excessive memory, but its use mainly depends on the model architecture.

Adam combines momentum and RMSProp concepts for effective parameter space navigation and fast convergence.

# Adam Optimizer

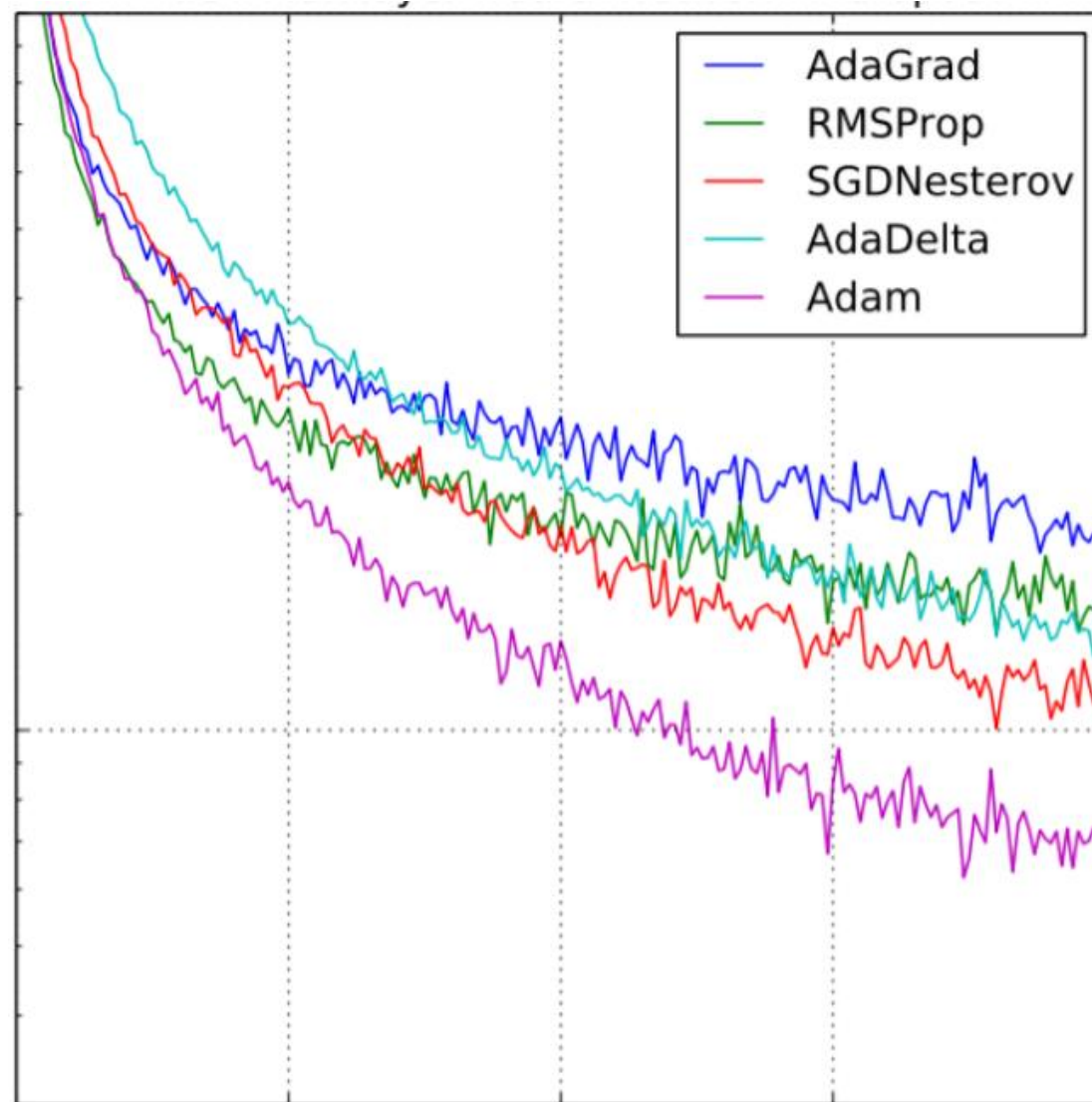
With Adam, limit the pace of gradient descent to minimize oscillation.



- It takes large enough steps (step size) to avoid local minimum barriers.
- As a result, it integrates the qualities of the preceding strategies to efficiently obtain the global minimum.
- It outperforms the optimization algorithm by a large margin to provide an optimized gradient descent.

# Adam Optimizer

Comparing the Adam optimizer with other optimizers for the MNIST dataset in the following graph:



Adam performed more efficiently than other optimizers.

# Adam Optimizer

It keeps an exponentially decaying average of the past gradient  $\mathbf{m}_t$ .

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2$$

$\mathbf{m}_t$ : Momentum of gradients

$\beta_1$ : Decay rate for momentum

$\mathbf{g}_t$ : Current gradient value

$\mathbf{v}_t$ : Moving average of squared gradients

$\beta_2$ : Decay rate for  $\mathbf{v}_t$

# Adam Optimizer

$m_t$  and  $v_t$  are the estimates of the first and second moments of the gradients, respectively.

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

The updated parameters are used like in Adadelta and RMSProp.  
This yields the following Adam update rule:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t.$$

# Adam Optimizer

The following are the parameters used in the Adam optimizer:

To prevent zero division, a tiny positive constant,  $\epsilon$  ( $10^{-8}$ ), is added.

$\beta_1$  = The decay rate for the momentum term

$\beta_2$  = The decay rate for the squared gradients

( $\beta_1 = 0.9$  and  $\beta_2 = 0.999$  recommended)

$\eta$  = learning rate (0.001)

# Adam Optimizer

The Adam optimization algorithm utilizes the following formulas for efficient training in machine learning models:

Update mean (first moment)

$$M_t = \beta_1 * M_{t-1} + (1-\beta_1) * G_t$$

Update variance (second moment)

$$V_t = \beta_2 * V_{t-1} + (1-\beta_2) * (G_t)^2$$

Update weights

$$W_t = W_{t-1} - \alpha * M_t / (\sqrt{V_t} + \epsilon)$$

# Assisted Practices



Let's understand the concept of Adam using Jupyter Notebooks.

- 7.13\_Implementation of Adam

**Note:** Please refer to the **Reference Material** section to download the notebook files corresponding to each mentioned topic.





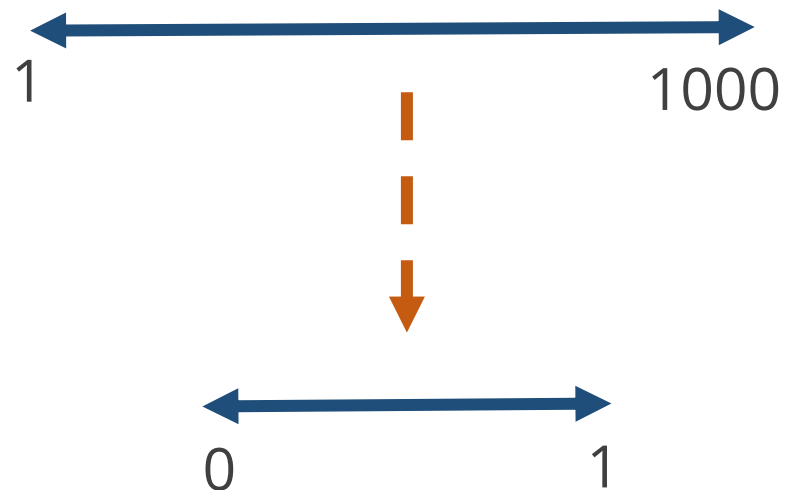
# Batch Normalization

# Data Preprocessing

In preprocessing, the data is generally normalized or standardized.

## Normalization

A typical normalization involves scaling down a large range of data into a smaller range.



## Standardization

A typical standardization is to subtract the mean of all the data points from each data point and then divide the difference by the standard deviation.

$$Z = \frac{X - m}{\sigma}$$

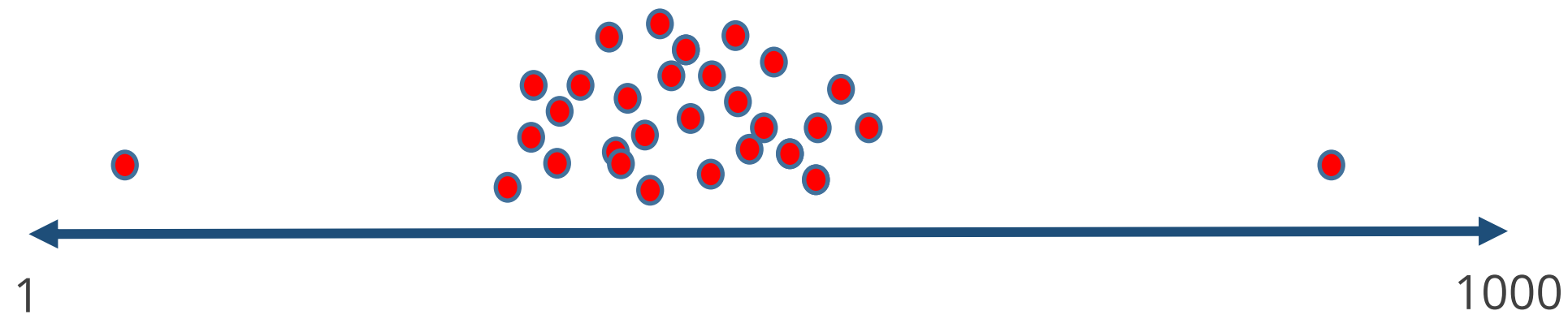
$X$  = Data points

$m$  = Mean

$\sigma$  = Standard Deviation

# Why Data Preprocessing?

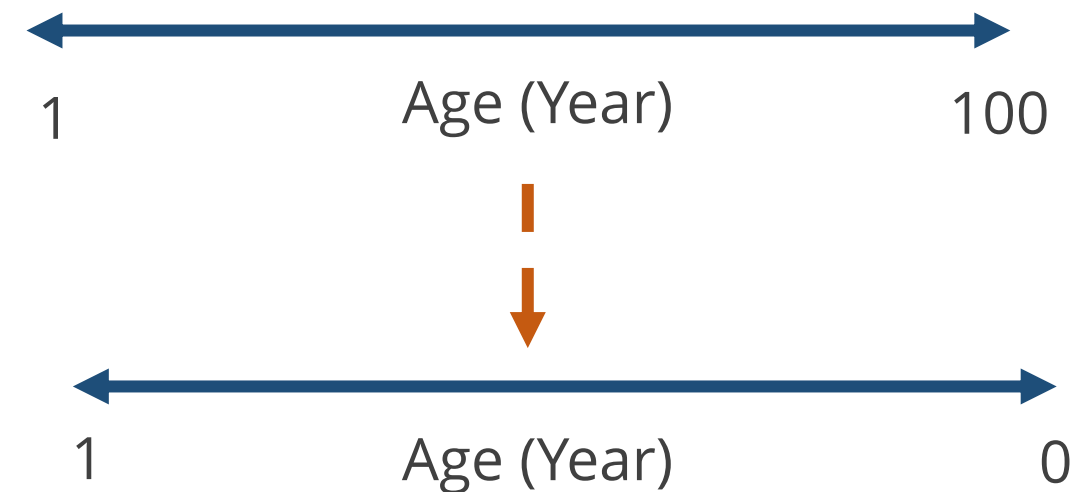
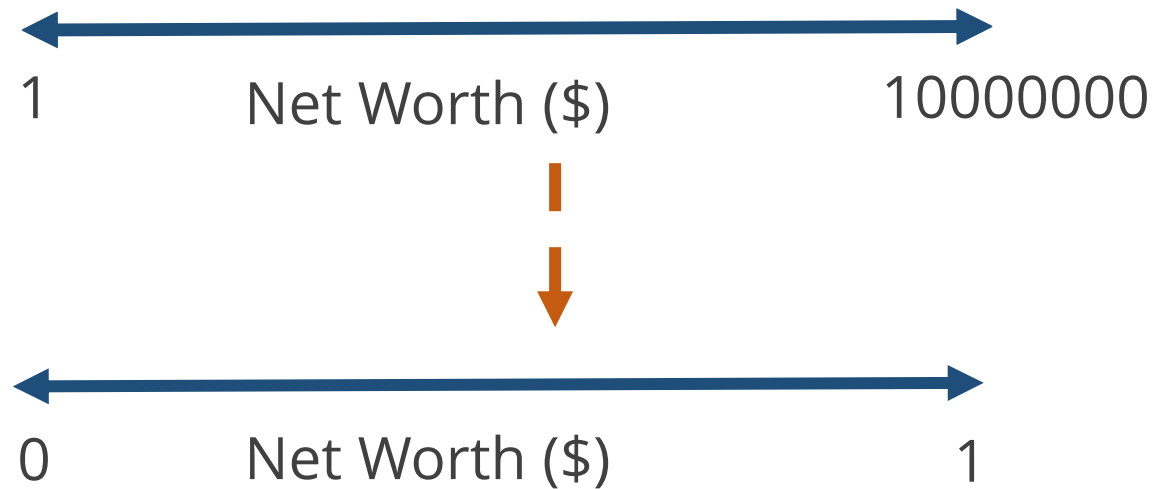
Data points can either be high or low. This leads to cascading effects in the network. Therefore, data preprocessing is needed.



The cascading effect in neural networks refers to a phenomenon where errors or variations in input data propagate through the layers of the network, potentially amplifying and affecting the final output.

# Why Data Preprocessing?

When there are multiple features, each with a different range of data points, the non-processed data creates instability and cascades through the neural network layers. Scaling the different ranges to a standard range leads to stability and better results.



# Normalization Techniques

During the preprocessing procedure, before training a neural network, the input is either normalized or standardized.

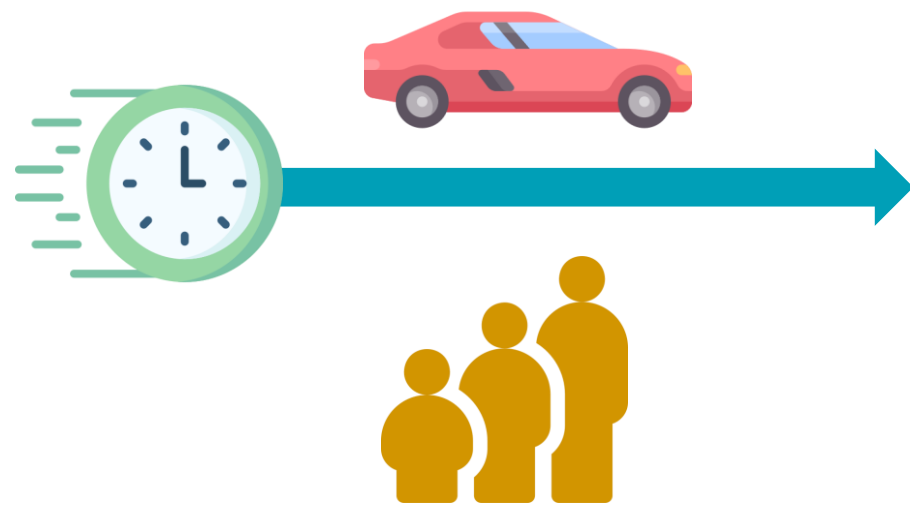


During the data preparation stage, the data is prepared for use in training.

This exercise alters the data so that all data points are on the same scale.

# Normalization Techniques: Example

Each of the characteristics in every sample differs significantly.



## Example

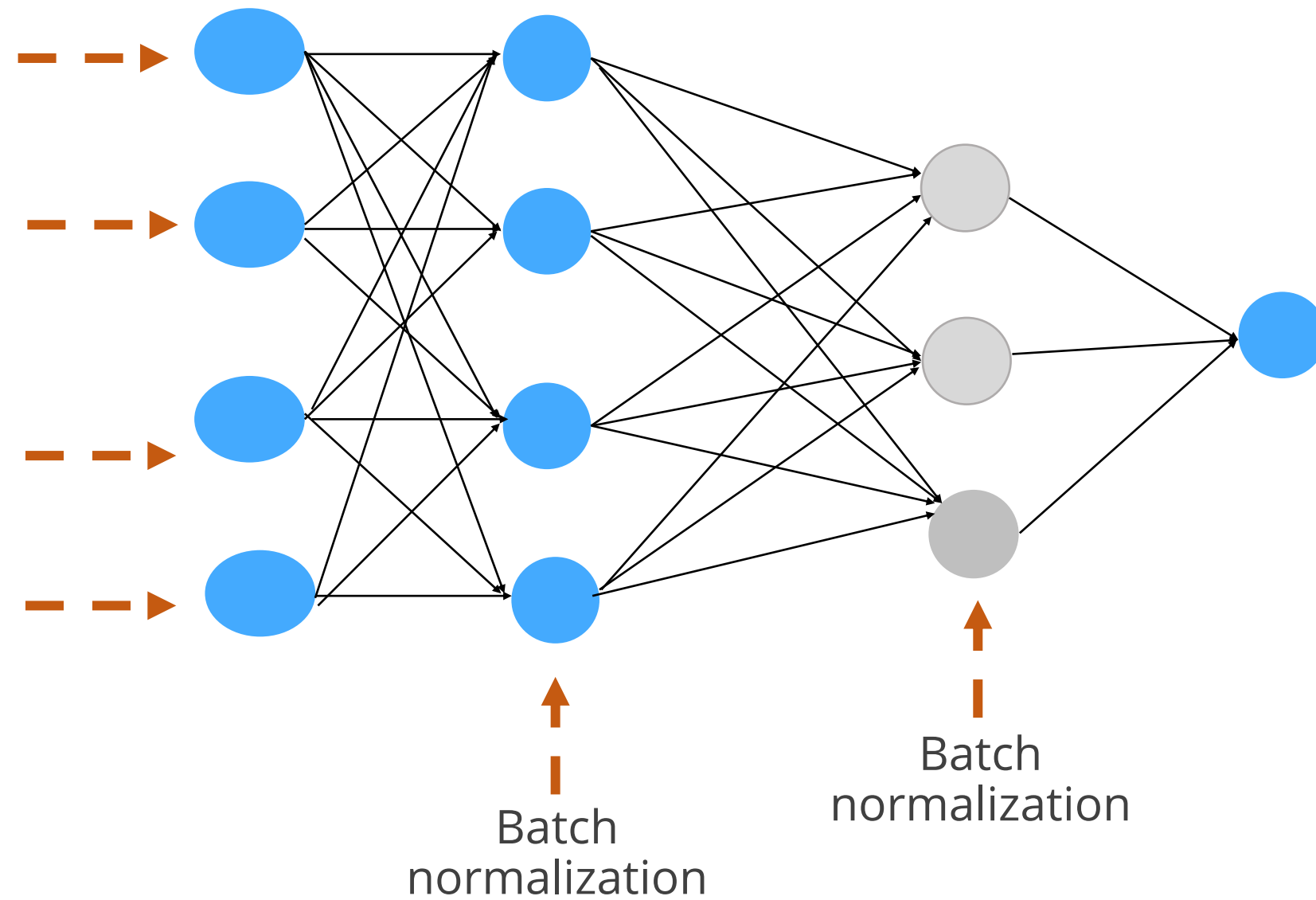
If one characteristic relates to an individual's age and the other to the individual's miles driven in a car in the last five years, these two data points, age and miles traveled, will not be on the same scale.



