

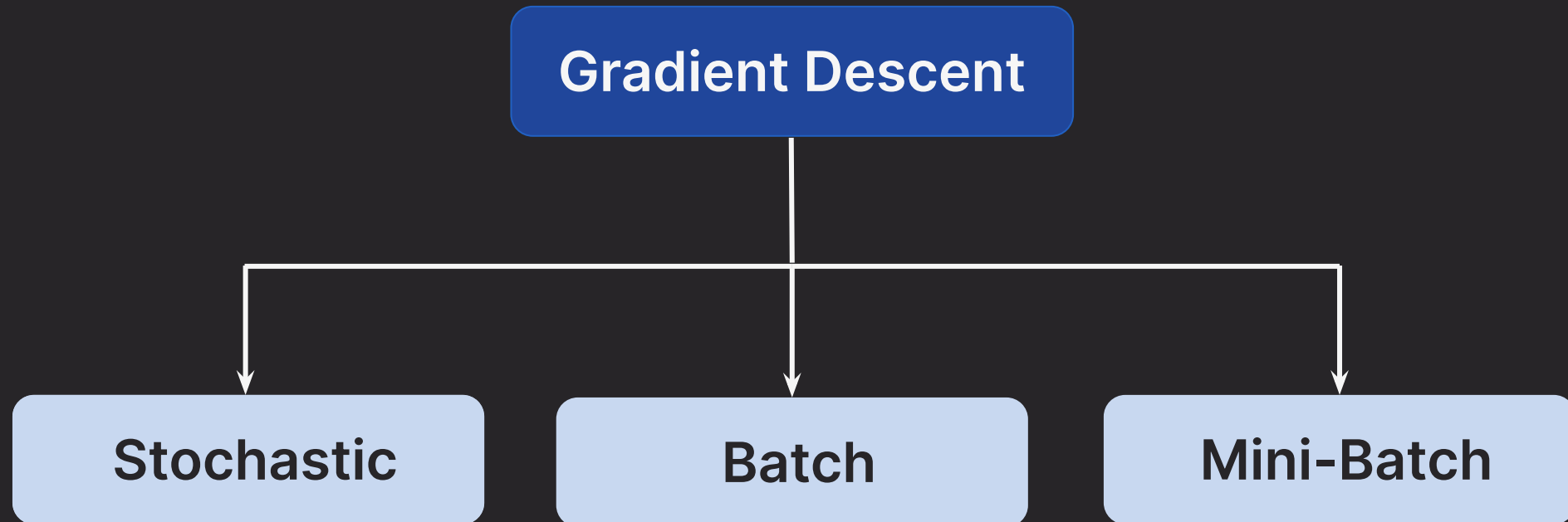


Understanding workings of Neural Networks

Video 5: Types of Gradient Descent

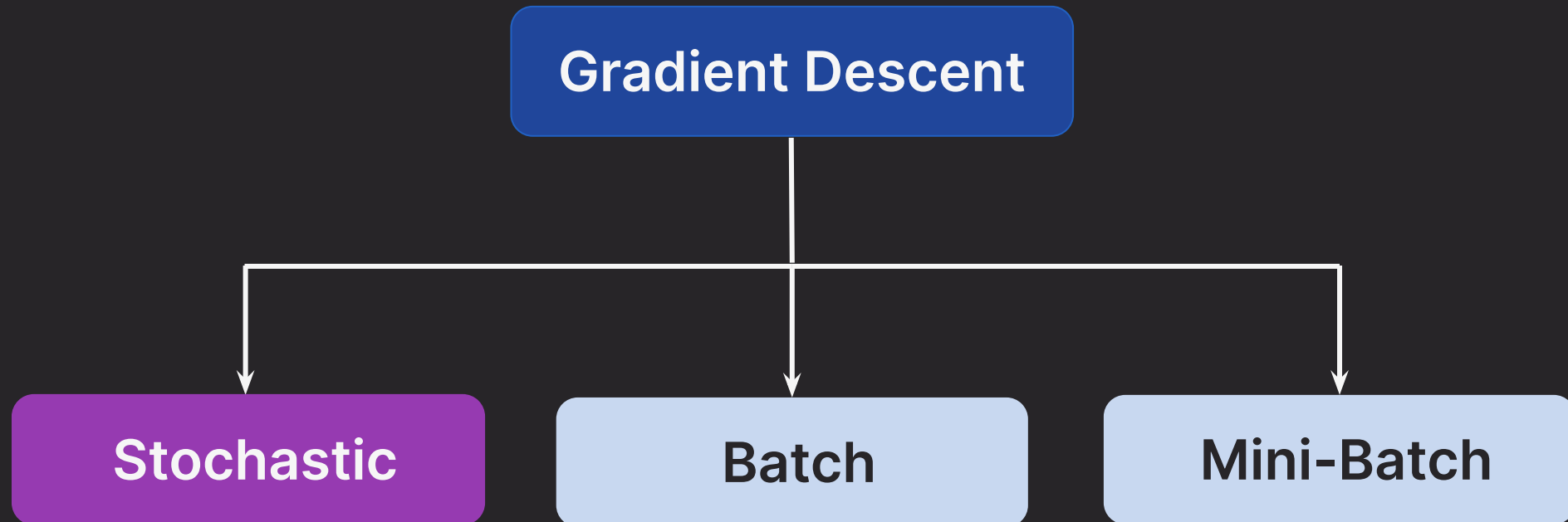
Gradient Descent: Classification

Categorized based on amount of data to compute the gradient of loss function



