# LECTURE 3:

# OPTIMIZATION IN DEEP LEARNING

University of Washington, Seattle
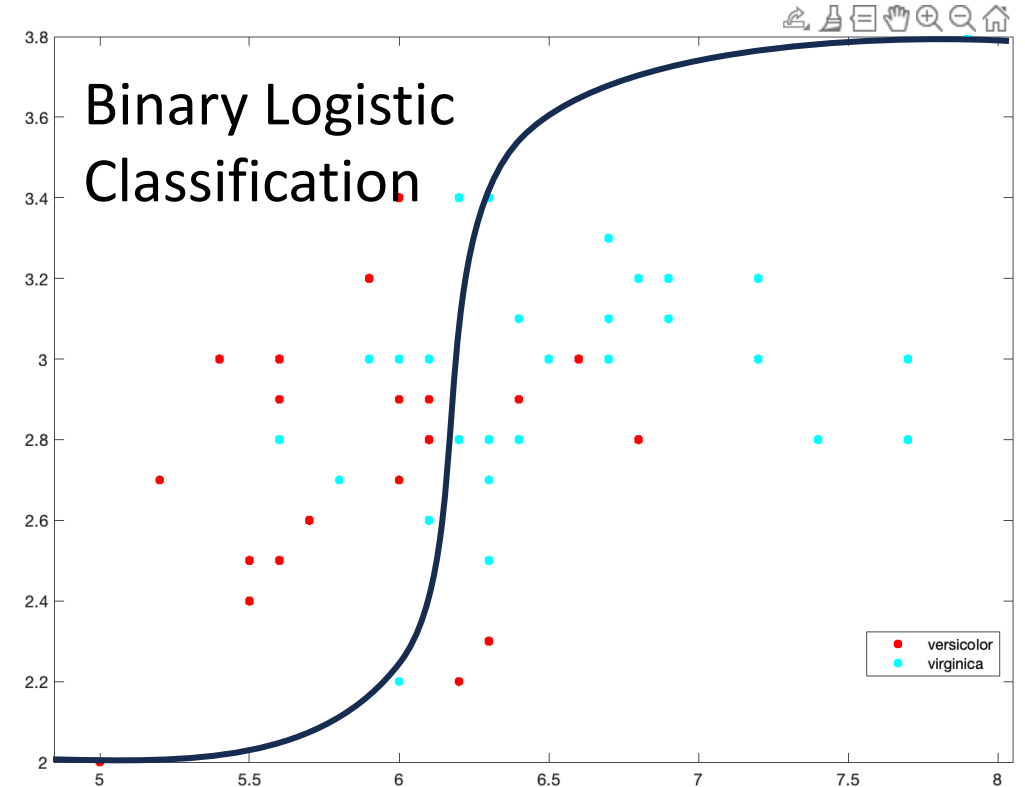
Fall 2024

# Previously in EEP 596…
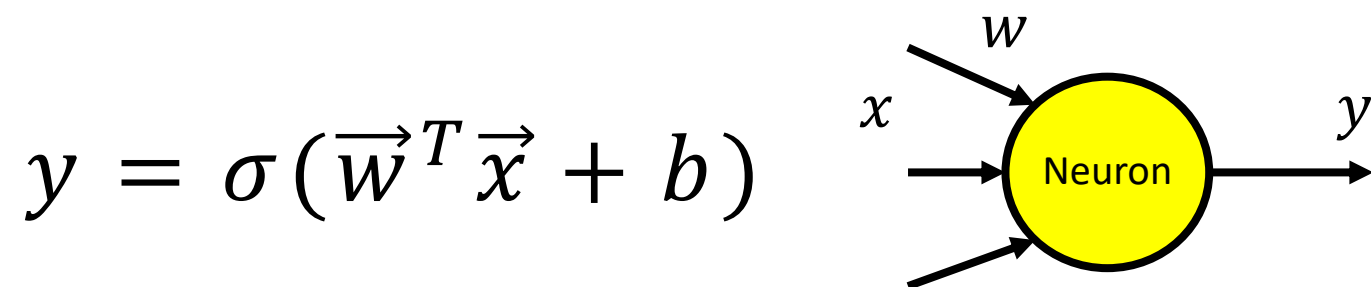


Linear regression

Binary Logistic Classification

$$J = E_2 = \sum_{i=1}^{n} (\vec{p} \cdot \vec{x}_i - y_i)^2 \quad \vec{p}^* = \left( X^\top X \right)^{-1} X^\top \vec{y}$$

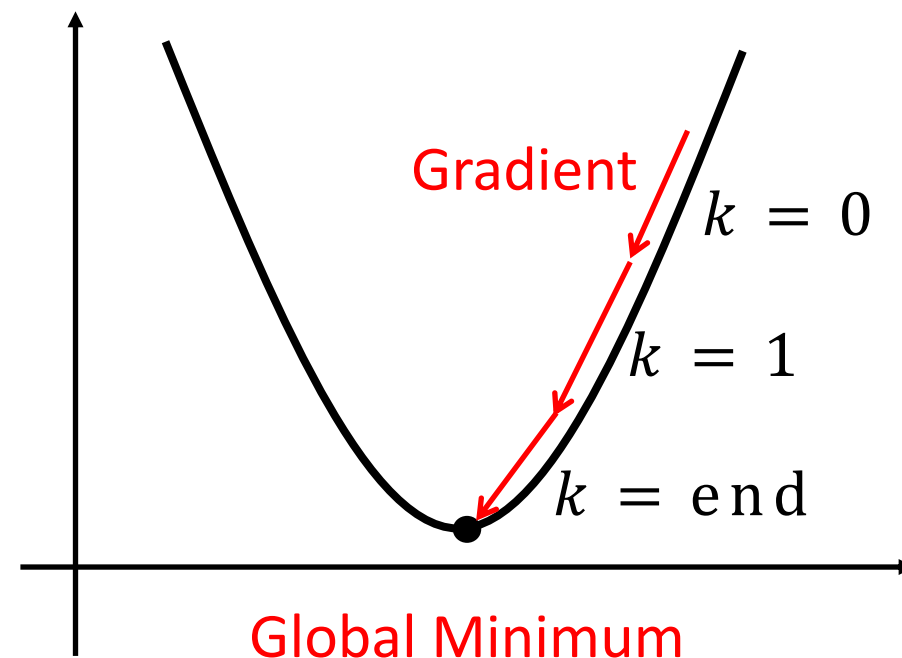$$p = \frac{1}{1 + e^{-x}} \quad 0 \le p \le 1$$