# **Batch Normalization Implementation**

# Batch Normalization

Normalization of data before feeding it into the network is not enough; the outputs from the neurons should also be normalized.





# Batch Normalization Process

The batch normalization process involves the following:

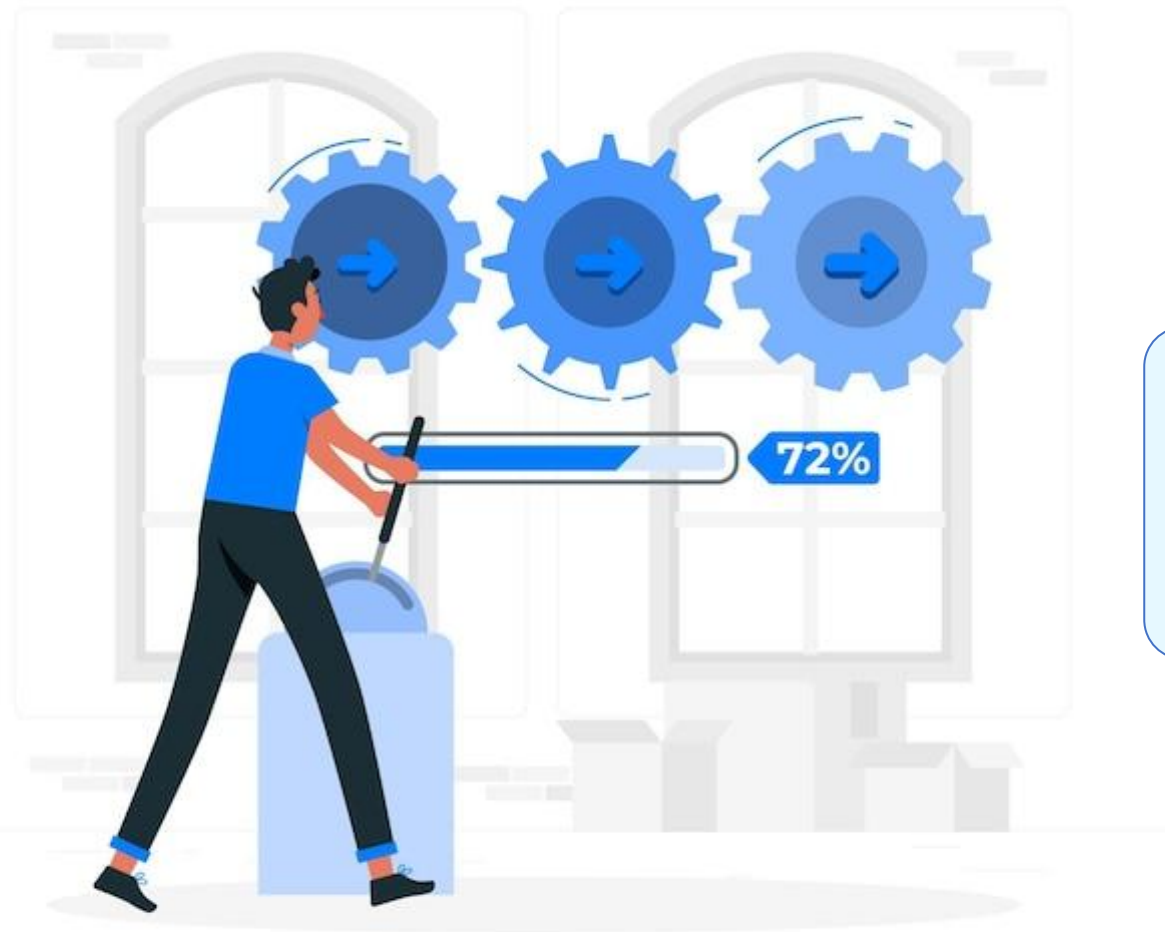
Steps	Expression	Description
1	$Z = \frac{X - m}{\sigma}$	Normalize output x from the activation function
2	$Z * g$	Multiply normalized output z by the arbitrary parameter g
3	$(Z * g) + b$	Add an arbitrary b to the resulting product (z * g)

**m**: Mean of the mini-batch, **s**: Standard deviation (or variance) of the mini-batch

**g**: Scale parameter (or gamma), **b**: Shift parameter (or beta)

# Batch Normalization Process

The data is given a new standard deviation and mean with two arbitrarily chosen, trainable parameters,  $g$  and  $b$ .



Since normalization is included in the gradient process, the weights inside the network do not become imbalanced.

By normalizing the inputs, batch normalization ensures that extreme weights do not dominate the training process. This reduces the risk of instability and overfitting caused by these outliers.

# Batch Normalization Process

It normalizes output data from the model's activation functions for specific layers.



Therefore, the data flowing in and the data within the model are both normalized.

The entire process occurs on a per-batch basis and is thus called the batch norm.

# Implementing Batch Normalization Using Keras

- Batch normalization is an additional imported library.
- Initialize batch normalization after the ReLU activation function.

Syntax:-

```
from keras.models import Sequential
from keras.layers import Dense, Activation, BatchNormalization

model = Sequential([
    Dense(16, input_shape=(1,5), activation='relu'),
    Dense(32, activation='relu'),
    BatchNormalization(axis=1),
    Dense(2, activation='softmax')])
```

# Implementing Batch Normalization Using Keras

The batches are determined by the batch size that is set when the model is trained.

Syntax:-

```
import tensorflow as tf

model = tf.keras.Sequential([
    tf.keras.layers.Dense(16, activation='relu'),
    tf.keras.layers.Dense(32, activation='relu'),
    tf.keras.layers.BatchNormalization(axis=1),
    tf.keras.layers.Dense(2, activation='softmax')])
```

# Components of Batch Normalization

## **Input layer**

Depends on the number of features in the dataset

## **Hidden layer 1 with 16 nodes**

Dense layer

Batch normalization (optional)

ReLU activation

## **Hidden layer 2 with 32 nodes**

Dense layer

Batch normalization (optional)

ReLU activation

## **Output layer with 2 nodes**

Dense layer

Softmax activation

# Components of Batch Normalization

There is a batch normalization layer between the last hidden layer and the output layer.

```
tf.keras.layers.BatchNormalization(axis=1)
```

In Keras, batch normalization is specified using the **BatchNormalization** layer, typically placed after the layer whose activation output needs normalization.

# Components of Batch Normalization

The axis parameter specifies the axis of the data to be normalized and is often the feature axis.

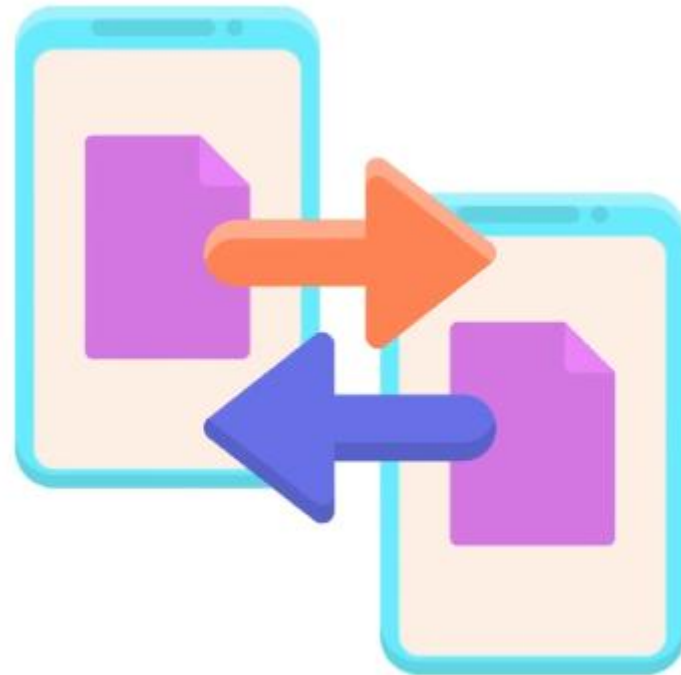
Many other parameters can also be specified and adjusted, including `beta_initializer` and `gamma_initializer`.

The `beta_initializer` and `gamma_initializer` parameters are commonly used in neural network architectures, especially in normalization layers like batch normalization.



## Applying Batch Norm to a Layer

When a batch norm is applied to a layer, it first normalizes the output from the activation function that goes as input to the next layer.



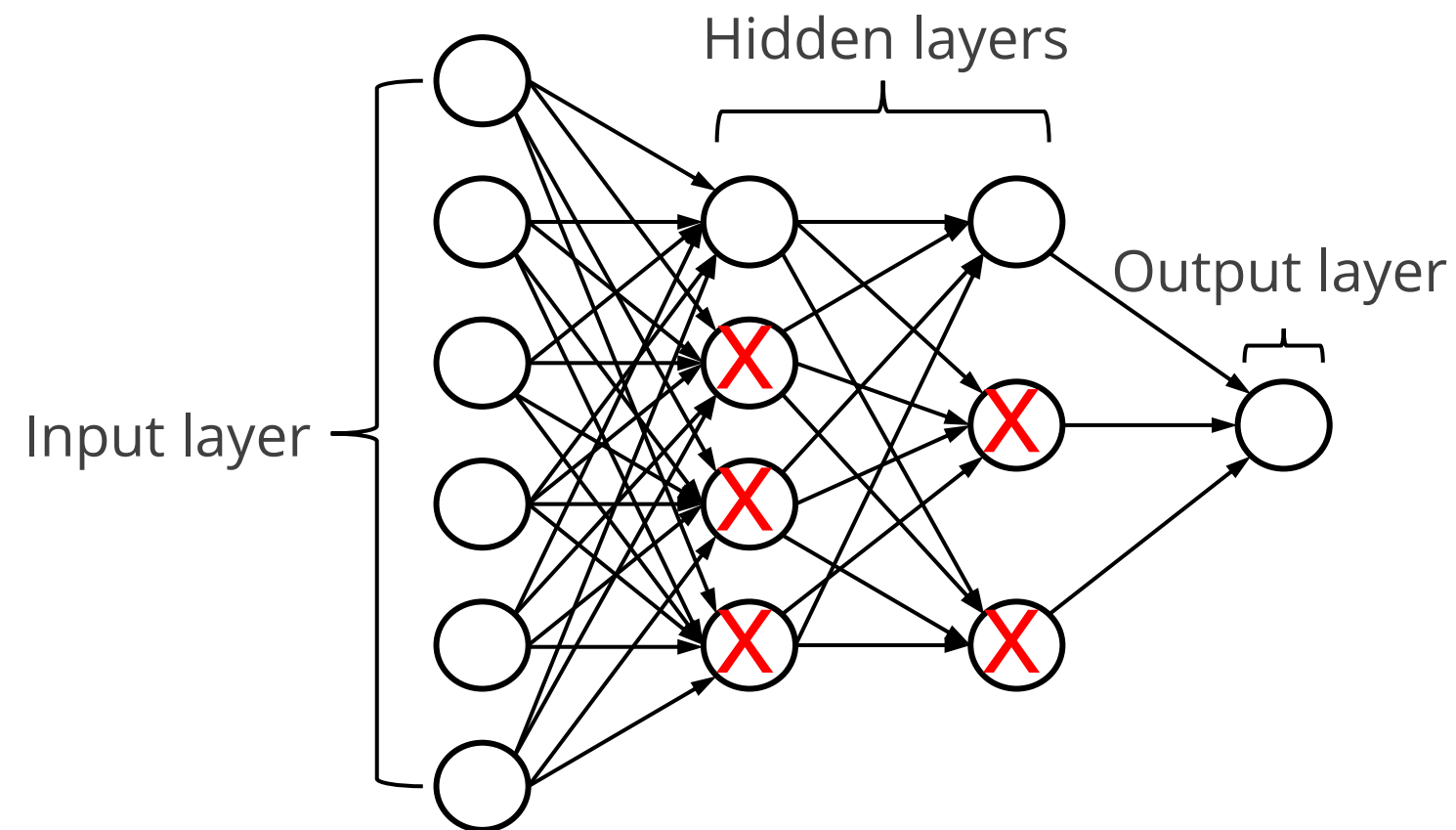
The purpose of batch normalization is to normalize this output, making it more standardized and easier to train.



# Regularization

# Regularization

Regularization is a technique that makes slight modifications to the learning algorithm so that the model generalizes unseen data more effectively.



Regularization helps reduce errors by fitting a function appropriately to the training set and avoiding overfitting.

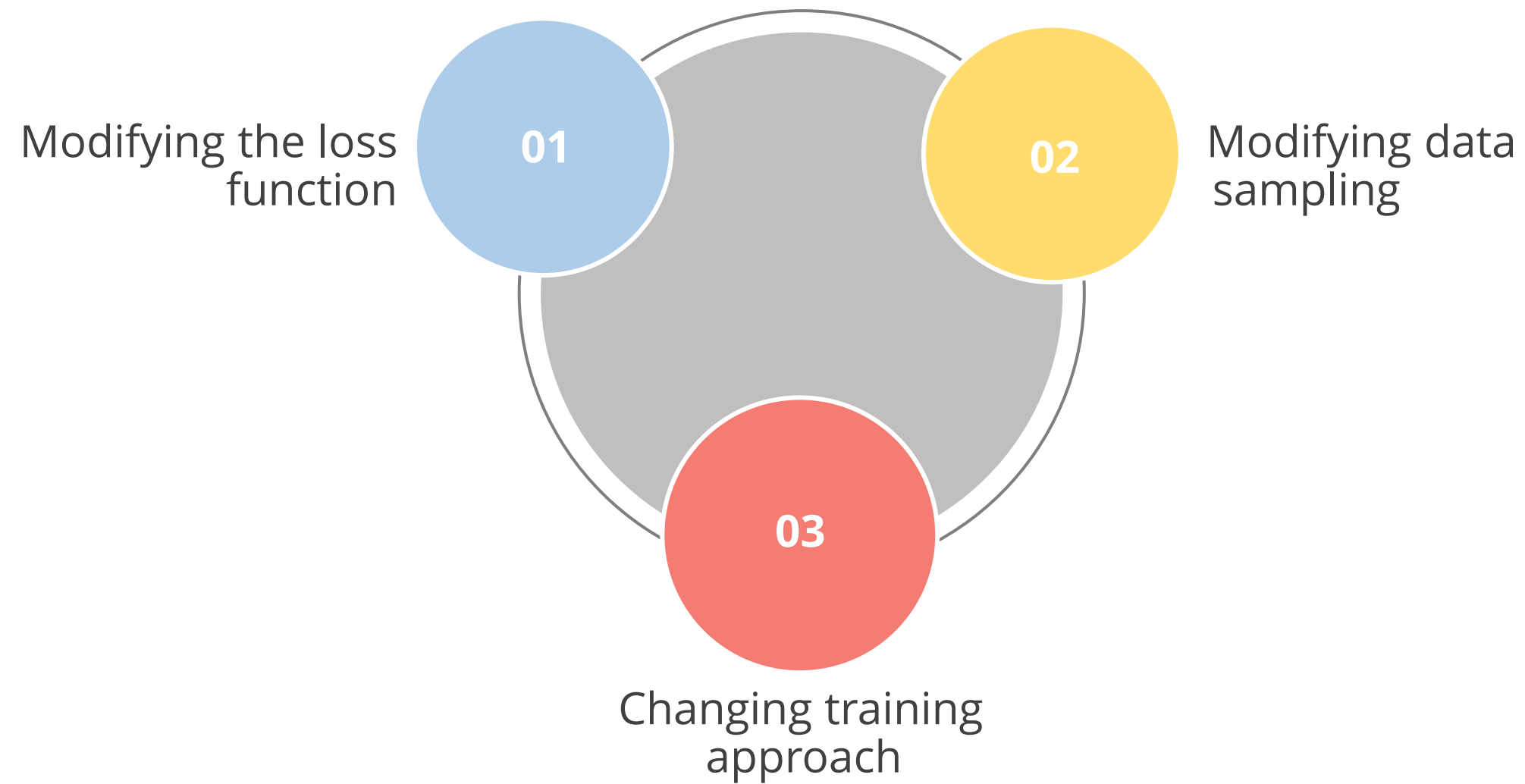
# Why Regularization?

In machine learning, regularization is frequently employed as a solution to the overfitting problem.



# Types of Regularization

The various forms of regularization employed in machine learning models are:



# Modifying the Loss Function

There are two components to modifying the loss function in the context of regularization strategies:

- In regularization strategies, the loss function is adjusted to directly consider the norm of the learned parameters or the output distribution to improve the model.
- Regularization itself involves modifying the loss function to penalize large weight values.

# Modifying the Loss Function: Strategies

The various strategies for modifying the loss function are:

01

## L2 Regularization

- Increases the model's complexity by adding more weights, which raises the risk of overfitting

02

## L1 Regularization

- Promotes weight sparsity by setting more weights to zero, instead of merely reducing their average magnitude

03

## Entropy

- Measures uncertainty in a probability distribution, where higher uncertainty corresponds to greater entropy

# Modifying Data Sampling

The two aspects related to modifying data sampling in regularization strategies are:

- The regularization approaches modify the available input data to accurately reflect the true input distribution.
- The strategies help overcome overfitting that results from the limited size of the available dataset.



# Modifying Data Sampling: Strategies

The two strategies for modifying the data sampling are:

01

## Data augmentation

- Create extra data from existing data by randomly cropping, dilating, rotating, and adding slight noise.

02

## K-fold cross-validation

- Separate the data into  $k$  groups, train with  $(k-1)$  groups, and test with the remaining  $k$  group.
- Experiment with all the possible  $k$  combinations.

# Change Training Approach

The change training approach in the context of regularization refers to altering how the training process is conducted to improve the model's performance and generalization ability.

There are two components of the change training approach in the context of regularization strategies:

## Algorithm modification

Adding regularization terms to the learning algorithm to prevent overfitting

## Data augmentation

Increasing dataset size and diversity through modifications of existing data to improve the model's generalization ability

# Change Training Approach: Strategies

The various strategies for changing training approach are:

01

## Injecting noise

- It improves generalization, prevents overfitting, and is widely used in the deep learning industry to enhance model performance on unseen data.