Gradient Descent: Classification

Categorized based on amount of data to compute the gradient of loss function

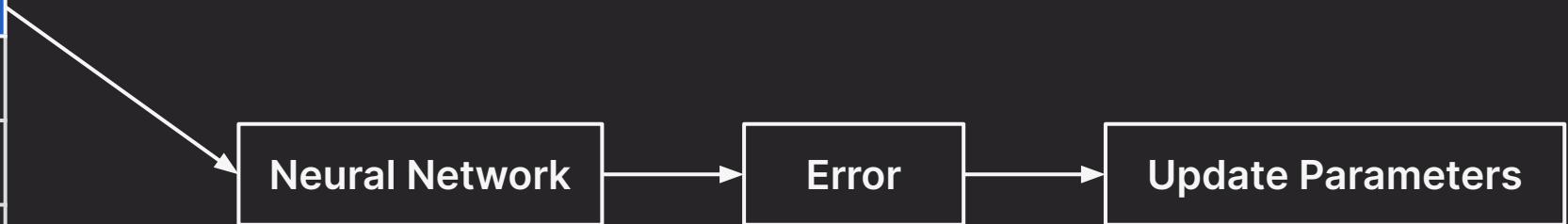


Stochastic Gradient Descent

Features	
Observations	sat_score
	727
	722
	716
	...
	702
Observations	hs_gpa
	3.40
	4.00
	3.75
	...
	3.53