$$y = \sigma(\vec{w}^T \vec{x} + b)$$



$$\vec{w}_{k+1} = \vec{w}_k - \alpha \nabla_{\vec{w}} J(\vec{w}_k; b)$$

$$b_{k+1} = b_k - \alpha \frac{\partial}{\partial b} J(\vec{w}; b_k)$$

$$J = L((\vec{w}, b), y)$$



Gradient
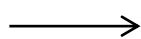
$k = 0$

$k = 1$

$k = \text{end}$

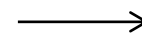Global Minimum

# Previously in EEP 596...

$$\nabla_{\vec{w}} L(\hat{y}, y) = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \nabla_{\vec{w}} z$$
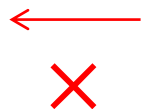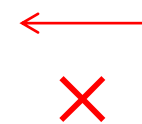
FWD

| $z = \vec{w}^T \vec{x} + b$ | $\longrightarrow$ | $\hat{y} = \sigma(z)$ | $\longrightarrow$ | $L(\hat{y}, y)$ |

BWD

| $\nabla_{\vec{w}} z$ | $\longleftarrow$ ✗ | $\frac{\partial \hat{y}}{\partial z}(z)$ | $\longleftarrow$ ✗ | $\frac{\partial L}{\partial \hat{y}}(\hat{y}, y)$ |

# OUTLINE

**Part 1: Stochastic Gradient Descent**

- GD vs SGD

- Convergence of SGD

- Learning rate and convergence

- Comparing GD variants

**Part 2: Optimizers**

- Variable learning rate

- Advanced methods

- Choosing optimizer

**Part 3: Optimization Techniques in DL**

- Cross validation

- Regularization

- Data Normalization

- Batch-normalization

- Network initialization
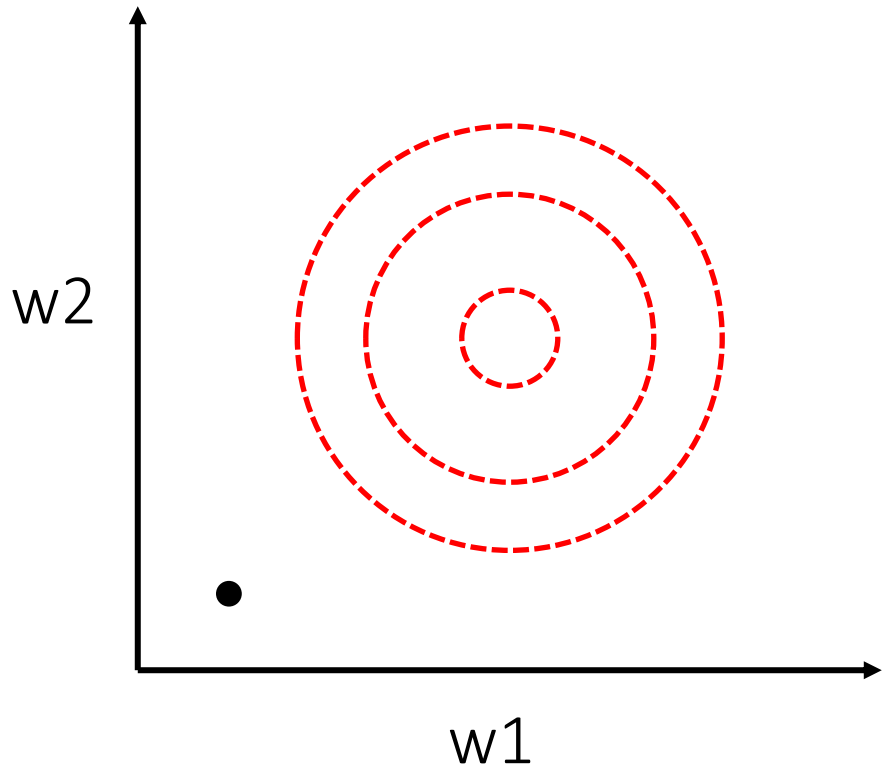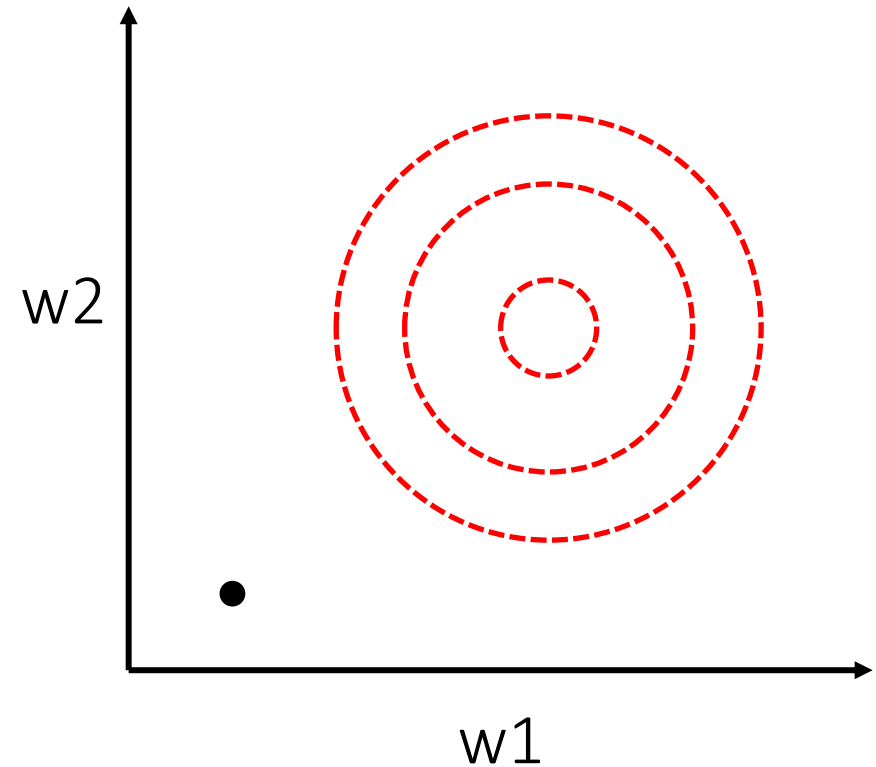
- Hyperparameter tunings