02

## Dropout

- It involves randomly setting a fraction of input units to zero at each update during training time, which helps prevent overfitting.



## **Dropout and Early Stopping**

# Dropout Layer

The dropout layer is a regularization technique used to prevent overfitting in neural networks.

Following are the steps to be performed for dropout:

Step 1: Choose a dropout rate, typically between 0.2 and 0.5.

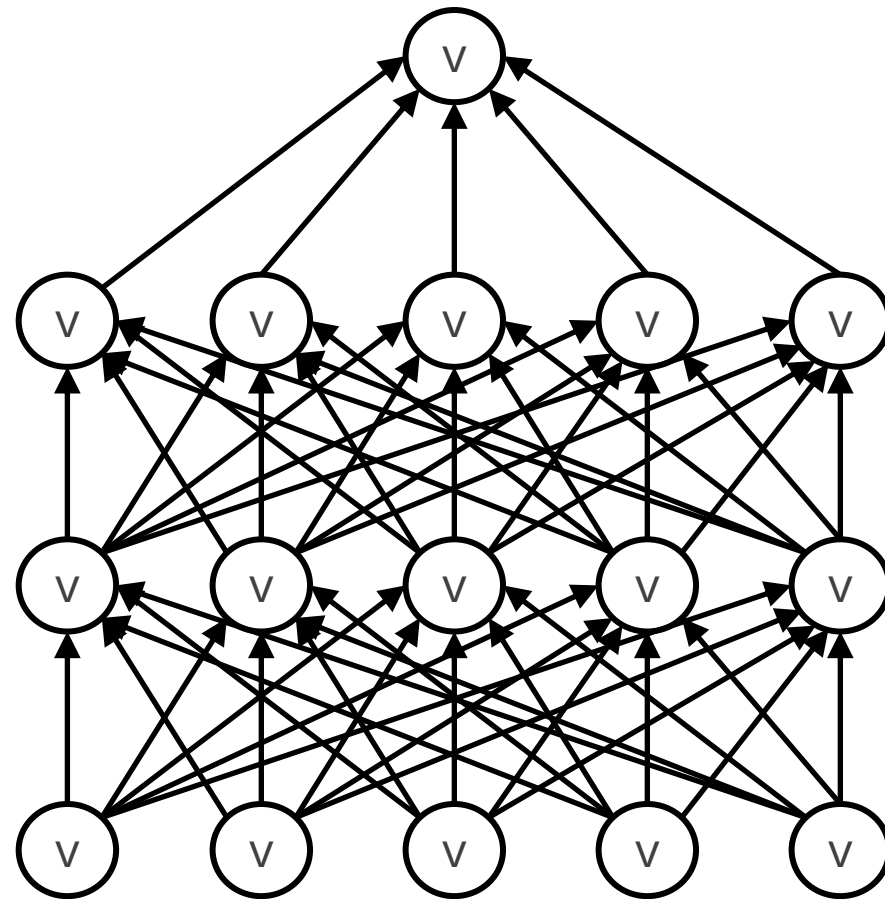
Step 2: Apply dropout during the forward pass by randomly setting a fraction of neurons to zero with the chosen dropout rate.

Step 3: Scale the remaining activations by dividing them by 1 (dropout rate).

Step 4: Complete the forward pass through the network and update weights based on gradients during backpropagation.

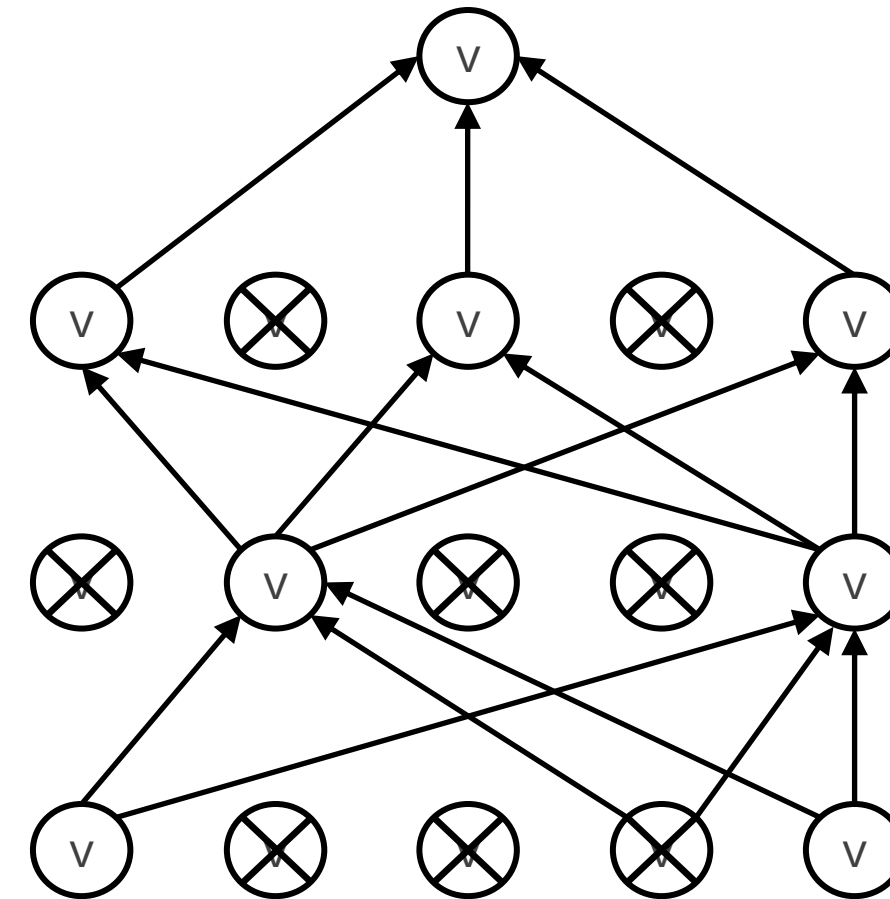
# Dropout Layer

Example: A typical neural network contains two hidden layers.



**(a) Standard neural net**

This figure depicts a standard neural network that is fully connected, with no dropout applied. Each node is connected to every other node in the next layer.

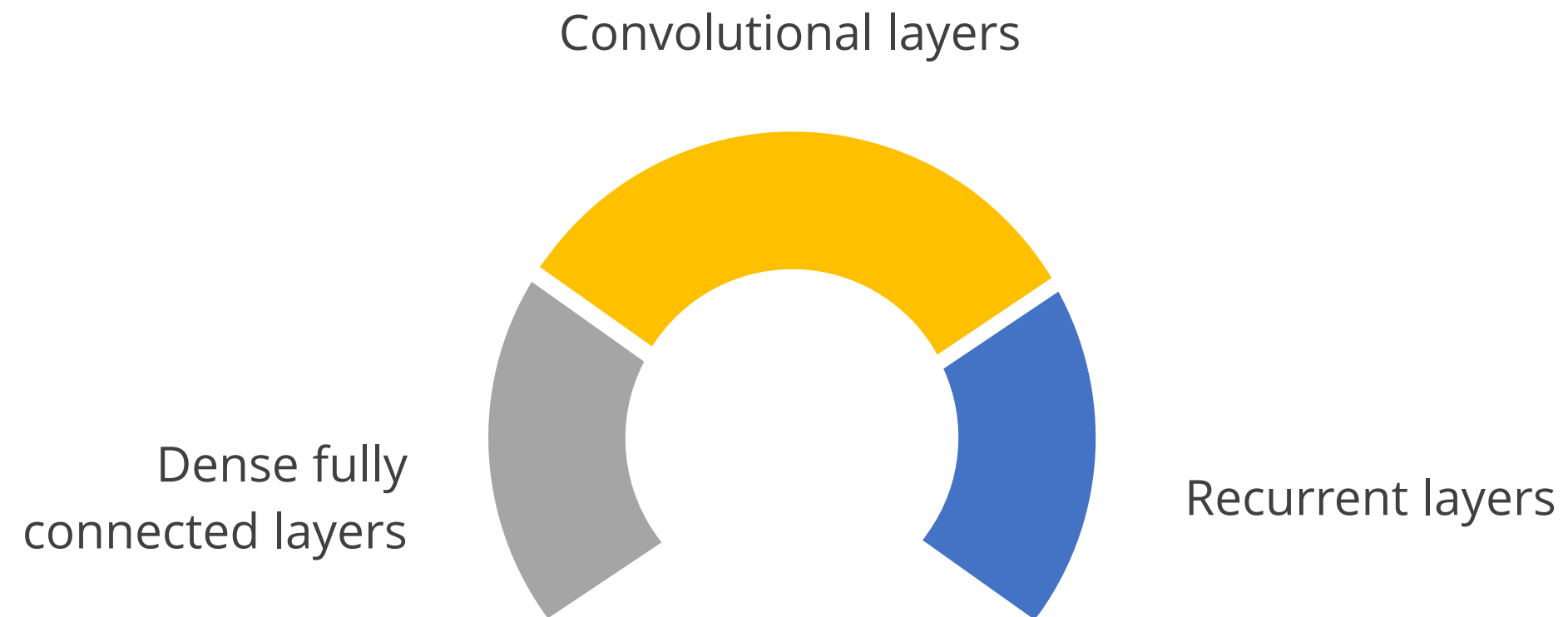


**(b) After applying dropout**

This figure depicts a thinned version of the previous network, achieved by applying dropout. Some nodes are turned off (dropped out), resulting in a thinned network.

# How to Use Dropout

In a neural network, dropout is implemented with types of layers, including:



# How to Use Dropout

The best practices for using dropout are as follows:

Use a value of 0.5 to maintain the output of every node in a hidden layer.

Use a value close to 1.0, such as 0.8, to retain inputs from the visible layer.

Dropout regularization reduces overfitting by randomly deactivating neurons during training, promoting the learning of robust features and improving generalization.



# Implementing Dropout Using Keras

Here is an implementation of dropout in a neural network for a 10-class classification problem using the Keras library:

Syntax:-

```
) # Example of output size for a 10-class classification
model = tf.keras.Sequential([
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(1024, activation=tf.nn.relu),
    tf.keras.layers.Dropout(0.2),
    tf.keras.layers.Dense(1024, activation=tf.nn.relu),
    tf.keras.layers.Dropout(0.2),
    tf.keras.layers.Dense(1024, activation=tf.nn.relu),
    tf.keras.layers.Dropout(0.2),
    tf.keras.layers.Dense(10, activation=tf.nn.softmax)
])
```

# Dropout and Early Stopping

Early stopping is a regularization strategy that reduces overfitting.

Following are the differences between dropout and early stopping:

## Dropout

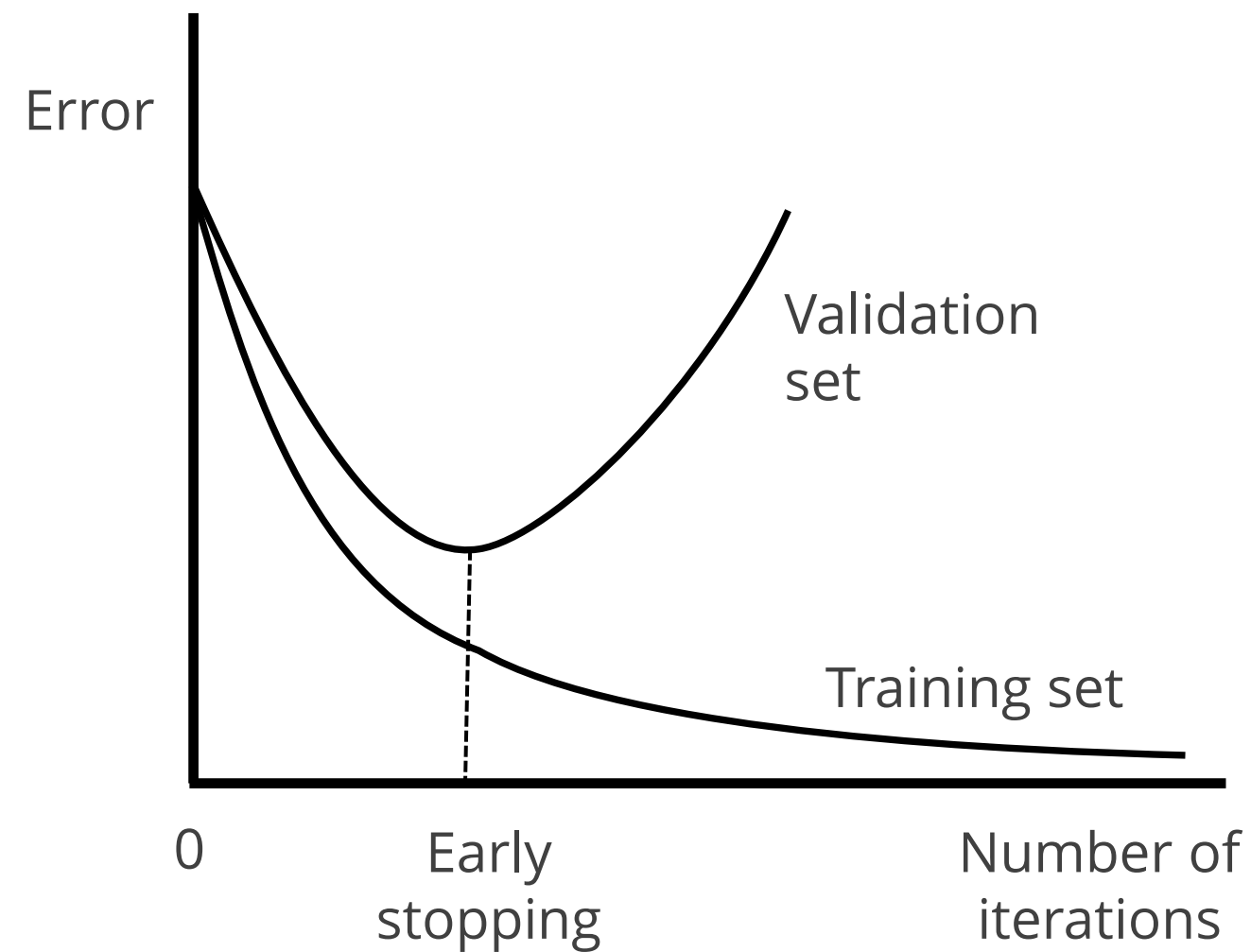
- It randomly drops nodes during training.
- It improves generalization error in deep neural networks, is computationally inexpensive, and is effective.

## Early stopping

- It allows one to specify an arbitrary number of training epochs.
- The training stops once the model's performance on a holdout validation dataset stops improving.

# Early Stopping

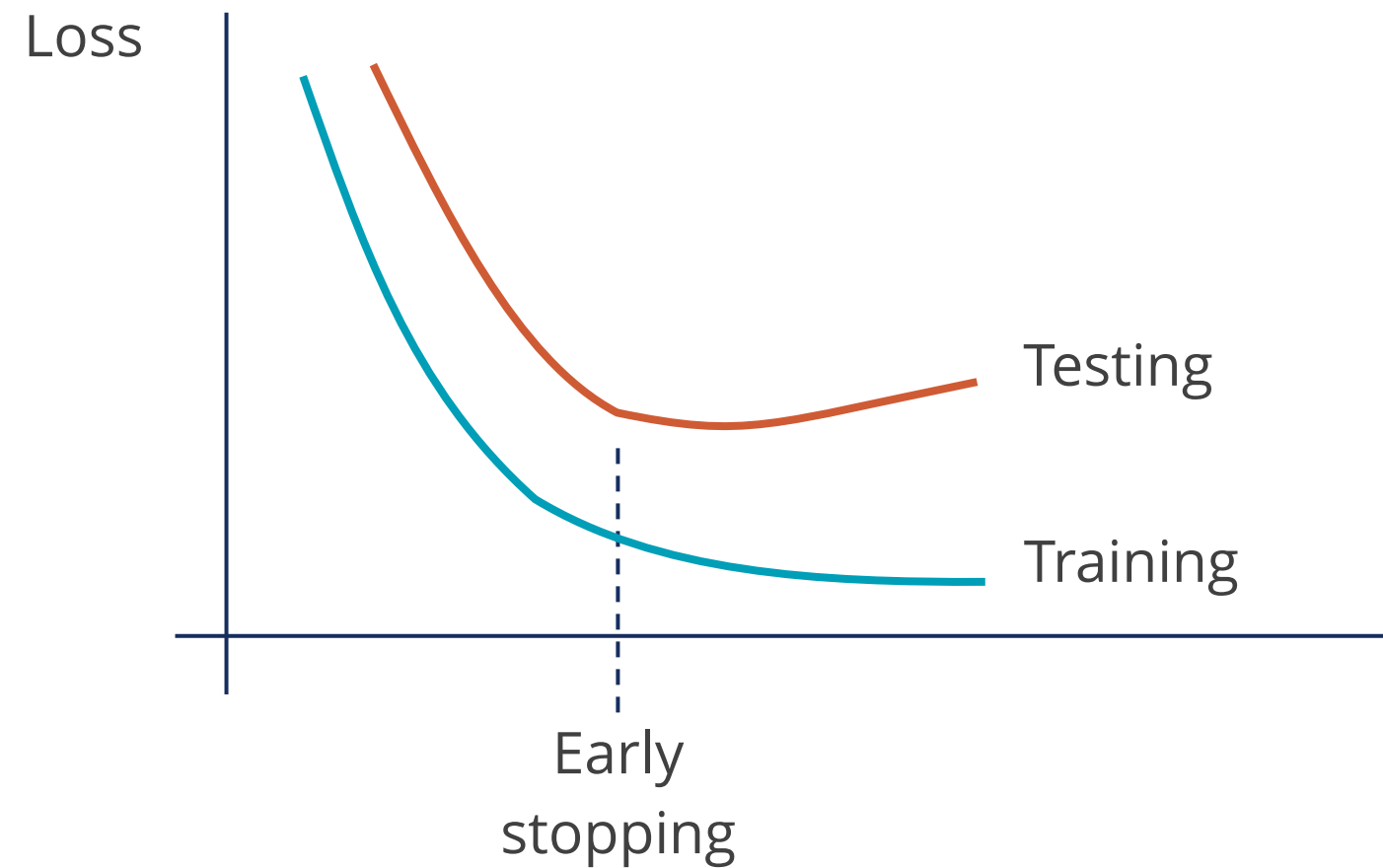
Early stopping is a regularization technique for deep neural networks that stops training when parameter updates no longer begin to yield improvements on a validation set.