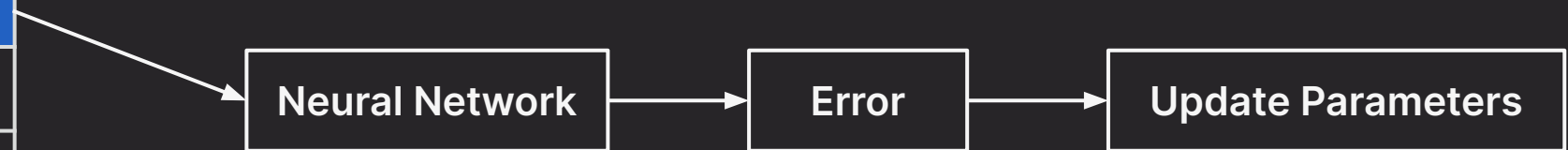
Stochastic Gradient Descent

sat_score	hs_gpa
727	3.40
722	4.00
716	3.75
...	...
702	3.53



Stochastic Gradient Descent

sat_score	hs_gpa
727	3.40
722	4.00
716	3.75
...	...
702	3.53

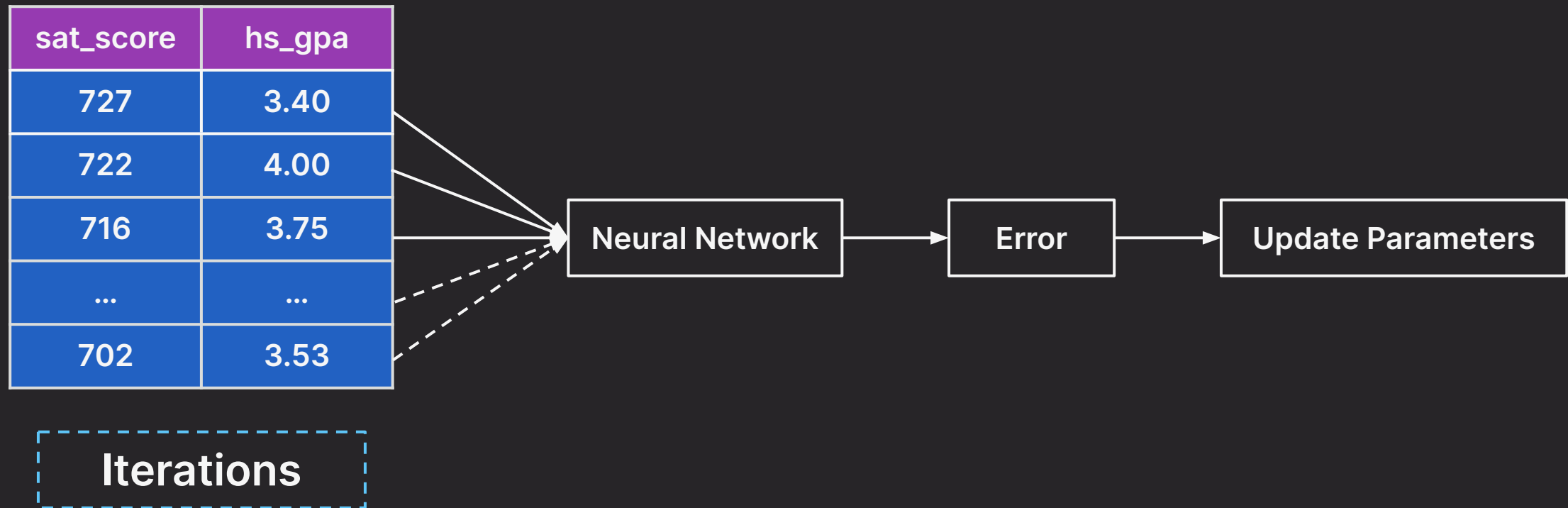


Stochastic Gradient Descent

sat_score	hs_gpa
727	3.40
722	4.00
716	3.75
...	...
702	3.53



Stochastic Gradient Descent



Stochastic Gradient Descent

'm' iterations per epoch

m _i Observations	sat_score	hs_gpa	fy_gpa
	727	3.40	3.18
	722	4.00	3.33
	716	3.75	3.25

	m _i	m _i	m _i

Stochastic Gradient Descent

One Epoch



Each sample has gone through one pass of forward and backward propagation.

Stochastic Gradient Descent

m_i Observations	sat_score	hs_gpa	fy_gpa
	727	3.40	3.18
	722	4.00	3.33
	716	3.75	3.25

	m_i	m_i	m_i

SGD Updation = $50 \times 10 = 500$ times

Stochastic Gradient Descent

Recommended Practice:



Add a shuffle parameter before selection of the records

Stochastic Gradient Descent: Advantage



Improved model performance because of the frequent updation of weights and biases