# PART 1:

# Stochastic Gradient Descent (SGD)

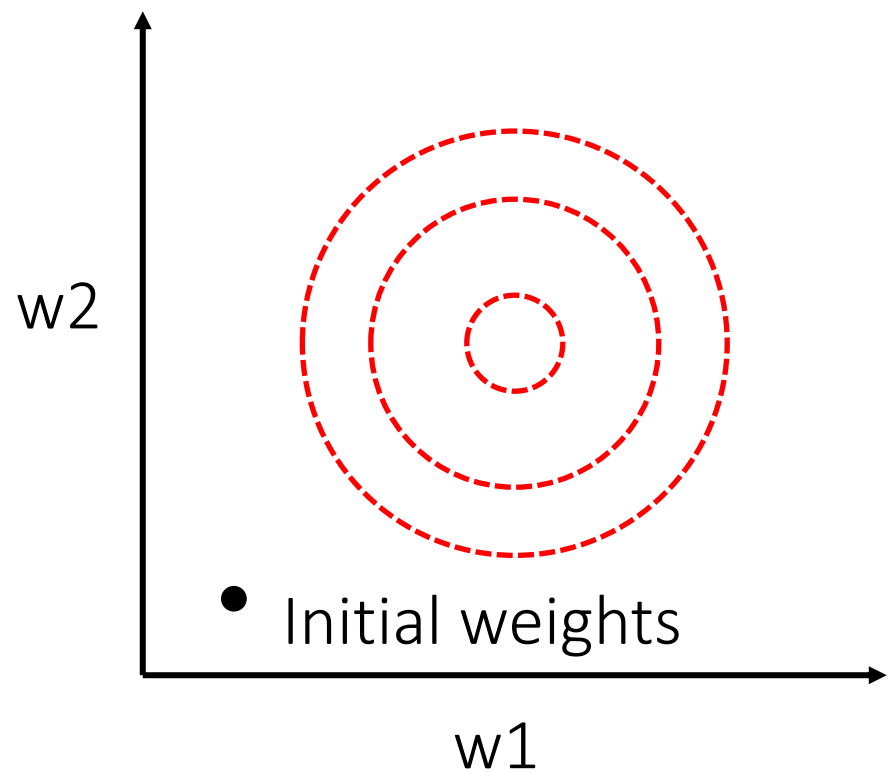# Gradient descent (GD) vs Stochastic Gradient Descent (SGD)



Batch Gradient Descent
1 iteration: FWD pass and BWD pass on
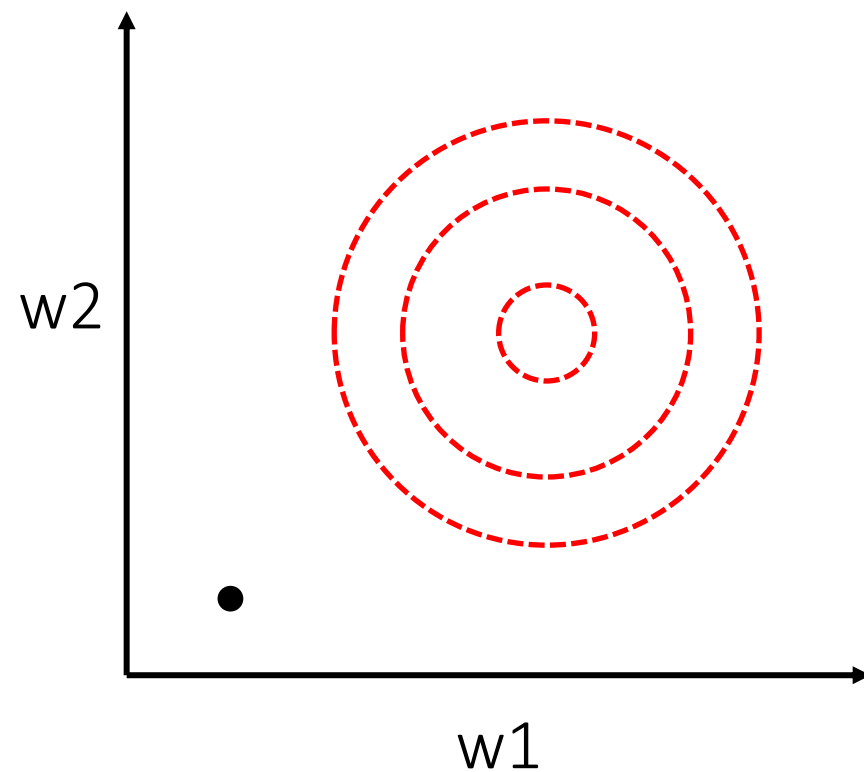**whole training set**

Stochastic Gradient Descent
1 iteration: FWD pass and BWD pass on
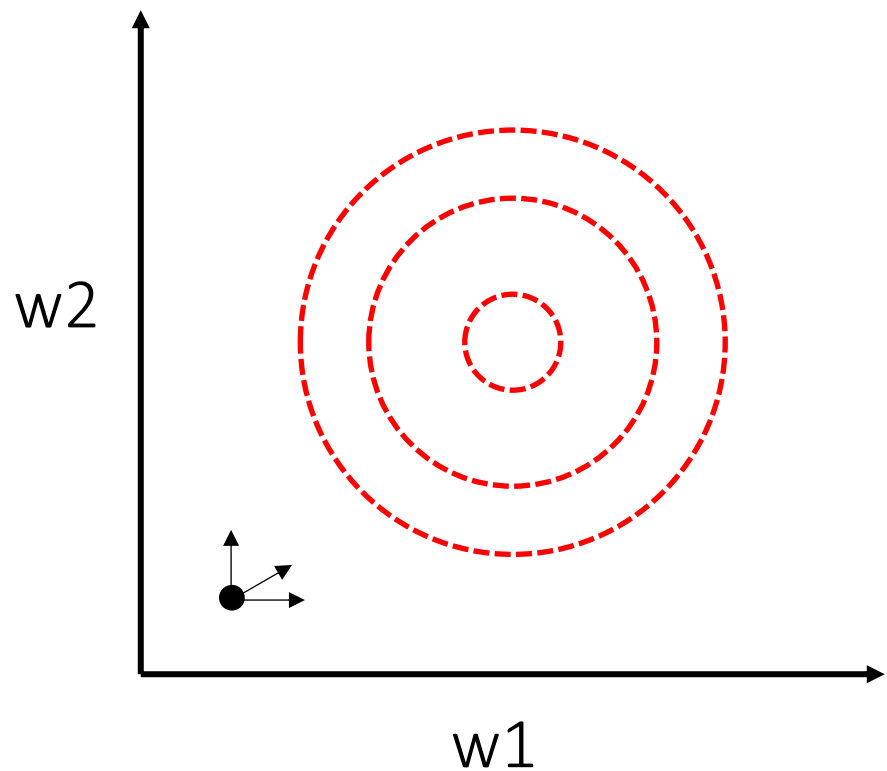**subset of training set**