By limiting the optimization approach to a smaller amount of parameter space, it functions as a regularized technique.

# Early Stopping in Keras

Keras has a callback named **EarlyStopping** that stops training early.



It lets one define the performance metric to track and a trigger to end the training.

# Early Stopping in Keras

Keras stops training when there's no more improvement in the chosen measure.

The first time with no improvement isn't always the best time to stop.

The model might not get better or might get worse before improving again.

# Early Stopping in Keras

The **patience** argument in Keras adds delay to the trigger for early stopping equal to the number of epochs where no improvement is desired.

This parameter allows the model more time to recover and improve before stopping the training process.



The required patience value varies depending on the model type and difficulty.

# Implementing Early Stopping Using Keras

Users may invoke various arguments to configure **Earlystopping**.

To stop training, use the **monitor** option to indicate the performance metric to track.

Syntax:-

```
callback = tf.keras.callbacks.EarlyStopping(monitor='val_accuracy', patience=3)

model = tf.keras.Sequential([
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(1024, activation=tf.nn.relu),
    tf.keras.layers.Dropout(0.2),
    tf.keras.layers.Dense(1024, activation=tf.nn.relu),
    tf.keras.layers.Dropout(0.2),
    tf.keras.layers.Dense(1024, activation=tf.nn.relu),
    tf.keras.layers.Dropout(0.2),
    tf.keras.layers.Dense(10, activation=tf.nn.softmax)
])
model.compile(optimizer=tf.keras.optimizers.Adam(),
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

history = model.fit(x_train, y_train, epochs=50, validation_data=(x_test, y_test), callbacks=[callback])
```

# Assisted Practices



Let's understand the concept of Dropout using Jupyter Notebooks.

- 7.18\_Implementation of Dropout

**Note:** Please refer to the **Reference Material** section to download the notebook files corresponding to each mentioned topic.

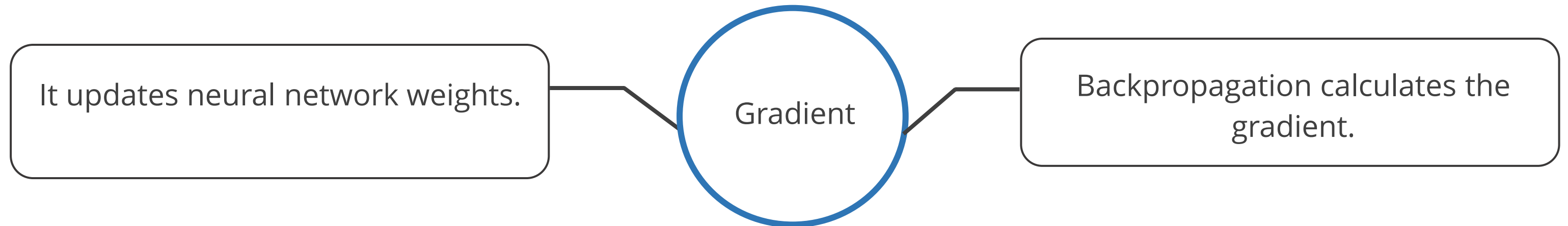




# Vanishing Gradient

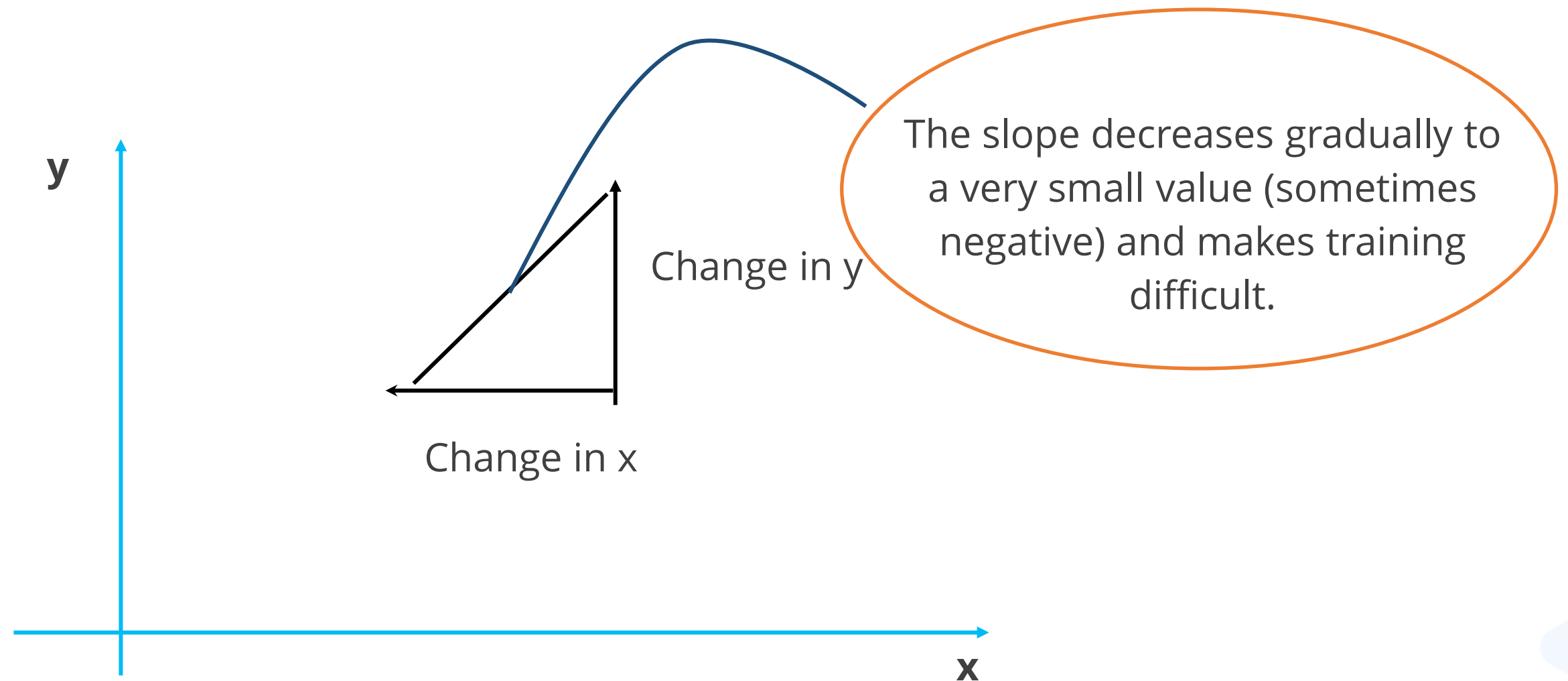
# What Is Gradient?

It refers to the derivative of loss with respect to weight.



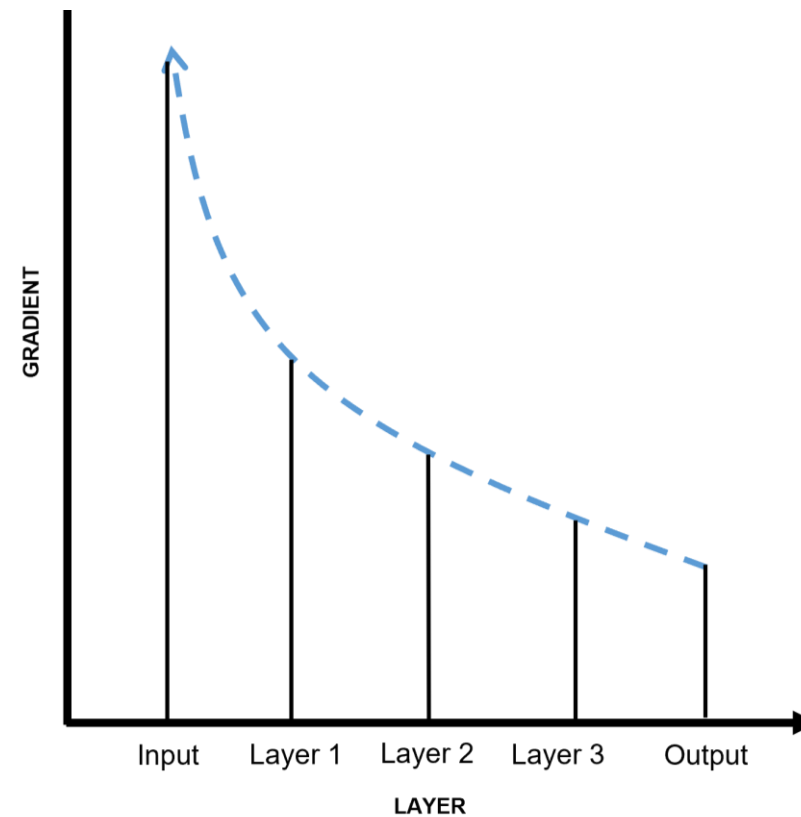
# What Is Vanishing Gradient?

When the gradient becomes very small, subtracting it from the weight doesn't change the previous weight. Therefore, the model stops learning. This problem of neural networks is called the vanishing gradient.



# Vanishing Gradients

In vanishing gradients, the gradients become smaller as the backpropagation method progresses backward from the output layer to the input layer.



The lower layer weights remain unchanged, and the gradient descent never reaches the optima.

# Vanishing Gradient Issue

The vanishing gradient issue refers to the challenge of updating the weights in the previous layers of a neural network during training.

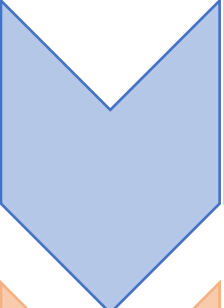
It occurs when the gradients become very small as they propagate backward, resulting in little or no update to the weights.



A vanishing gradient issue is a difficulty with weights in the network's previous layers during training.

The network may struggle to learn complex patterns and exhibit reduced performance.

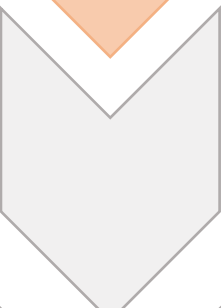
# How to Prevent Vanishing Gradient?



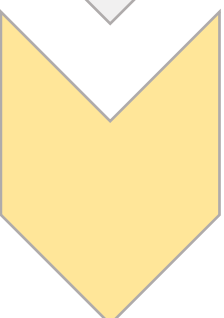
Residual Networks (ResNets) use shortcut connections to effectively address the vanishing gradient problem and enable the training of deep neural networks.



GPUs play a crucial role in mitigating the vanishing gradient issue through parallel processing, faster training, and more frequent weight updates during backpropagation.



Choosing the right activation function, like rectified linear (ReLU), helps prevent the vanishing gradient problem by avoiding input compression into small ranges.



The switch from CPUs to GPUs has significantly improved the feasibility of standard backpropagation, even for low-cost models, due to their faster compilation times.



# Exploding Gradient

# What Is Exploding Gradient?

“

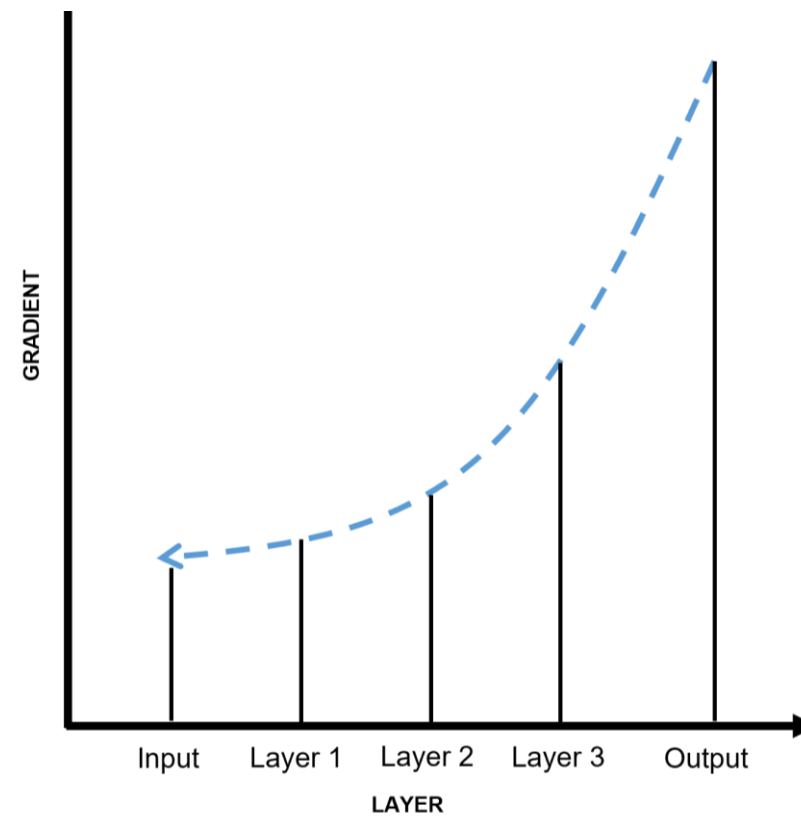
The issue of exploding gradients arises when there is a significant accumulation of error gradients. It leads to excessively large weight updates in a neural network during the training process.

”



# Exploding Gradients

When the gradients grow larger with the advancement of the backpropagation method, massive weight updates occur, causing the gradient descent to diverge.



# How to Fix Exploding Gradients?

It can be fixed by redesigning the neural network with fewer layers and mini-batch sizes.

Using Long Short-Term Memory (LSTM) networks reduces exploding gradients.

Gradient clipping limits gradient size to effectively fix exploding gradients.



# Hyperparameter Tuning

# What Are Parameters and Hyperparameters?

## Parameters

- They are found during model training.
- These internal variables are adjusted to make predictions based on the input data.
- For example, in K-means clustering, the positions of the centroids are learned during training. These positions are the model's parameters.

## Hyperparameters

- They are determined before training.
- These are configurations or settings that govern the learning process but aren't learned from the data.
- For example, the value of K in K-means clustering is decided before creating the model. This value, representing the number of clusters, is a hyperparameter.

# Hyperparameters of Deep Learning Models

## Learning rate

It is the most important hyperparameter that helps the model get an optimized result.

## Number of hidden units

It is a classic hyperparameter that specifies the representational capacity of a model.

## Convolutional kernel width

It determines the size of the filters in a convolutional neural network, which influences the receptive field and the capacity of the model.

## Mini-batch size

It affects the training process, training speed, and number of iterations in a deep learning model.

## Number of epochs

It is partly responsible for the weight optimization in a neural network.

# Hyperparameter Tuning

“

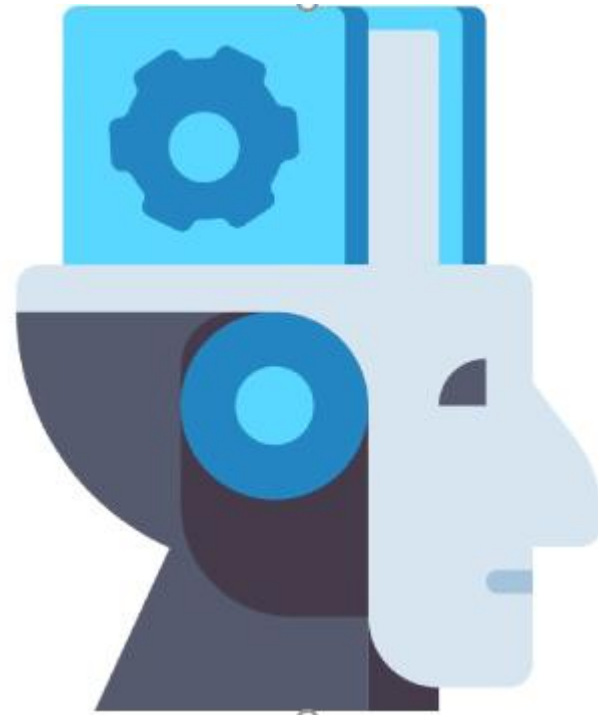
Hyperparameter tuning is the process of selecting a set of optimal hyperparameters for a learning algorithm.

”

# Hyperparameter Tuning

Hyperparameter tuning involves finding the optimal values for hyperparameters, which are manually set parameters in a machine-learning model.

The goal of iteratively adjusting hyperparameter values and evaluating model performance is to enhance the model's effectiveness and generalization.



Grid search, random search, and Bayesian optimization are popular hyperparameter tuning techniques.

# Hyperparameter Tuning

Unlike parameters, hyperparameters are not estimated directly from training data.

Hyperparameters define the model's complexity and learning efficiency.

Set hyperparameters before training to control the model's behavior.

Tune hyperparameters correctly to improve performance and efficiency.



# Hyperparameter Tuning

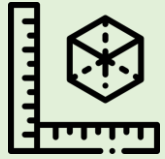
Some hyperparameters, like the momentum hyperparameter, have specific values that are typically used.



Generalizing the values of hyperparameters might not be practical for real-world data.

Therefore, finding the best set of hyperparameters is a quintessential search problem.

# How to Tune the Hyperparameters?



## **Choose the parameters wisely**

Select the most influential parameters, as it is not possible to tune all of them.



## **Understand the training process**

Know the training process and how exactly it can be influenced.

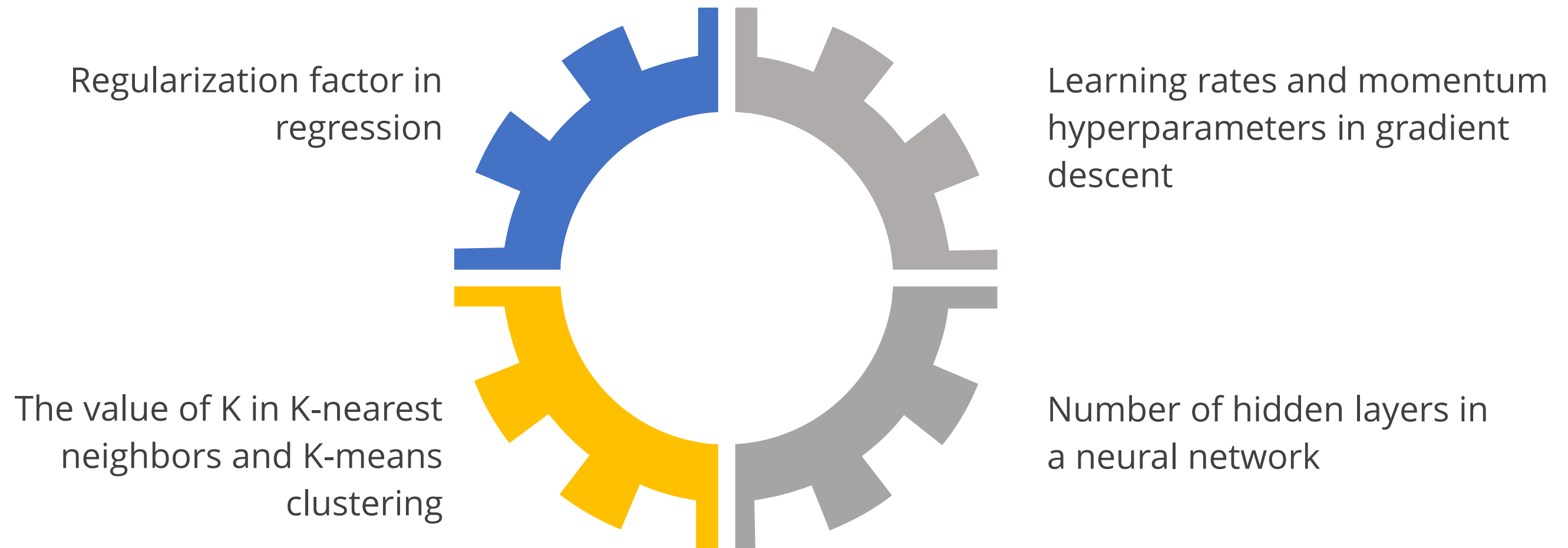


## **Perform a grid search**

Search over the defined range of hyperparameters using grid or random search methods.