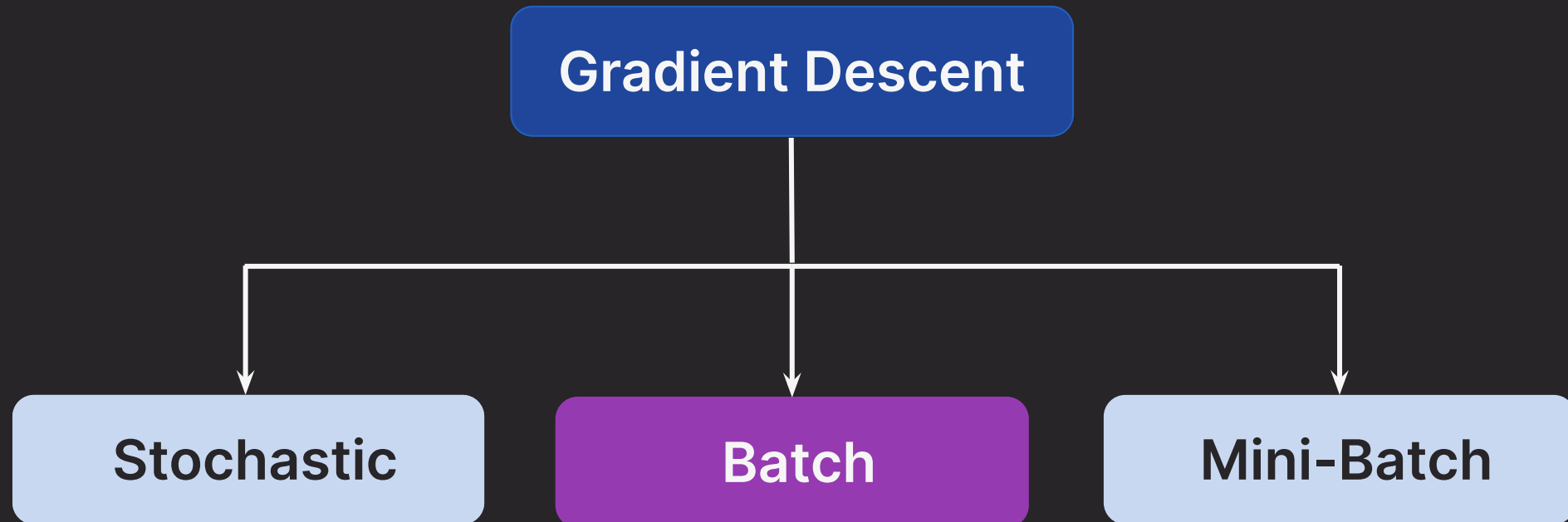
Stochastic Gradient Descent: Advantage



Updating the weights based on just one sample
can result in fluctuations in loss values

Gradient Descent: Classification

Categorized based on amount of data to compute the gradient of loss function



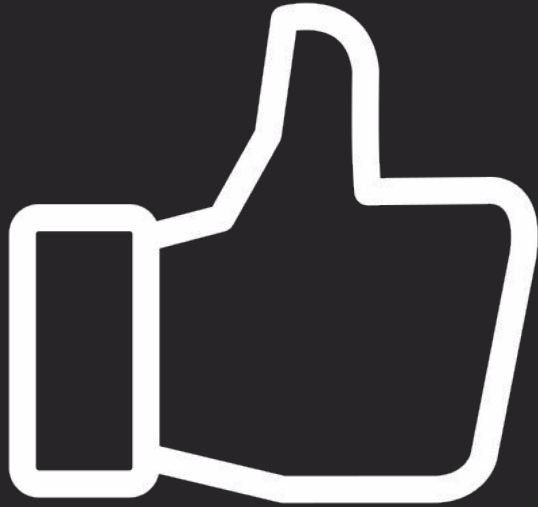
Batch Gradient Descent

50 Observations

sat_score	hs_gpa	fy_gpa
727	3.40	3.18
722	4.00	3.33
716	3.75	3.25
...
...
710	3.15	3.05