# GD vs SGD



Batch Gradient Descent

Stochastic Gradient Descent

# GD vs SGD



Batch Gradient Descent

Stochastic Gradient Descent

# GD vs SGD

$\nabla \overrightarrow{w}_1$

$\nabla \overrightarrow{w}_2$

$\nabla \overrightarrow{w}_3$

w2

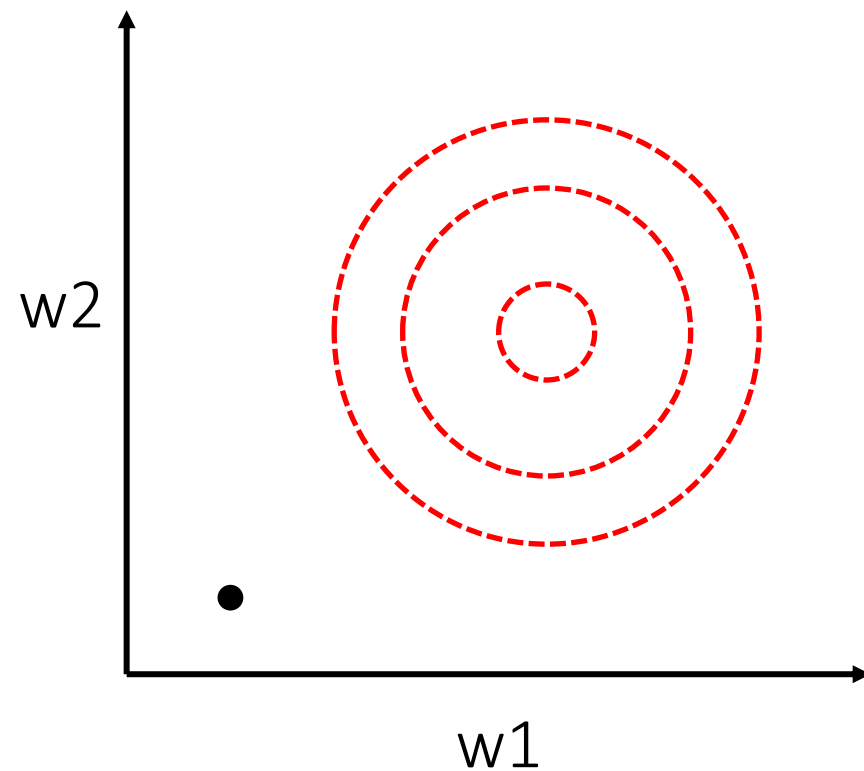w1

Batch Gradient Descent
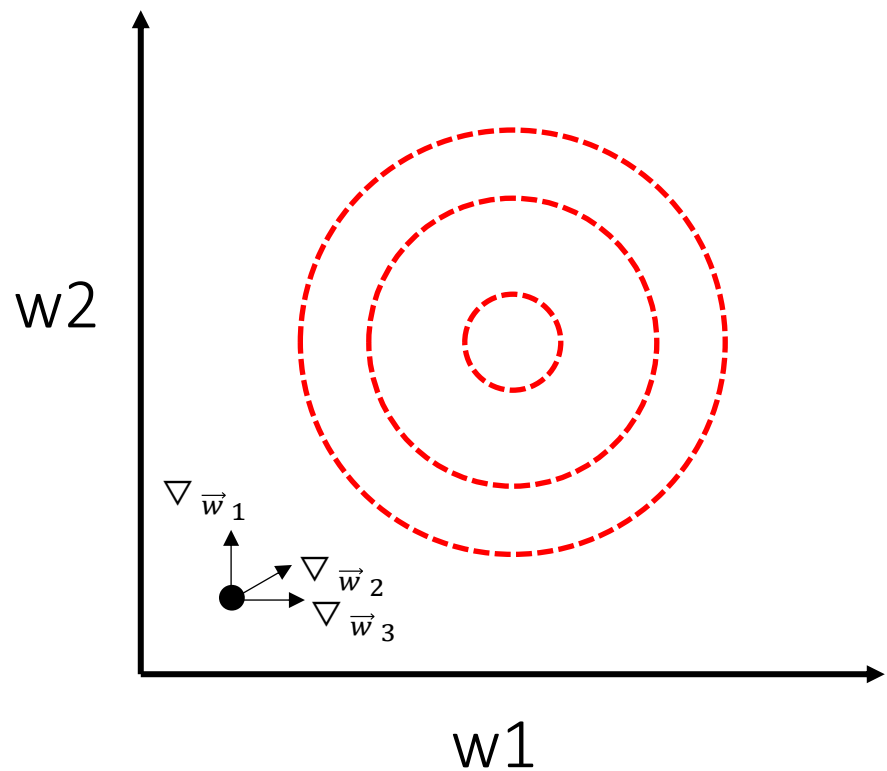
w2

w1

Stochastic Gradient Descent

# GD vs SGD



Batch Gradient Descent

Stochastic Gradient Descent