# Hyperparameter: Examples

Some examples of hyperparameters are:



# Selection of Hyperparameters

There are two approaches to selecting hyperparameters:

Manual  
hyperparameter tuning

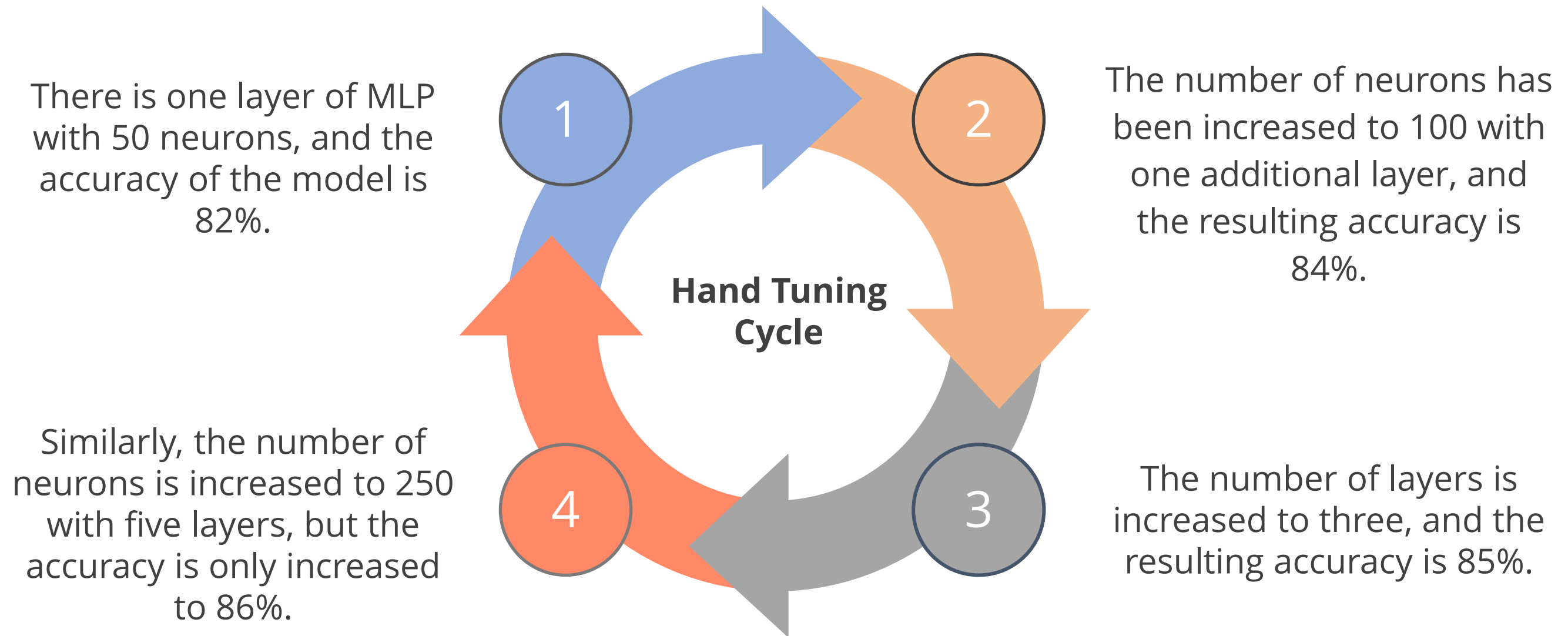
This involves manually selecting and tuning hyperparameters based on intuition, experience, and trial and error.

Automatic  
hyperparameter tuning

Use algorithms to systematically search the hyperparameter space to find the optimal values.

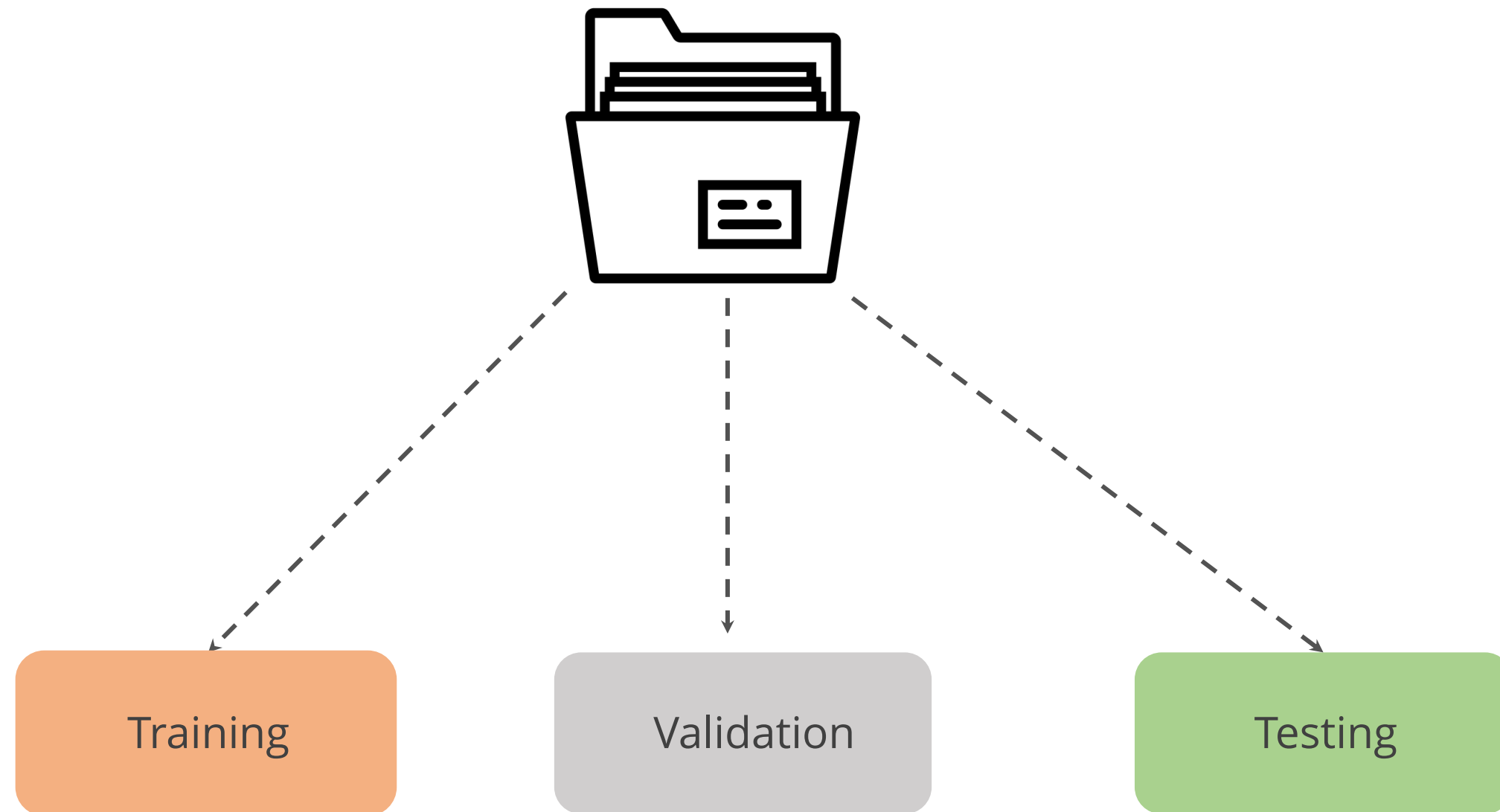
# Manual Tuning Approach

The following example illustrates that manual tuning is not efficient, as even a fivefold increase in the neurons resulted in an improvement in accuracy of only 4%.



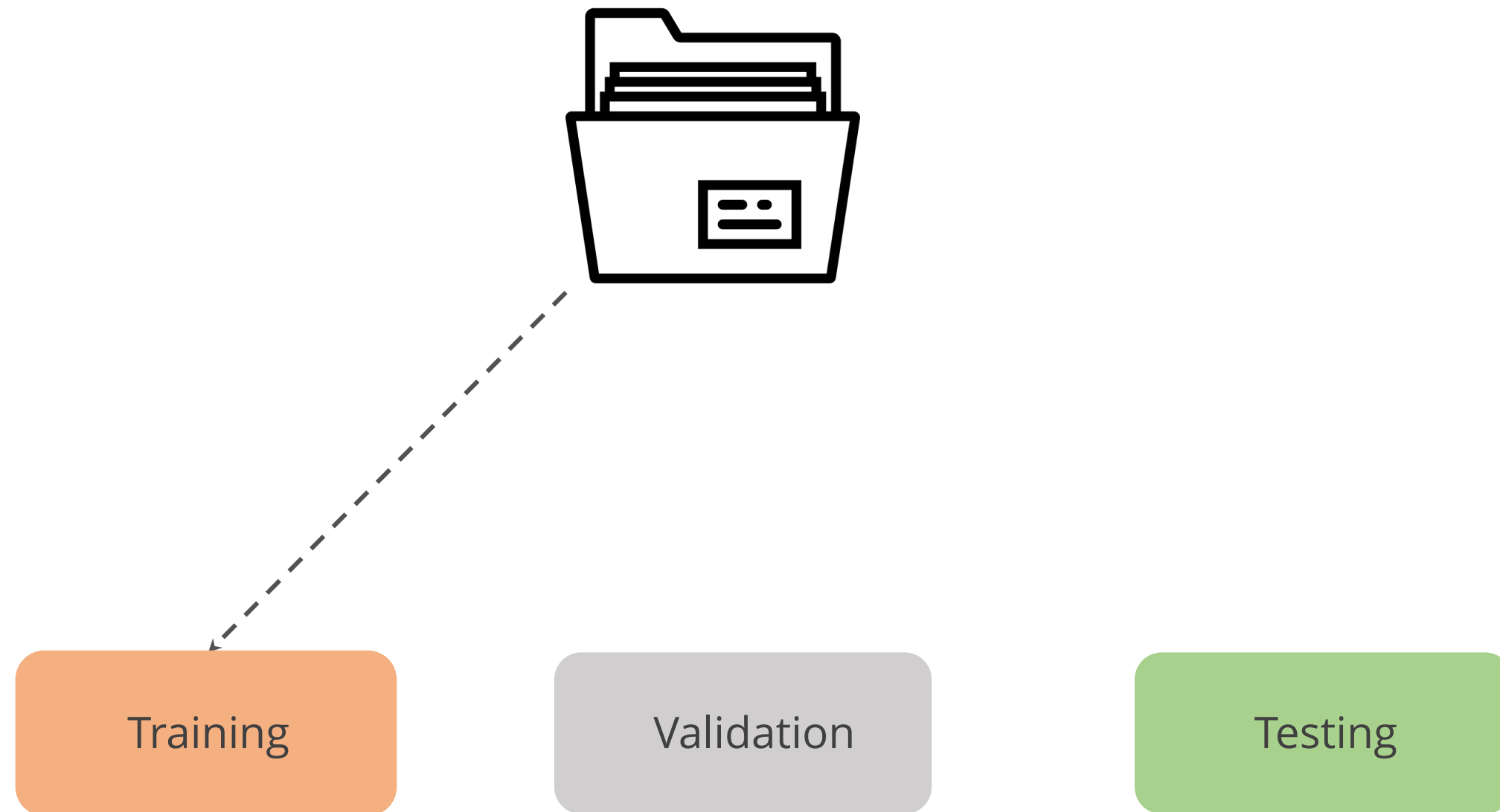
# Data Partitioning for Hyperparameter Selection

There are three categories of data used for selecting hyperparameters:



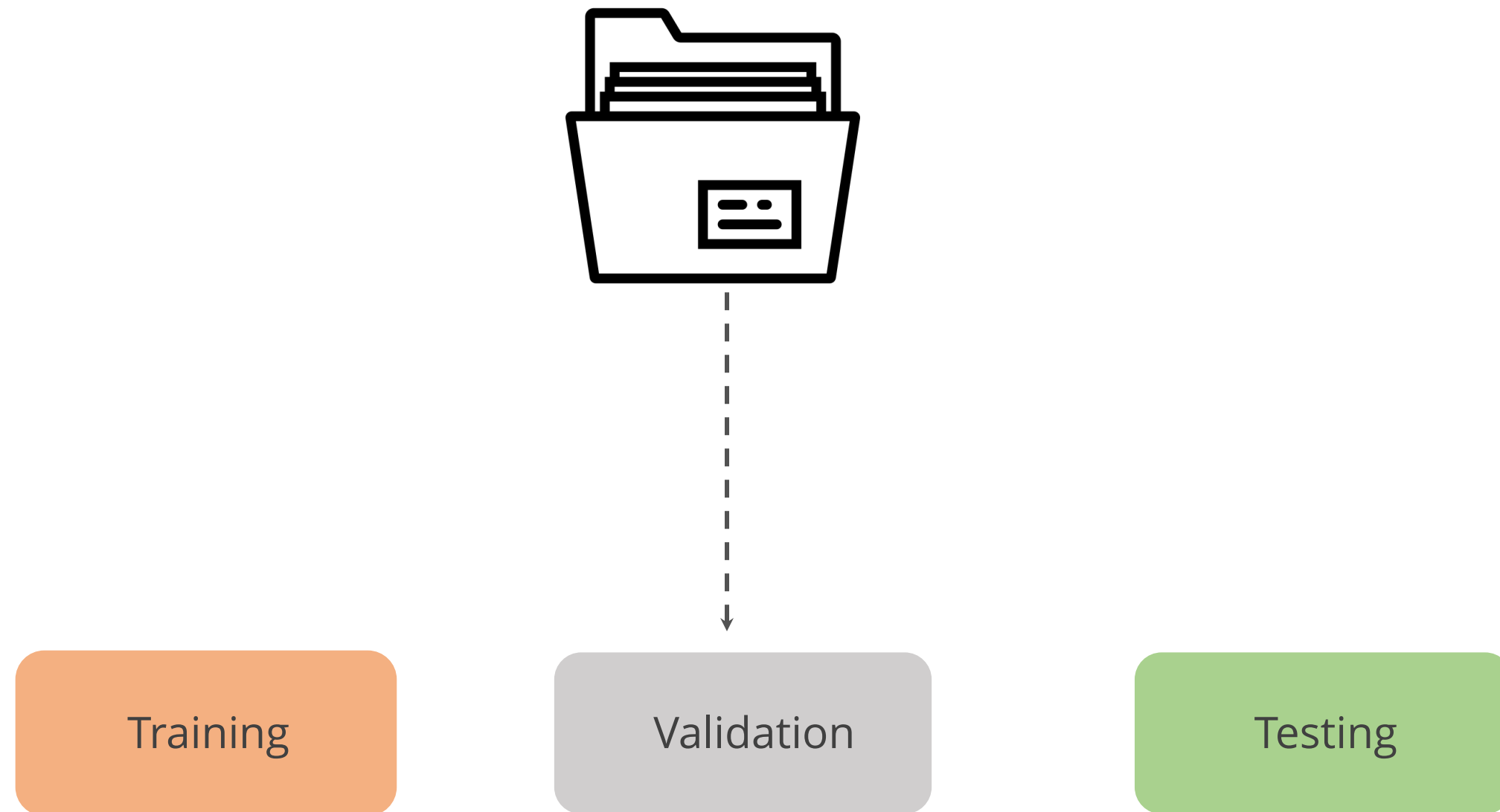
# Data Partitioning for Hyperparameter Selection

The training set trains the model with different hyperparameter combinations.