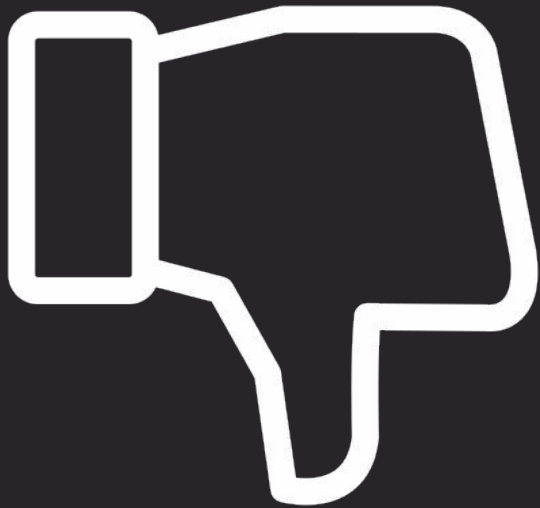
Batch Gradient Descent Updation = 10 times for 10 epochs

Batch Gradient Descent: Advantage



- Computational Efficiency
- Stable Error Gradient
- Faster Convergence

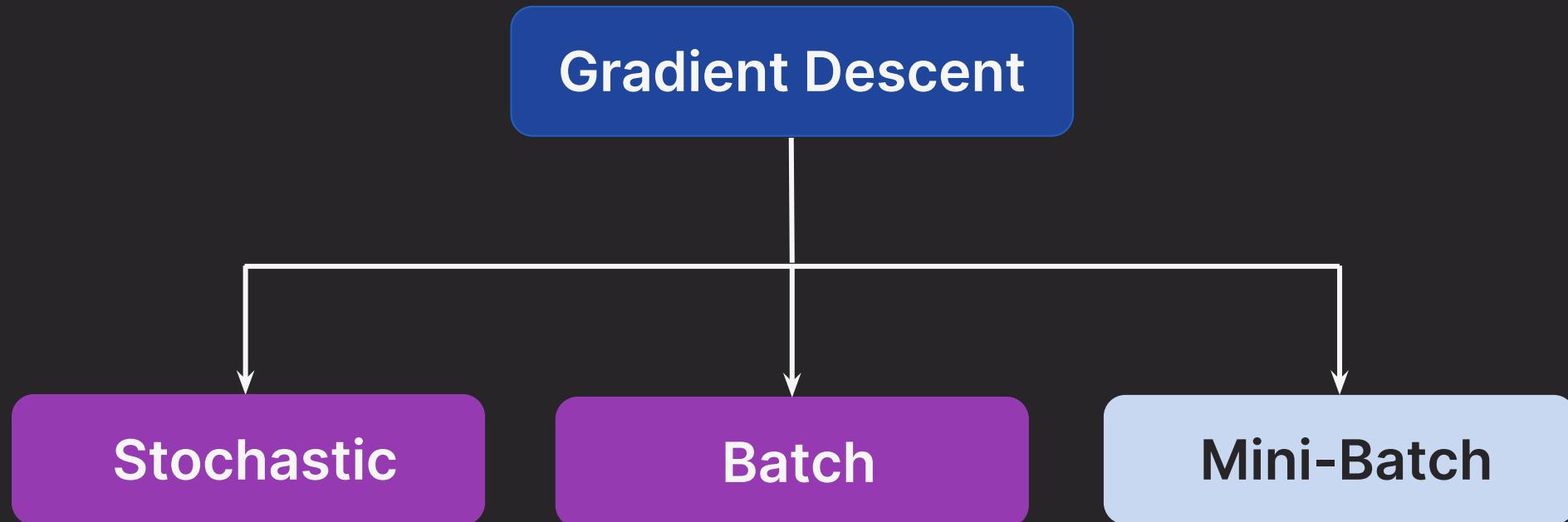
Batch Gradient Descent: Drawback



- Would need several epochs for training
- Requires the full dataset in memory, limiting scalability