# GD vs SGD



Batch Gradient Descent

Stochastic Gradient Descent

# GD vs SGD



Batch Gradient Descent

Stochastic Gradient Descent

# GD vs SGD



Batch Gradient Descent

Stochastic Gradient Descent
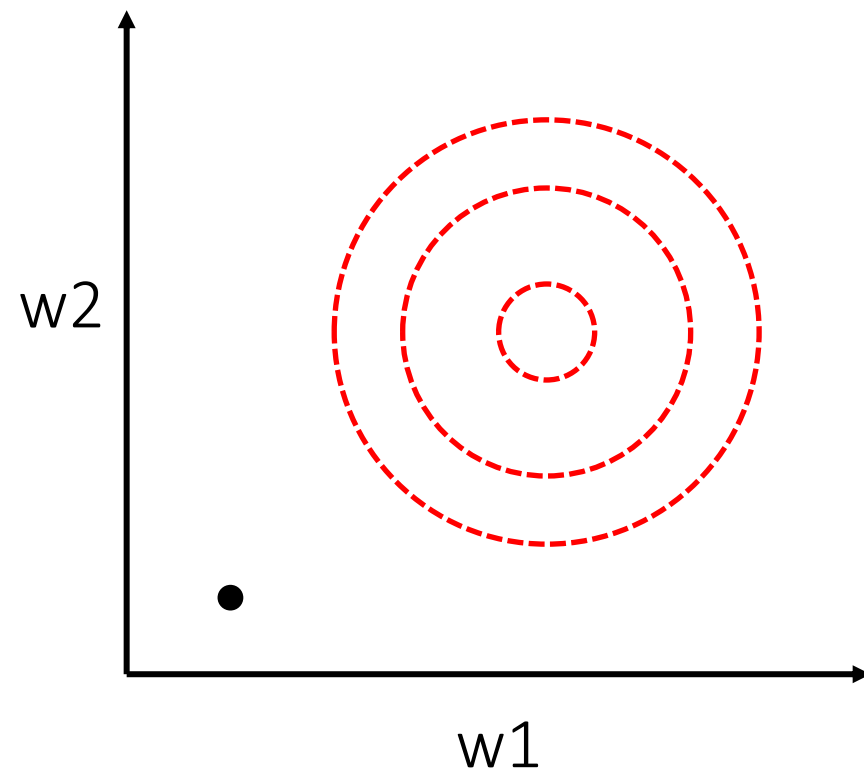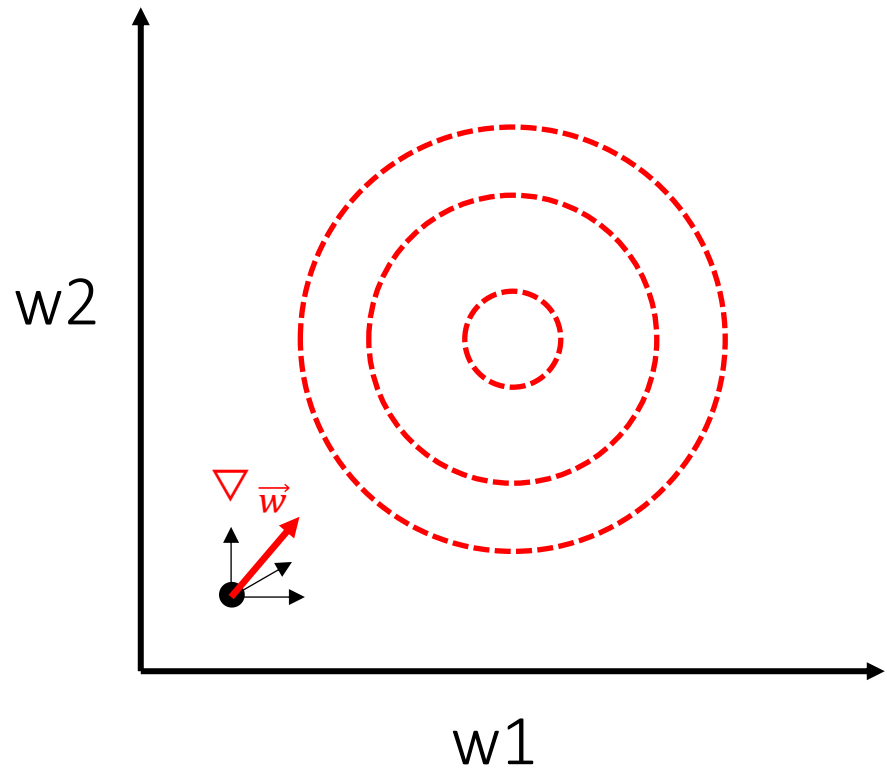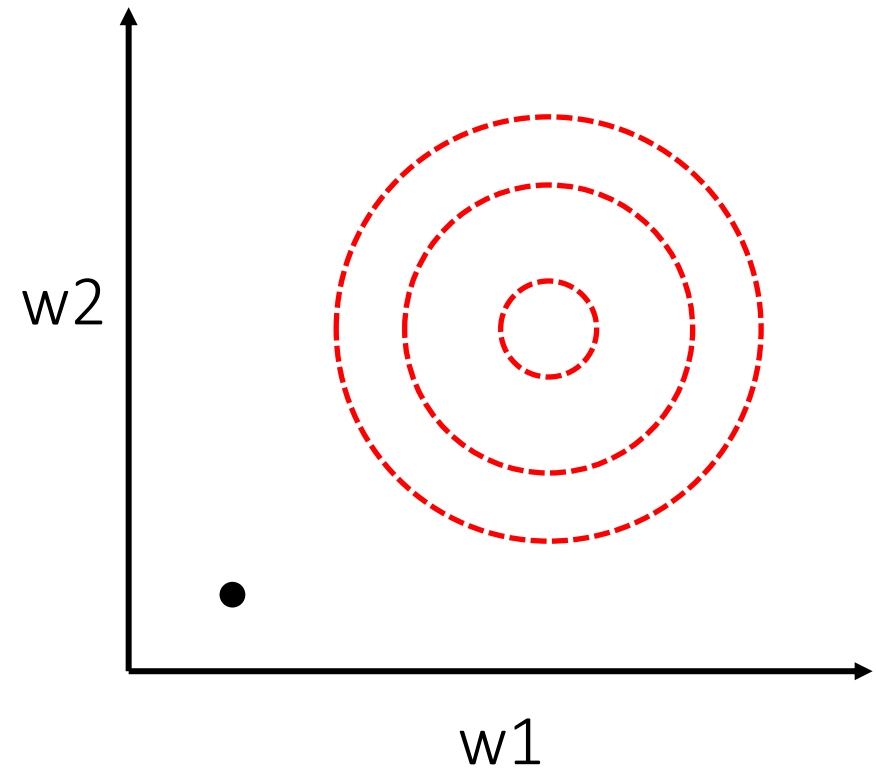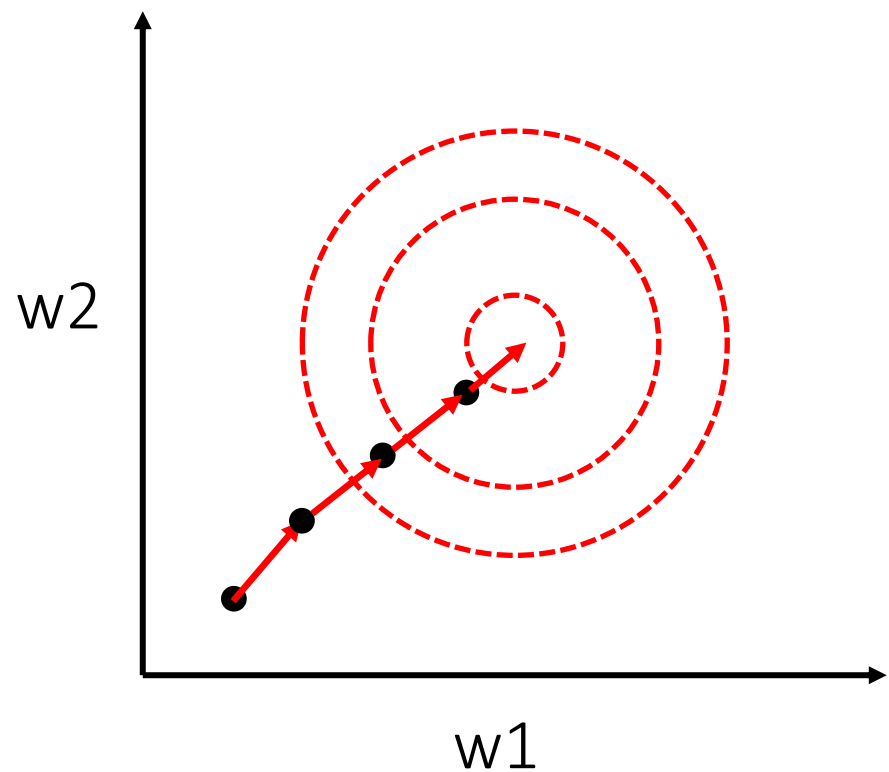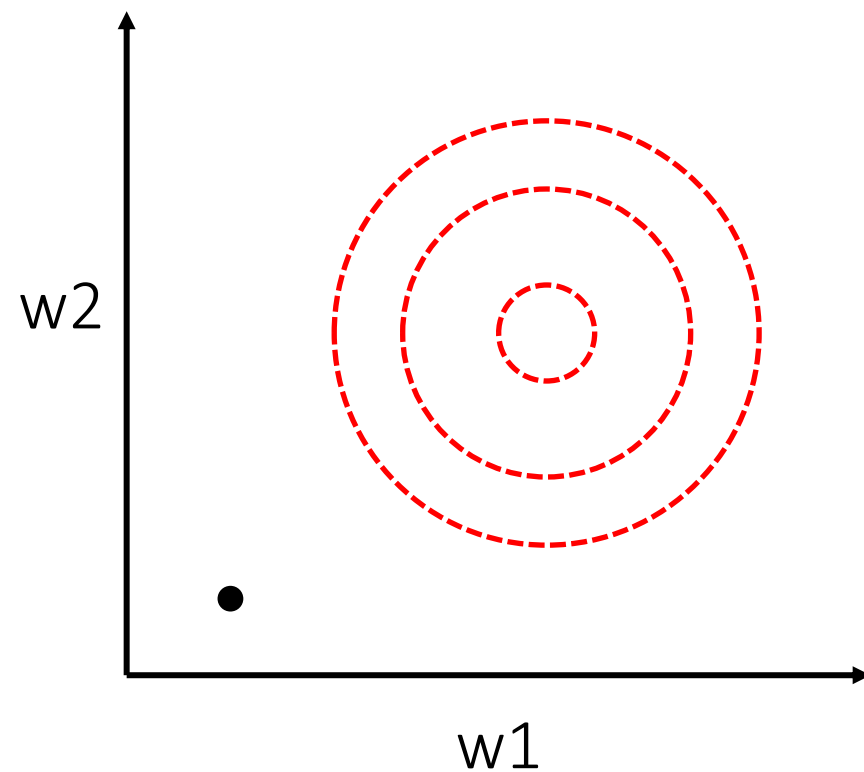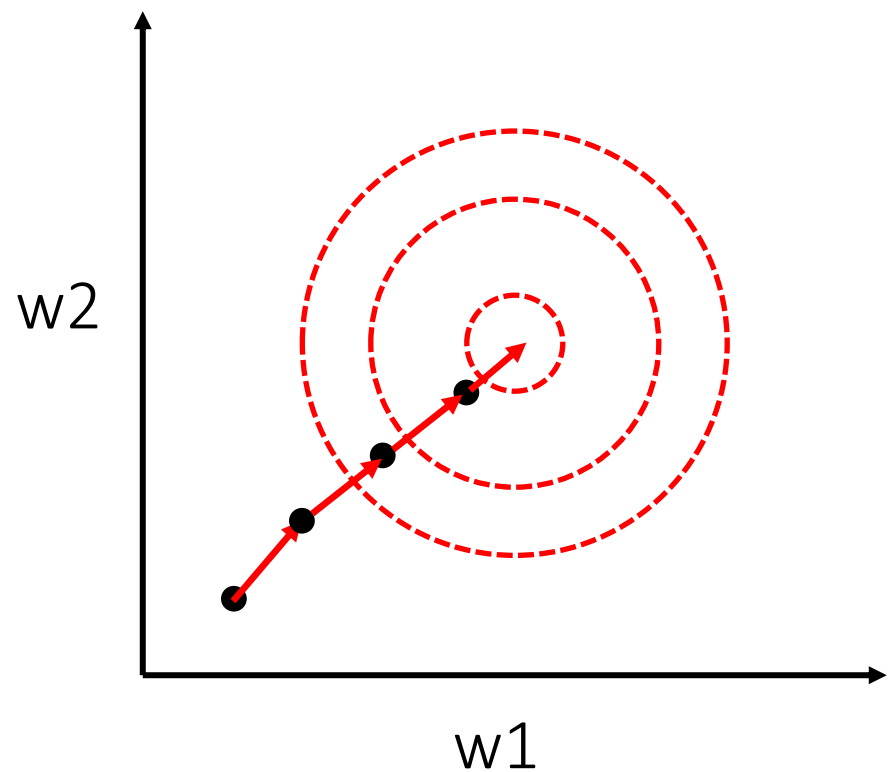
# GD vs SGD