# Data Partitioning for Hyperparameter Selection

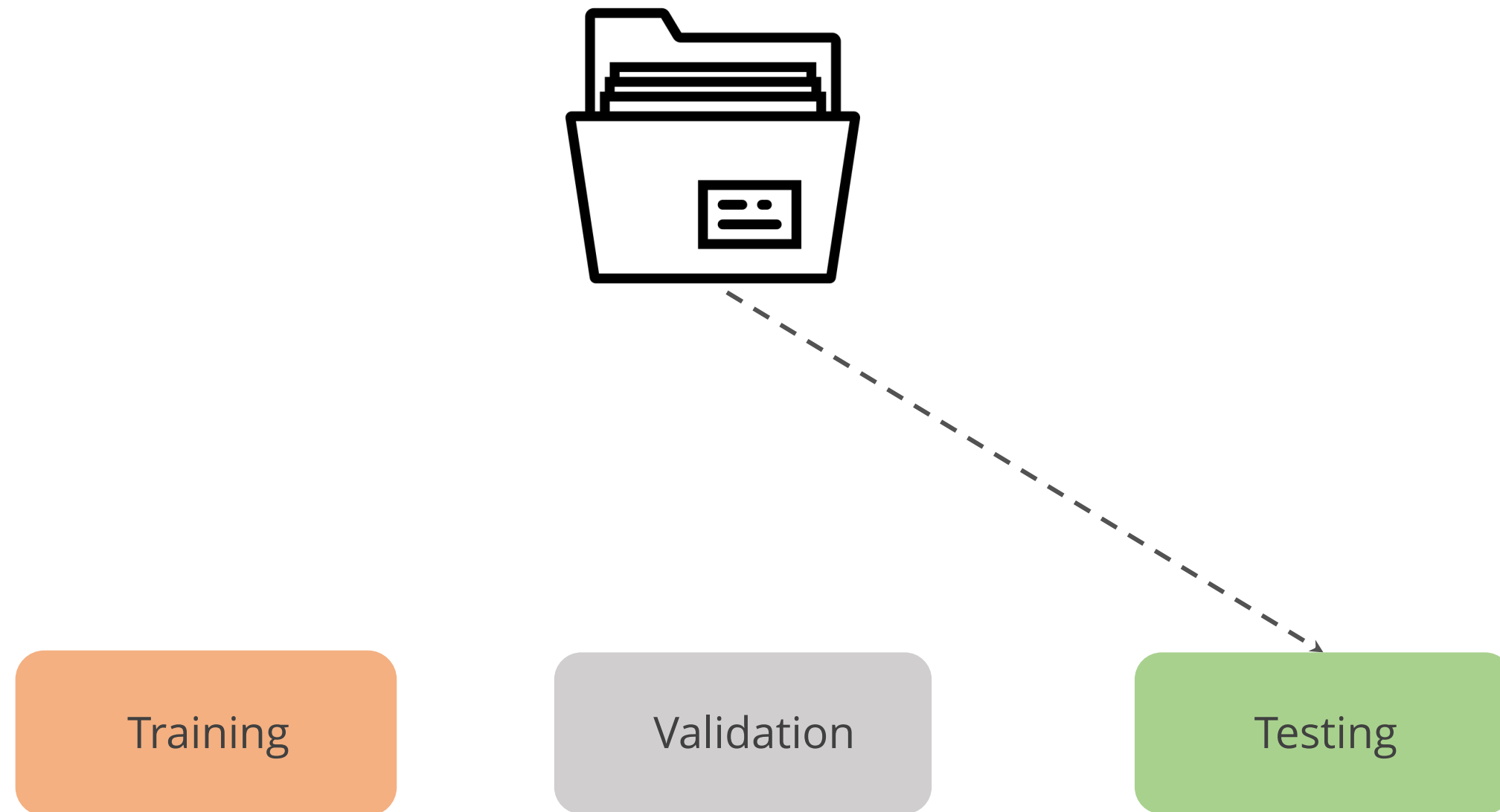
The validation set identifies which parameters minimize error.





# Data Partitioning for Hyperparameter Selection

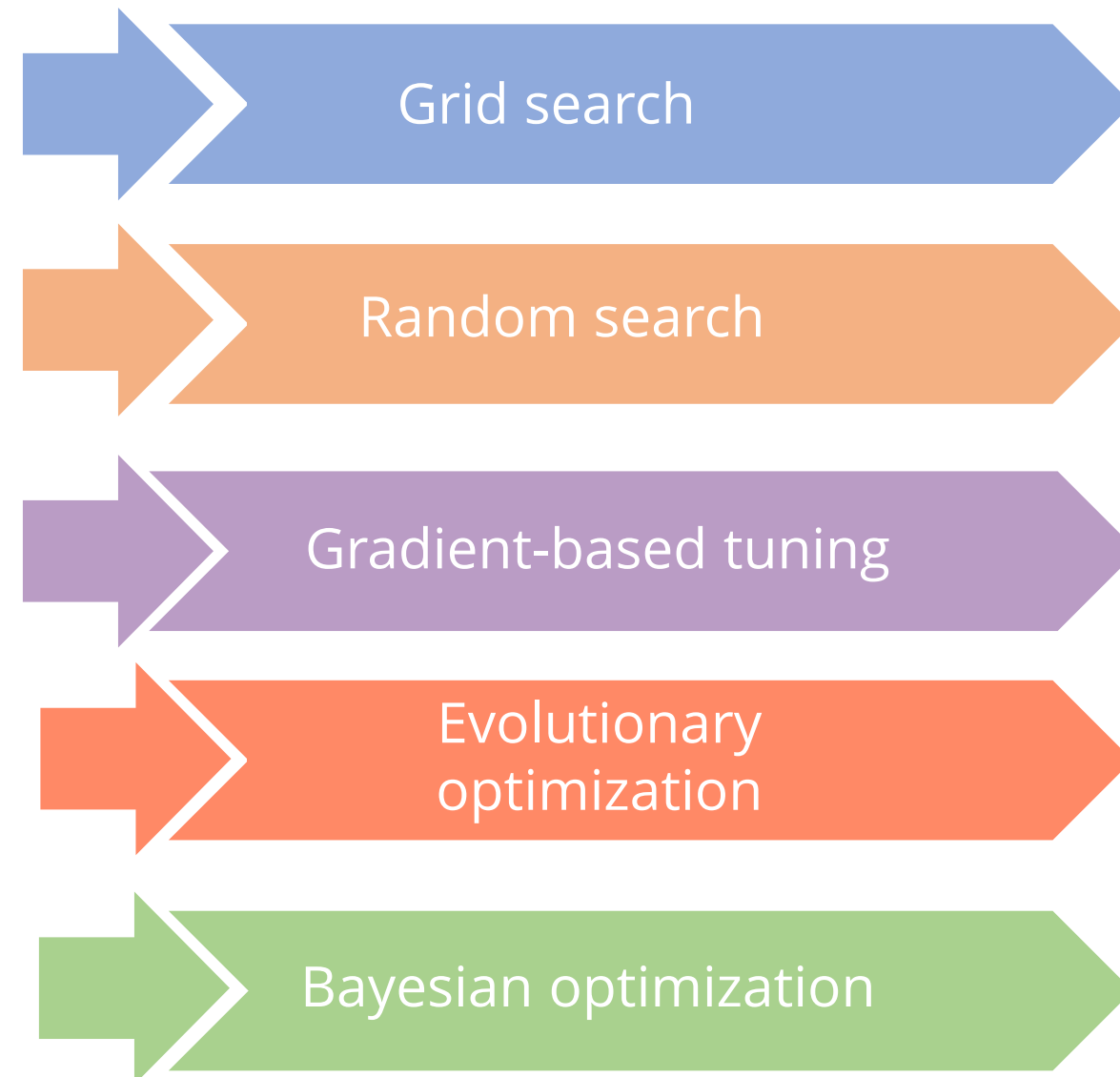
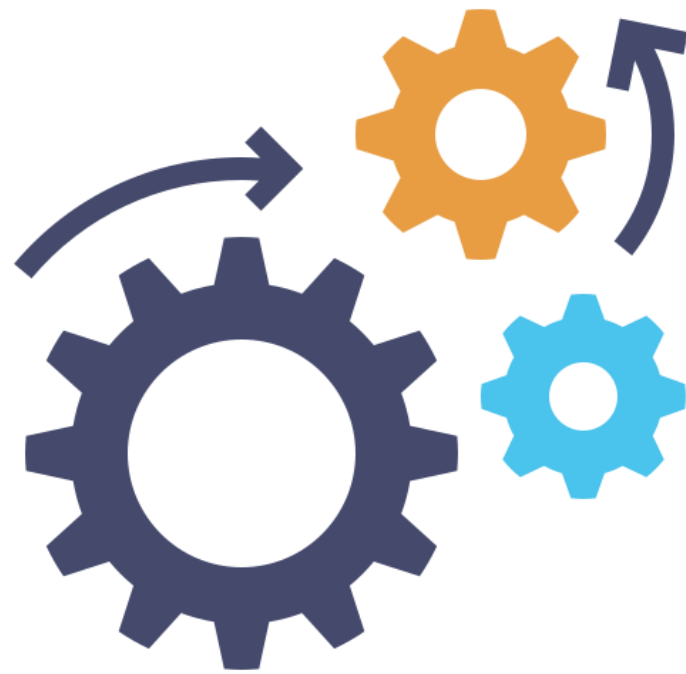
The test set assesses performance with selected parameters.



# Automatic Hyperparameter Tuning

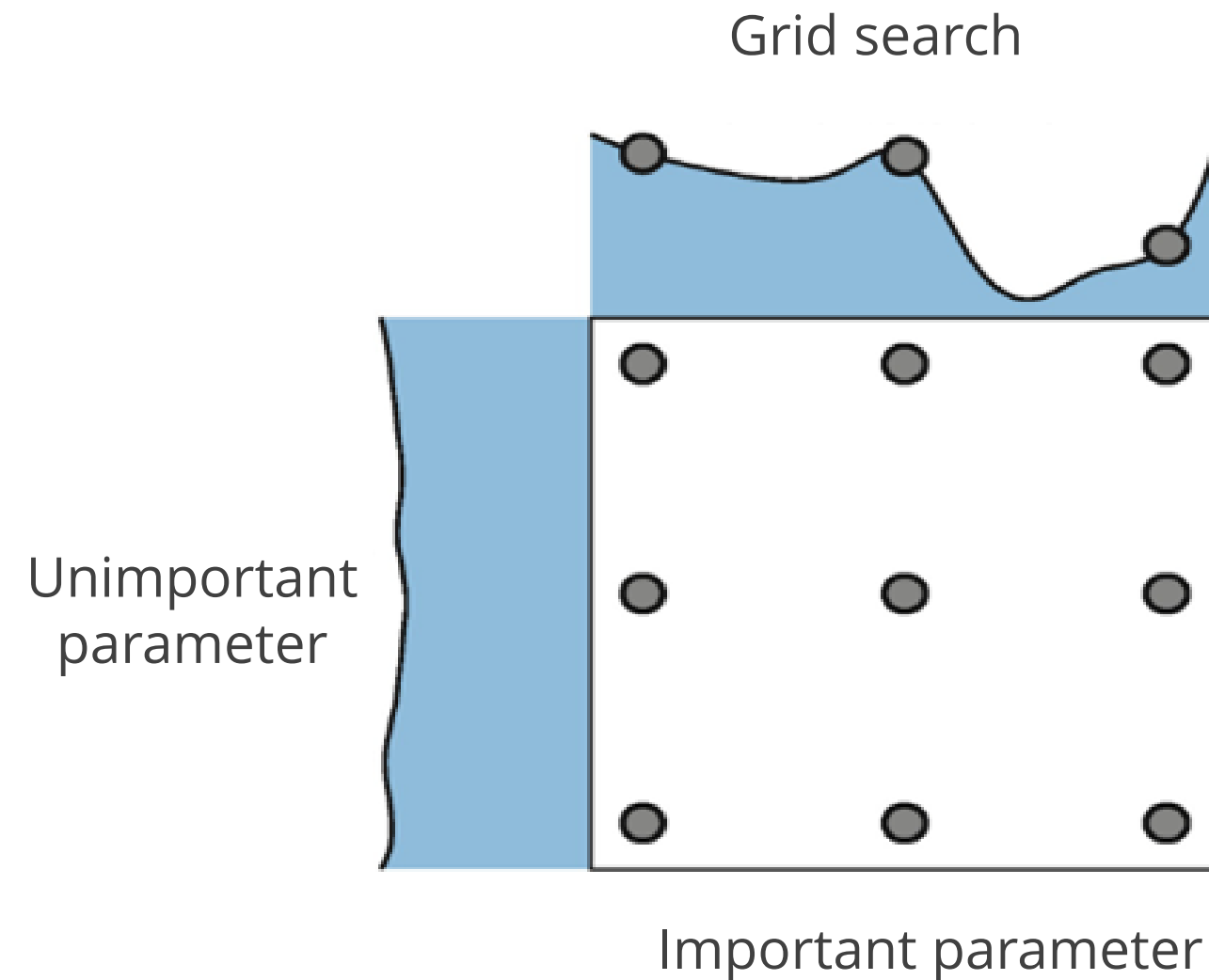
The automatic selection approach is preferred over the manual approach, as the latter is a very rigorous method. The automatic approach is the process of tuning the hyperparameters with the help of algorithms.

The popular tuning approaches are as follows:



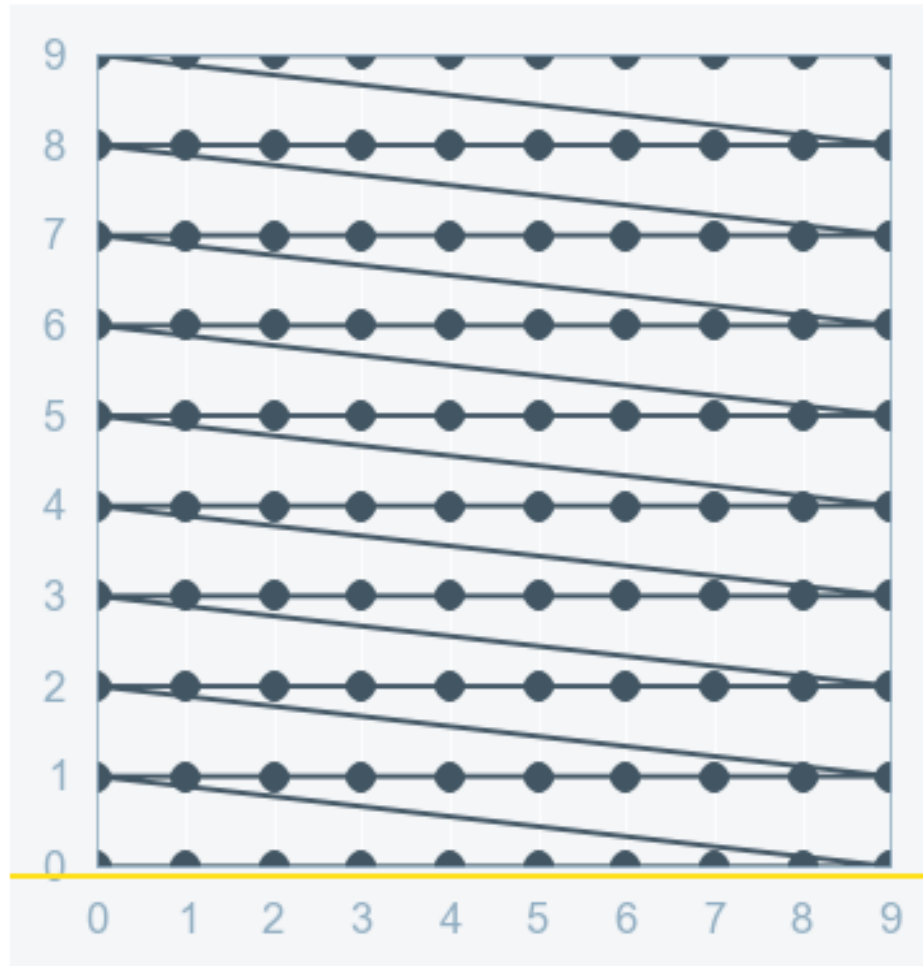
# Hyperparameter Tuning Techniques: Grid Search

Iterating over given hyperparameters using cross-validation is called **Grid Search**.



# Hyperparameter Tuning Techniques: Grid Search

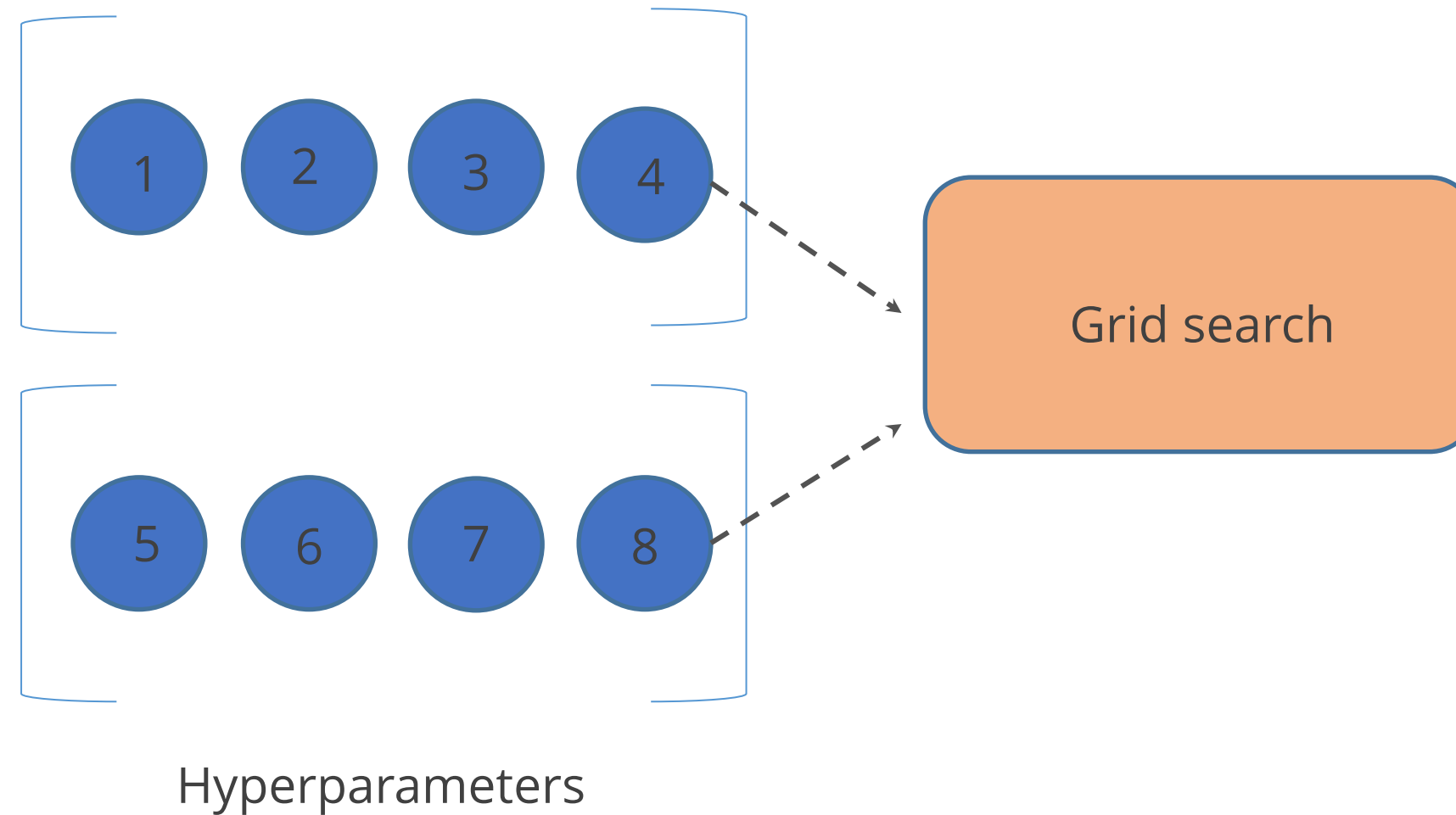
The process of grid search for hyperparameter tuning is as follows:



- **Parameter grid construction:** Arranges all potential hyperparameter combinations in a grid layout
- **Matrix conversion:** Represents each unique hyperparameter combination in a matrix for systematic processing
- **Performance evaluation:** Trains and assesses models for each distinct set of parameters, typically using a validation set
- **Best model identification:** Chooses the model with the highest performance score (based on metrics like accuracy, precision, and various others) as the grid search's optimal outcome