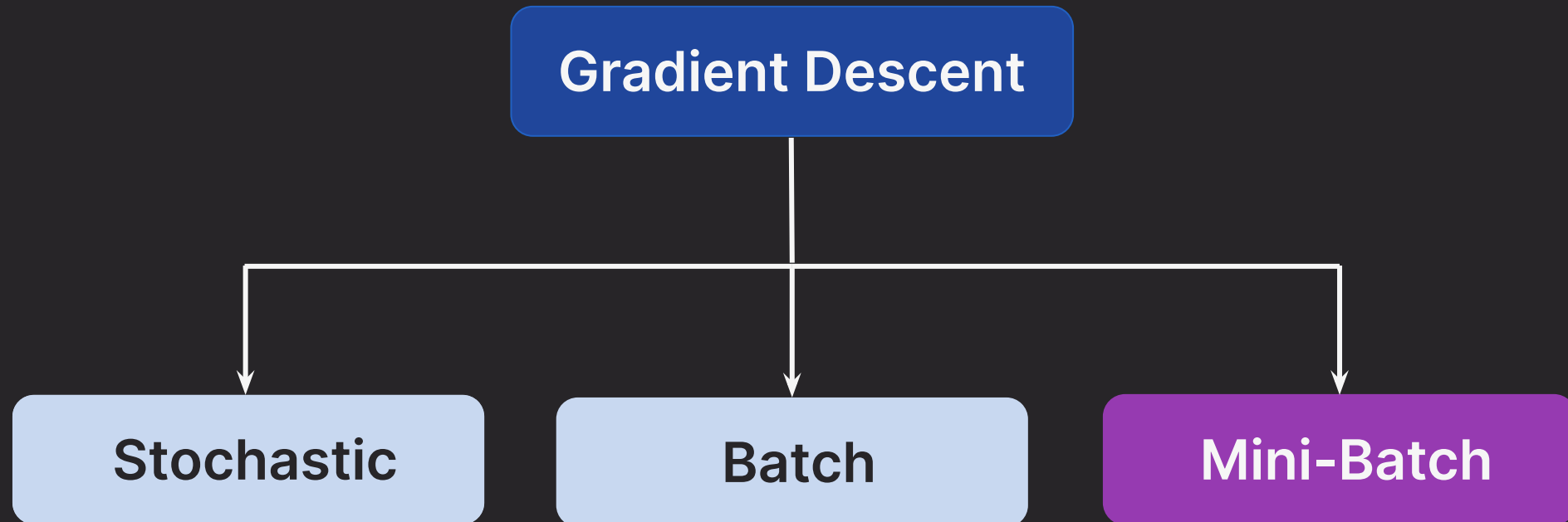
Gradient Descent: Classification

Categorized based on amount of data to compute the gradient of loss function



Gradient Descent: Classification

Categorized based on amount of data to compute the gradient of loss function



Mini-Batch Gradient Descent

sat_score	hs_gpa	fy_gpa	Subset 1
727	3.40	3.18	
722	4.00	3.33	
716	3.75	3.25	
...	Subset n
...	
710	3.15	3.05	

Mini-Batch Gradient Descent

sat_score	hs_gpa	fy_gpa	Subset 1
727	3.40	3.18	
722	4.00	3.33	
716	3.75	3.25	
...	Subset n
...	
710	3.15	3.05	

50 samples / 10 = 5 subsets

Number of epochs = 5

Mini-Batch Gradient Descent

800 Observations

sat_score	hs_gpa	fy_gpa
727	3.40	3.18
722	4.00	3.33
716	3.75	3.25
...
...
710	3.15	3.05

800 samples ; batch_size = 80

How many times will the weights be updated in 10 epochs?

Mini-Batch Gradient Descent

800 Observations	sat_score	hs_gpa	fy_gpa
	727	3.40	3.18
	722	4.00	3.33
	716	3.75	3.25

	710	3.15	3.05

800 samples ; batch_size=80 → 10 subsets

10 subsets * 10 epochs = 100 times

Mini-Batch Gradient Descent

Observations	sat_score	hs_gpa	fy_gpa	Subset 1
	727	3.40	3.18	
	722	4.00	3.33	
	716	3.75	3.25	Subset n
	
	
	710	3.15	3.05	

Mini batch size of 2^n is preferred
eg: 16, 32, 64, 128, 256, 512, 1024.....

Hands-on