Batch Gradient Descent
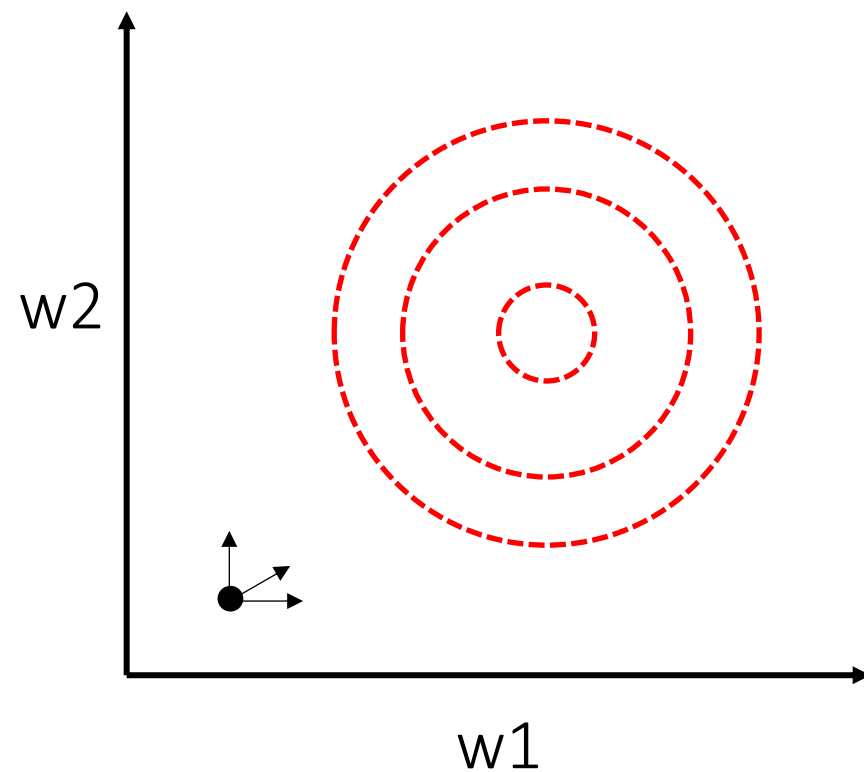
Stochastic Gradient Descent