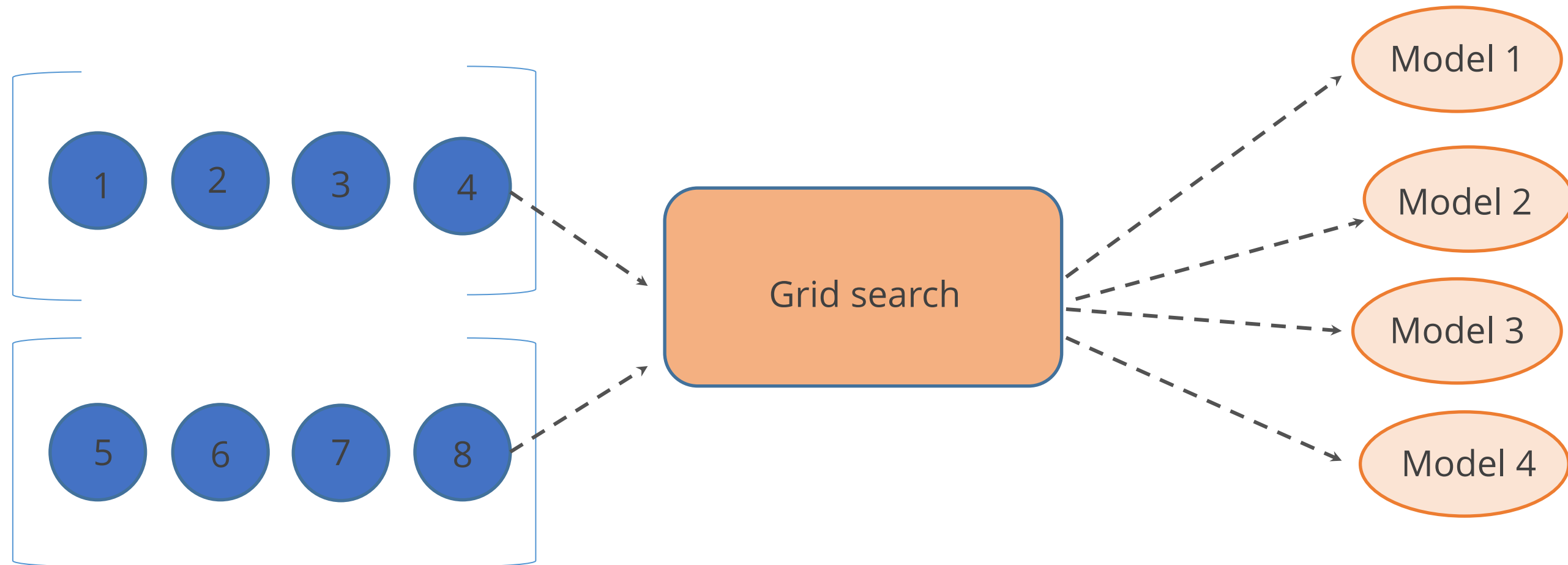
# Hyperparameter Tuning Techniques: Grid Search

Out of the eight hyperparameters given below, Grid search evaluates different combinations of hyperparameter values to find the optimal configuration.



# Hyperparameter Tuning Techniques: Grid Search

It creates four models using the available hyperparameters.



Grid search selects the model with the lowest error as the most efficient and finalizes those hyperparameters.

# Hyperparameter Tuning Techniques: GridSearchCV

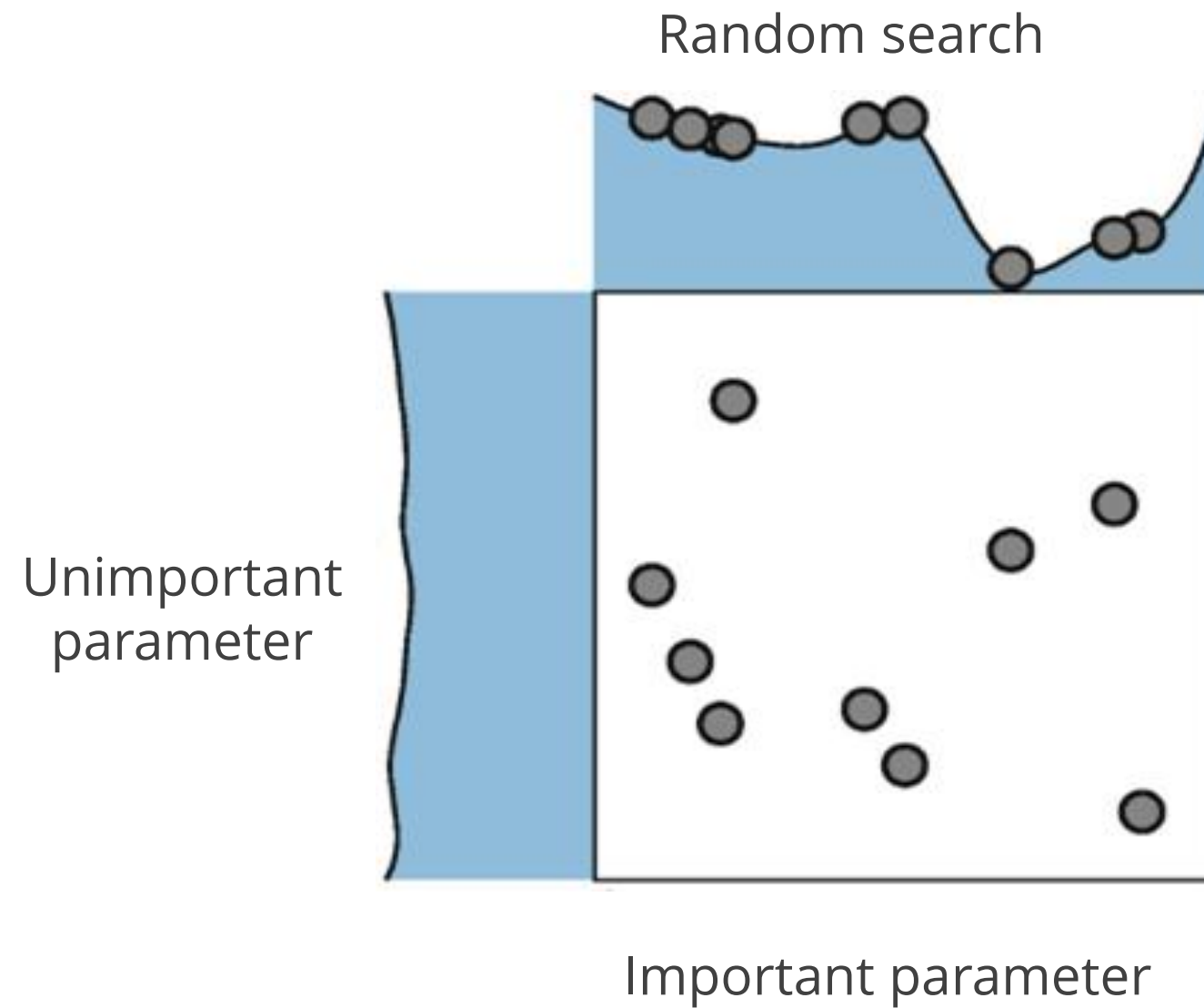
The GridSearchCV algorithm extends grid search by adding cross-validation, evaluating all hyperparameter combinations with different data splits.

- GridSearchCV constructs multiple versions of the machine learning model by systematically testing all possible combinations of the specified hyperparameter values.
- By incorporating cross-validation, GridSearchCV evaluates each possible set of hyperparameters, fitting the model on different training and validation splits to ensure reliable performance metrics. While this method is thorough, it is also computationally intensive.

This thorough approach offers better model validation but at the cost of higher computational effort.

# Hyperparameter Tuning Techniques: Random Search

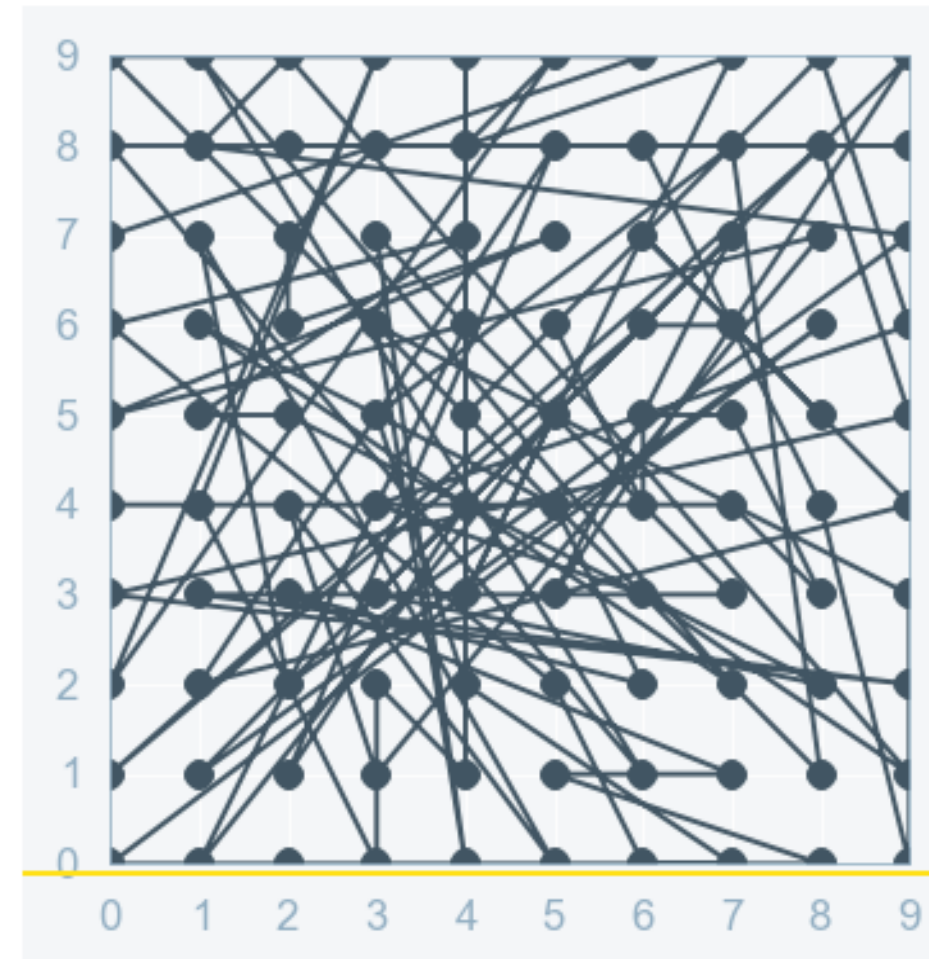
Random search is commonly used as a hyperparameter tuning method for functions that are non-differentiable or discontinuous, including those with complex, nonlinear behavior.





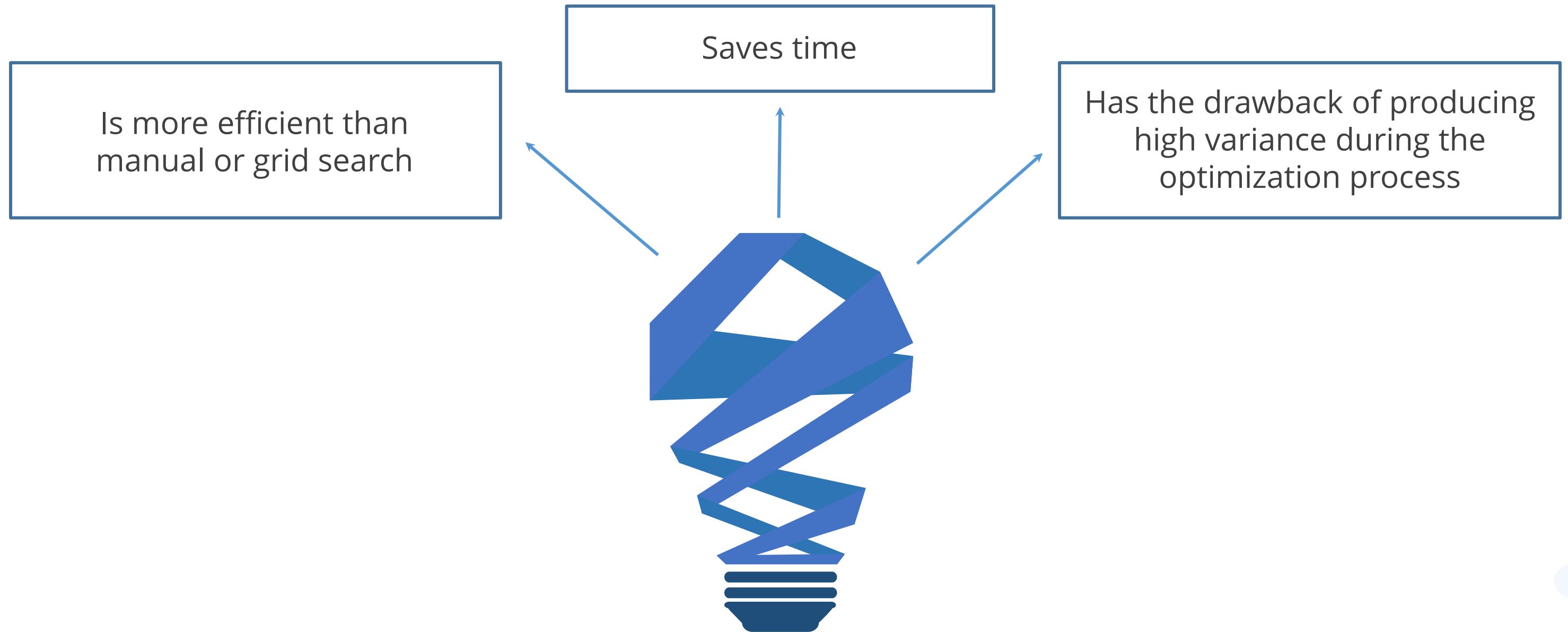
# Hyperparameter Tuning Techniques: Random Search

- Produces a random value at each instance
- Covers every combination of instances
- Considers a random combination of parameters at every iteration
- Finds the optimized parameter through the performance of models



# Hyperparameter Tuning Techniques: Random Search

The key characteristics of using random search for hyperparameter tuning include the following:



# Hyperparameter Tuning Techniques: RandomSearchCV

RandomSearchCV is an advanced version of random search, efficiently sampling hyperparameter combinations to reduce computational burden in machine learning tuning.

- It explores a set number of randomly chosen hyperparameter combinations, constructing multiple versions of the model while reducing the computational burden.
- It incorporates cross-validation to assess model performance with a more efficient search process.
- This balanced approach often leads to nearly as good results with significantly reduced computational effort.

# Gradient-Based Tuning

“

Gradient-based tuning is used for algorithms where it is possible to compute the hyperparameter for the gradient, and optimization of the hyperparameter is done by the gradient descent.

”

# Evolutionary Optimization

Evolutionary algorithms mimic natural evolution to optimize solutions.

They use processes like selection, crossover, mutation, and replacement to explore and adapt solutions.

These algorithms can be applied to find optimal hyperparameters in various models.

They are especially useful in black box functions with noise, aiding in global optimization.



# Bayesian Optimization



Bayesian optimization is an advanced method for hyperparameter tuning that uses a probabilistic model to predict the performance of different hyperparameter combinations and iteratively selects the most promising ones to evaluate.



It inherently studies the trend in each dataset, which is not possible for a human.



# Interpretability

# What Is Interpretability?

“

Interpretability is the degree of a human's ability to predict the model's result consistently.

”



# Importance of Interpretability

Fairness

Ensures that predictions are unbiased

Privacy

Ensures that sensitive information in the data is well-protected

Reliability

Ensures that small changes in the data do not lead to big changes in the prediction

Causality

Checks that all the relationships picked up are causal

Trust

Ensures it explains its decisions less like a machine so is easily trustable for humans

# When Is Interpretability Not Needed?

Interpretability is not used in these situations:

For an insignificant model

For a well-studied and researched problem

For scenarios where people or the program might manipulate the model

# Classification of Interpretability Methods



Intrinsic or post-hoc

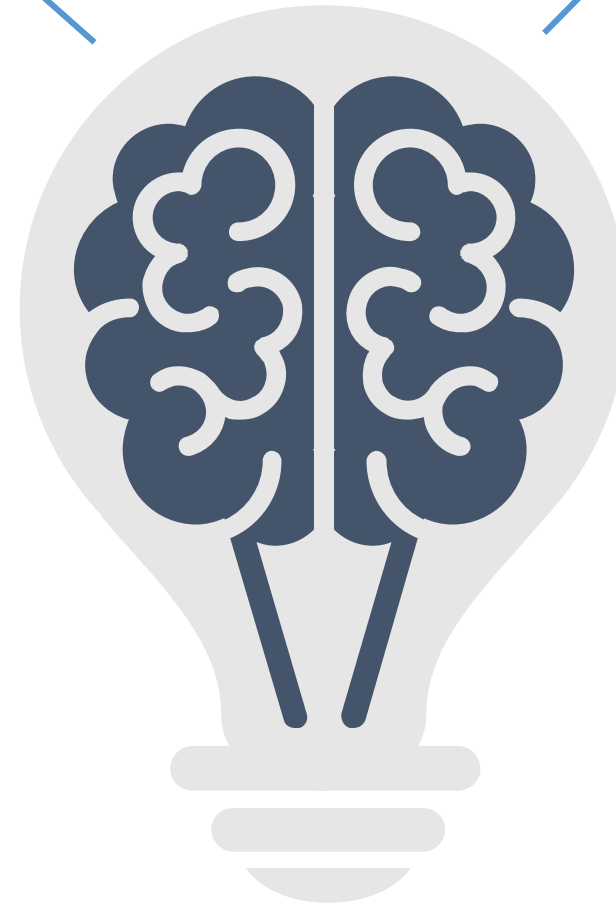


Model-specific or model-agnostic

## Intrinsic or Post-Hoc

Achieves interpretability by simplifying the machine learning model and analyses the method after training