# Pros and Cons of SGD



Pros:
Still consistently converges to minimum
May take shortcut to minimum

# Pros and Cons of SGD



Pros:
Still consistently converges to minimum
May take shortcut to minimum

Cons:
Not useful when we are already close to minimum
Hard to parallelize

# GD vs SGD

Almost surely convergence

Pros:
Still consistently converges to minimum
May take shortcut to minimum

Cons:
Not useful when we are already close to minimum
Hard to parallelize

$$\vec{w}_{k+1} = \vec{w}_k - \boxed{\alpha} \nabla_{\vec{w}} J(\vec{w}_k; b)$$

$$b_{k+1} = b_k - \boxed{\alpha} \frac{\partial}{\partial b} J(\vec{w}; b_k)$$

# Effects of learning rate ($\alpha$) on SGD





Loss curve is typically noisy with SGD

$$\vec{w}_{k+1} = \vec{w}_k - \boxed{\alpha} \nabla_{\vec{w}} J(\vec{w}_k; b)$$

$$b_{k+1} = b_k - \boxed{\alpha} \frac{\partial}{\partial b} J(\vec{w}; b_k)$$

# Effects of learning rate on SGD

|  | SGD | Mini-batch GD | Batch GD |
|---|---|---|---|
| **Data batch size per iteration** | 1 |  |  |
|  | (-) Can loose speedup from oscillations<br><br>(-) hard to parallelize |  |  |

**N = Total # of datapoints in training set**

**m = Number of mini-batches for training set**

# Effects of learning rate on SGD

| | SGD | Mini-batch GD | Batch GD |
|---|---|---|---|
| **Data batch size per iteration** | 1 | N/m | |
| | (-) Can loose speedup from oscillations<br><br>(-) hard to parallelize | (+) The whole mini-batch is evaluated in parallel<br><br>(+) Mostly consistent convergence | |

**N = Total # of datapoints in training set**

**m = Number of mini-batches for training set**

# Effects of learning rate on SGD

|  | **SGD** | **Mini-batch GD** | **Batch GD** |
|---|---|---|---|
| **Data batch size per iteration** | 1 | N/m | n |
|  | (-) Can loose speedup from oscillations<br><br>(-) hard to parallelize | (+) The whole mini-batch is evaluated in parallel<br><br>(+) Mostly consistent convergence | (+) Consistent convergence<br><br>(+) Maximum parallelization<br><br>(-) Too long per iteration<br><br>(-) Hardware memory limit (RAM, VRAM) |

**N = Total # of datapoints in training set**

**m = Number of mini-batches for training set**
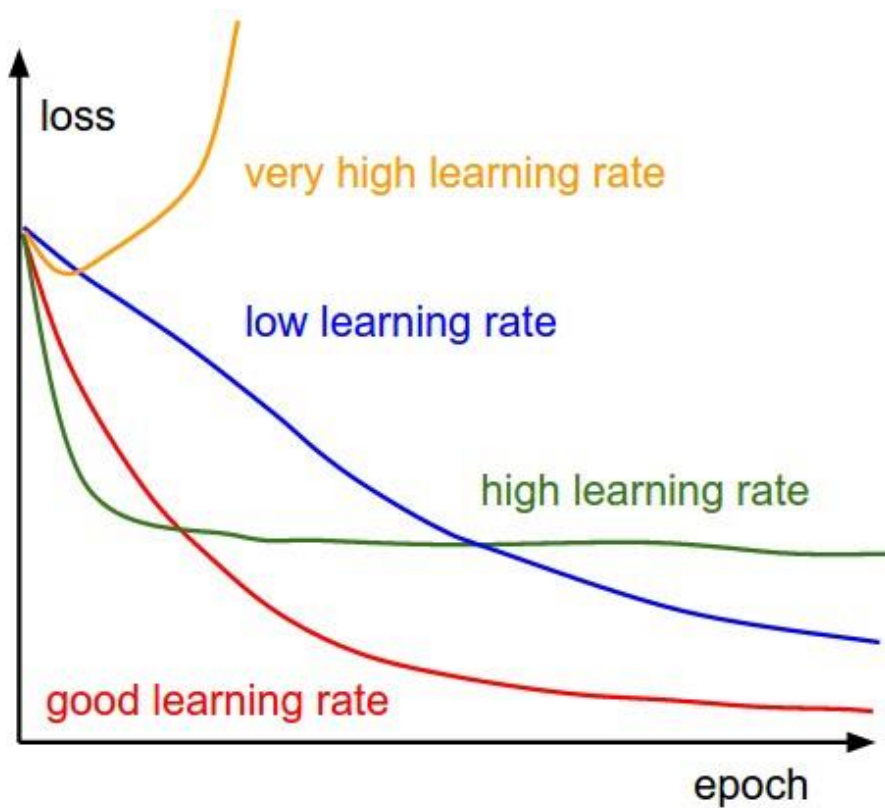
# PART 2:
# Optimizers in Deep Learning