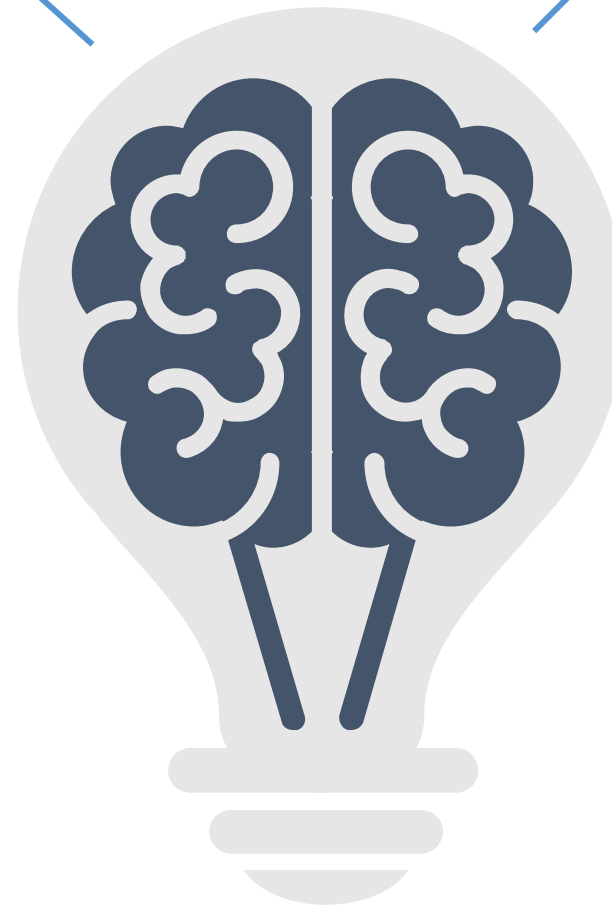
Refers to models that are considered interpretable due to their simple structure



## Model-Specific or Model-Agnostic

Can be used on any model and are applied after the model has been trained

Works by analyzing feature input and output pairs



# Scope of Interpretability



Algorithm transparency

Global, holistic model interpretability

Global model interpretability on a modular level

Local interpretability for a single prediction

Local interpretability for a group of predictions

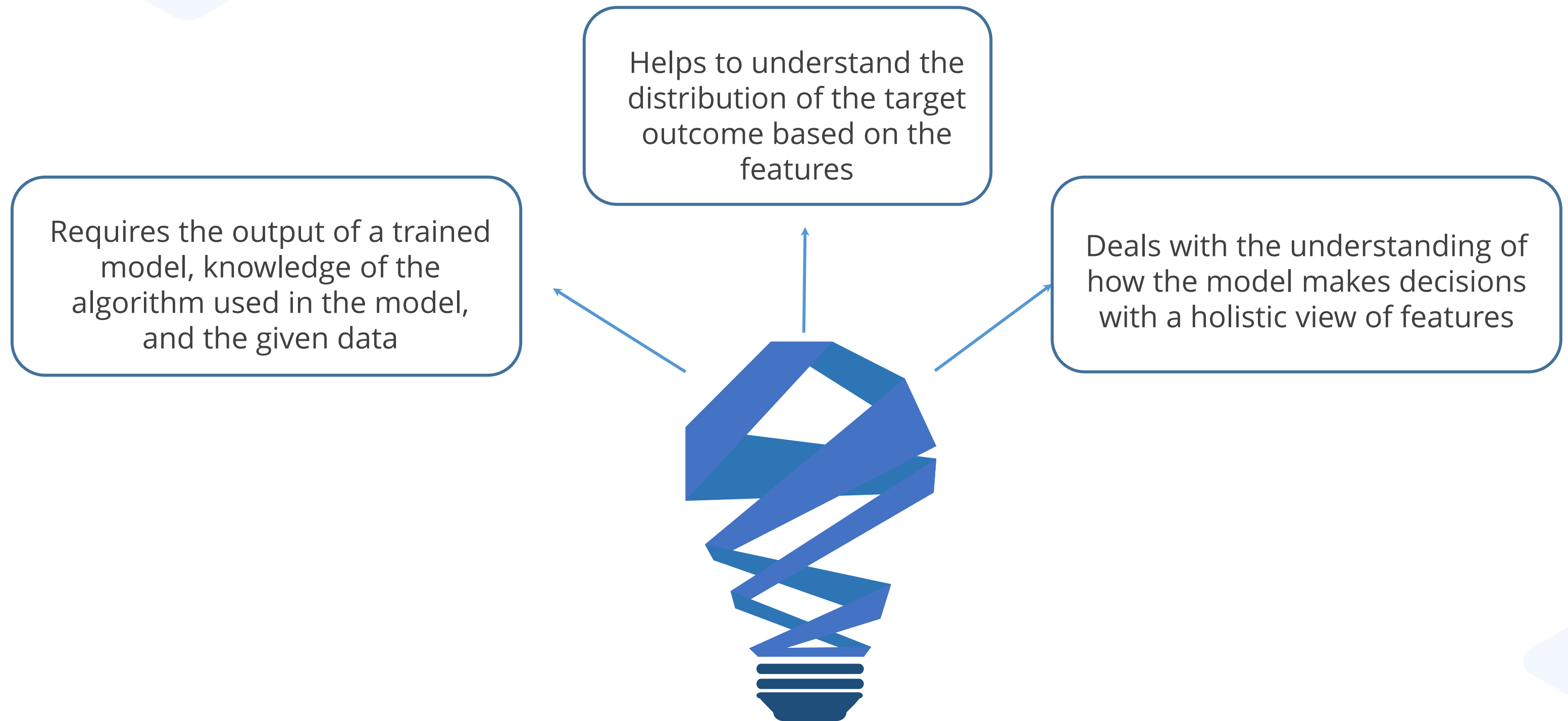
# Algorithm Transparency

Deals with how the algorithm learns a model and the types of relationships identified in the given data

Requires knowledge of the algorithm and not of the data or trained model



# Global, Holistic Model Interpretability





# Global Model Interpretability on a Modular Level

Can be used when there is difficulty achieving global model interpretability

Can be understood through the average effects of parameters and features on predictions



# Local Interpretability for a Single Prediction

Examines a single instance of a model and its prediction for the specific input

The accuracy of local interpretability is more important than the prediction of global interpretability



# Local Interpretability

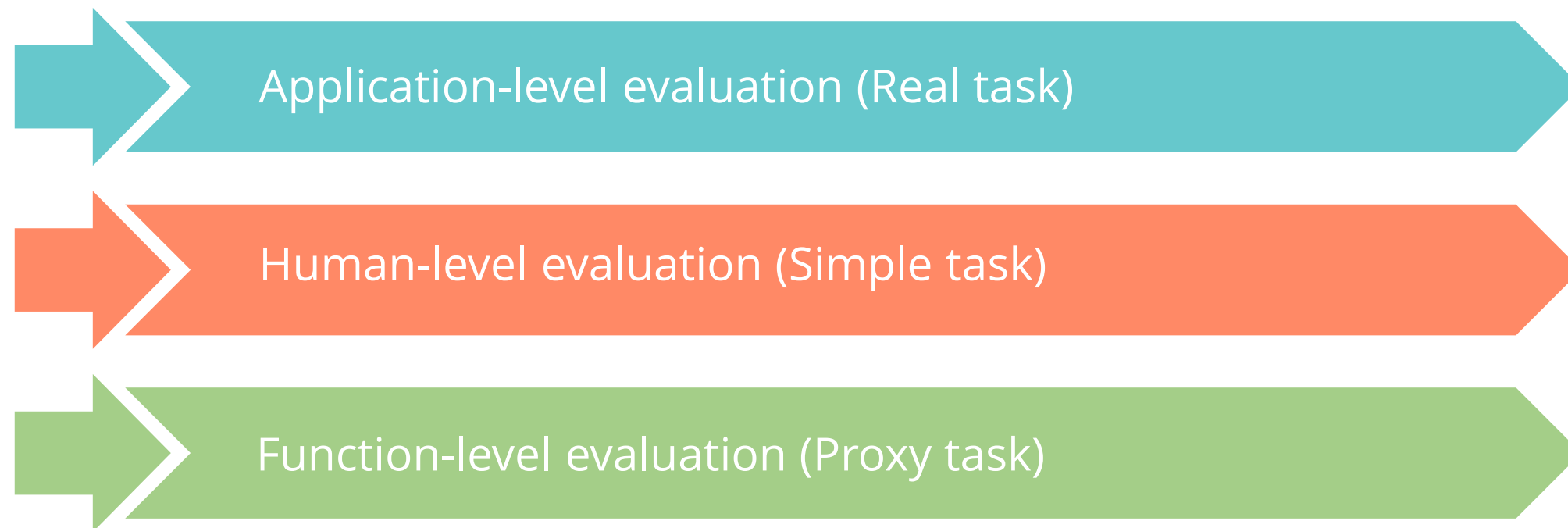
Applies global methods to a group of instances, considering the group as a complete dataset

Uses individual explanation methods on each instance and aggregates the entire group of instances



# Evaluation of Interpretability

Doshi-Velez and Kim (2017) propose three main levels for the evaluation of interpretability, which are:



# Application-Level Evaluation (Real Task)

Is the assessment of interpretability as an outcome by domain experts

Requires a good experimental setup and an understanding of quality assessment



# Human-Level Evaluation (Simple Task)

Is a simplified version of application-level evaluation

Is an inexpensive method as the evaluation does not require technical expertise



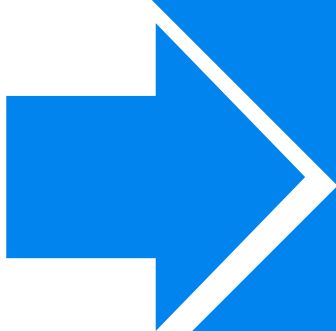
## Function-Level Evaluation (Proxy Task)

Does not require human  
expertise

Is generally performed after  
human-level evaluation, which  
leads to enhanced results



# Explainability



It pertains to the ability to explain the decision-making process of an AI model in terms understandable to the end user. An explainable model provides a clear and intuitive explanation of the decisions made



Model Explainability is crucial while debugging a model during the development phase.



# Effectiveness of Explainability

These properties measure the effectiveness of the explanation method:

Expressive power	Is a language structure generated from the explanation method
Translucency	Describes how the model relies on its parameters
Probability	Describes the explanation method suitable for the range of models
Algorithmic complexity	Checks if all relationships picked up are causal

# Effectiveness of Explainability

Accuracy

Assesses how the explanation predicts the unseen data

Fidelity

Checks how the explanation approximates the prediction

Consistency

Helps differentiate among models trained on the same dataset with the same procedure

Stability

Highlights the similar parameters in a fixed model

Comprehensibility

Helps in making the explanation understandable

# Interpretability v/s Explainability

Interpretability is the degree to which an observer can understand the cause of a decision. It is the success rate that humans can predict for the result of an AI output, while explainability goes a step further and looks at how the AI arrived at the result.

By distinguishing between interpretability and explainability, one can better evaluate and select appropriate methods for model assessment, ensuring transparency, trust, and compliance in machine learning applications.

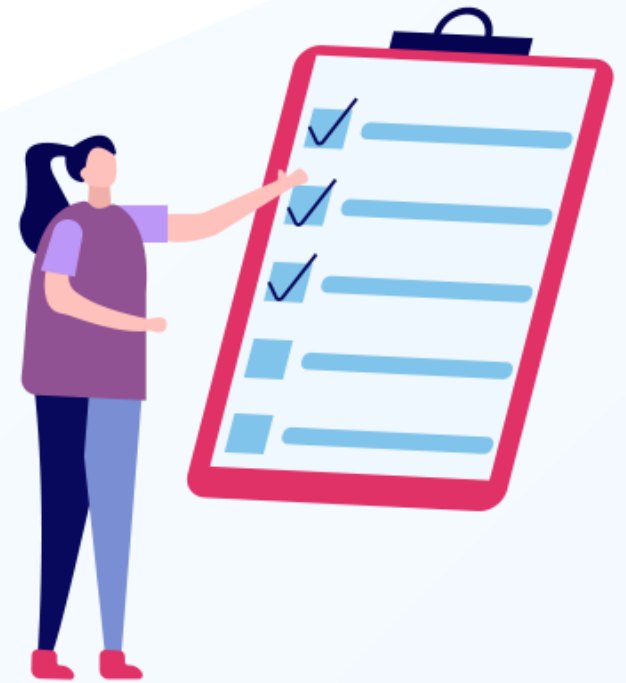
# Key Takeaways

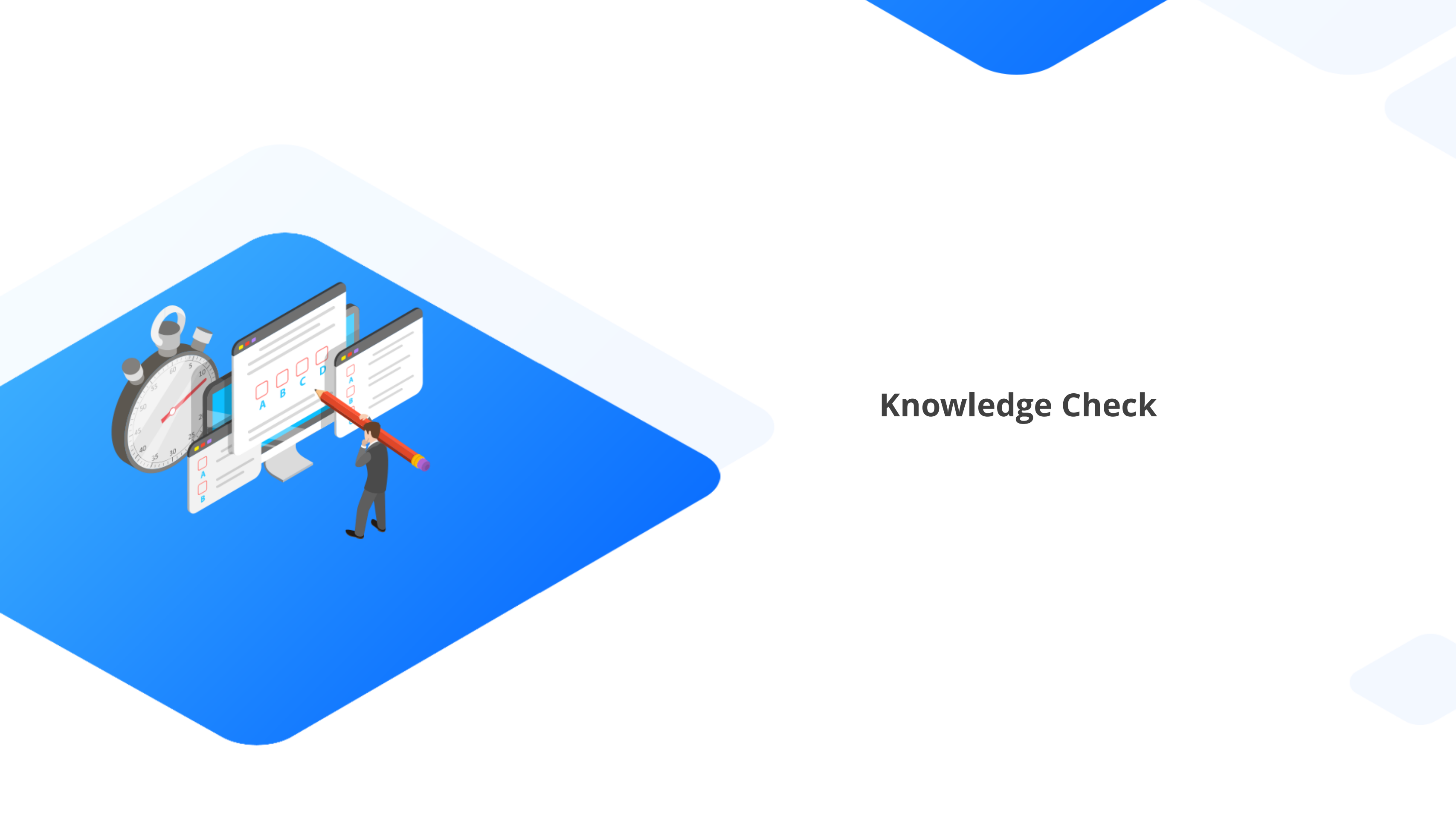
- Optimization algorithms change the attributes of the neural network to reduce losses.
- The standard gradient descent algorithm updates the parameters by evaluating the loss and gradient over the entire training dataset, leading to optimal solutions.
- Adaptive gradient (AdaGrad) optimization iteratively updates different learning rates for each parameter without manual tuning.
- Adadelat alters the custom step size calculation and removes the need for an initial learning rate hyperparameter.



# Key Takeaways

- Batch normalization normalizes output data from the model's activation functions for layers.
- Dropout and early stopping are the regularization strategies that reduce overfitting.
- Interpretability is the degree of a human's ability to predict the model's result consistently.





## Knowledge Check

## Knowledge Check

1

### What are optimizers in deep learning?

- A. Algorithms or methods used to change the attributes of the neural network, such as weights and learning rate, to reduce losses
- B. Algorithms or methods used to generate synthetic data for training
- C. Algorithms or methods used to evaluate the accuracy of the model
- D. Algorithms or methods used to initialize the weights of the model



## Knowledge Check

1

### What are optimizers in deep learning?

- A. Algorithms or methods used to change the attributes of the neural network, such as weights and learning rate, to reduce losses
- B. Algorithms or methods used to generate synthetic data for training
- C. Algorithms or methods used to evaluate the accuracy of the model
- D. Algorithms or methods used to initialize the weights of the model



---

The correct answer is **A**

---

**Optimizers are algorithms or methods to change the attributes of the neural network, such as weights and learning rates, to reduce losses.**



## Knowledge Check

2

**What are hyperparameters in a machine learning model?**

- A. The parameters that are learned by the model during training
- B. The parameters that are used to evaluate the model after training
- C. The parameters that are set manually before training the model
- D. The parameters that are used to test the model's accuracy



## Knowledge Check

2

**What are hyperparameters in a machine learning model?**

- A. The parameters that are learned by the model during training
- B. The parameters that are used to evaluate the model after training
- C. The parameters that are set manually before training the model
- D. The parameters that are used to test the model's accuracy



---

The correct answer is **C**

---

**Hyperparameters are the parameters that are set manually before training the model.**

# Lesson-End Project: Hyperparameter Tuning with MNIST Dataset



**Problem statement:** A classification model has been created using inbuilt optimizers of deep learning frameworks, but it is not producing the desired output. You are encountering the same issue and need to perform tuning using the MNIST dataset.

## **Objective:**

To build a single-layer dense neural network to perform hyperparameter tuning using random search with a **KerasTuner**

**Access:** Click on the **Lab** tab on the left side of the LMS panel. Copy the generated username and password. Click on the **Launch Lab** button. On the new page, enter the username and password you copied earlier into the respective fields. Click **Login** to start your lab session. A full-fledged Jupyter lab opens, which you can use for your hands-on practice and projects.



**Thank You!**