# Variable Learning Rates

$$\vec{w}_{k+1} = \vec{w}_k - \alpha \nabla_{\vec{w}} J(\vec{w}_k; b)$$

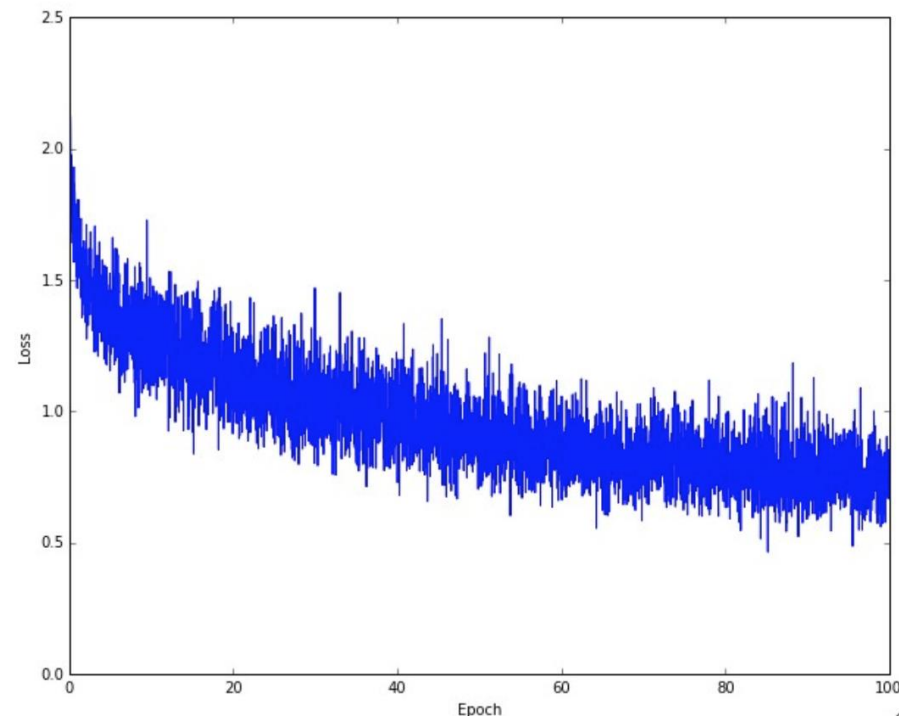$$b_{k+1} = b_k - \alpha \frac{\partial}{\partial b} J(\vec{w}; b_k)$$

$$J = L((\vec{w}, b), y)$$

Gradient

$k = 0$

$k = 1$

$k = \text{end}$

Global Minimum
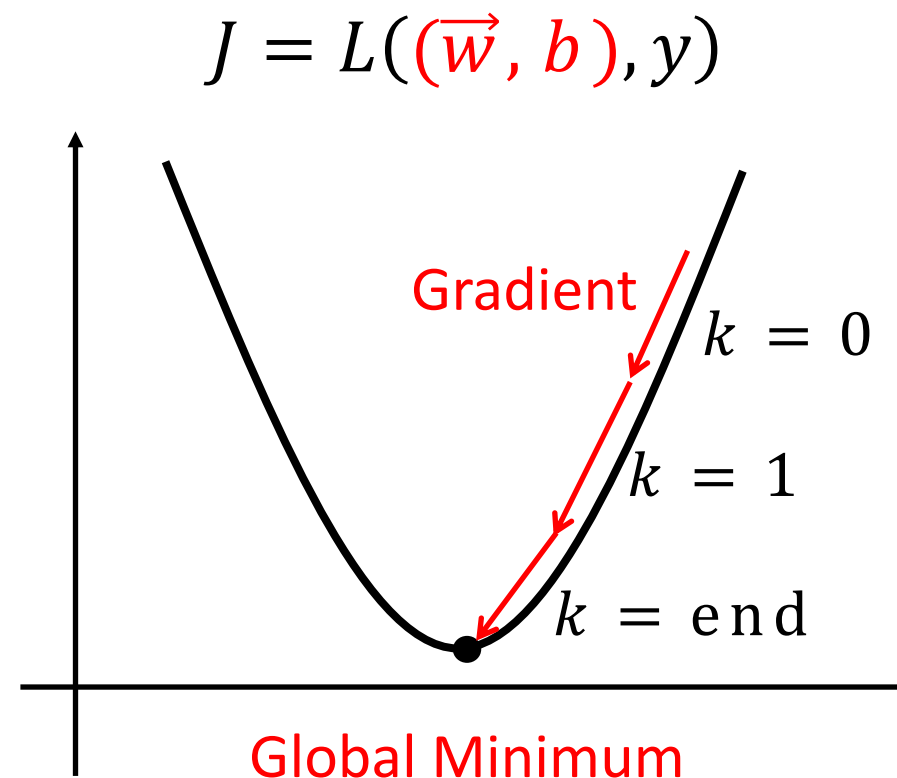
# Variable Learning Rates

$$\vec{w}_{k+1} = \vec{w}_k - \alpha \nabla_{\vec{w}} J(\vec{w}_k; b)$$

$$b_{k+1} = b_k - \alpha \frac{\partial}{\partial b} J(\vec{w}; b_k)$$

$$\alpha = f(hp_1, hp_2, \dots)$$

$$J = L((\vec{w}, b), y)$$

Gradient

$k = 0$

$k = 1$

$k = \text{end}$

Global Minimum

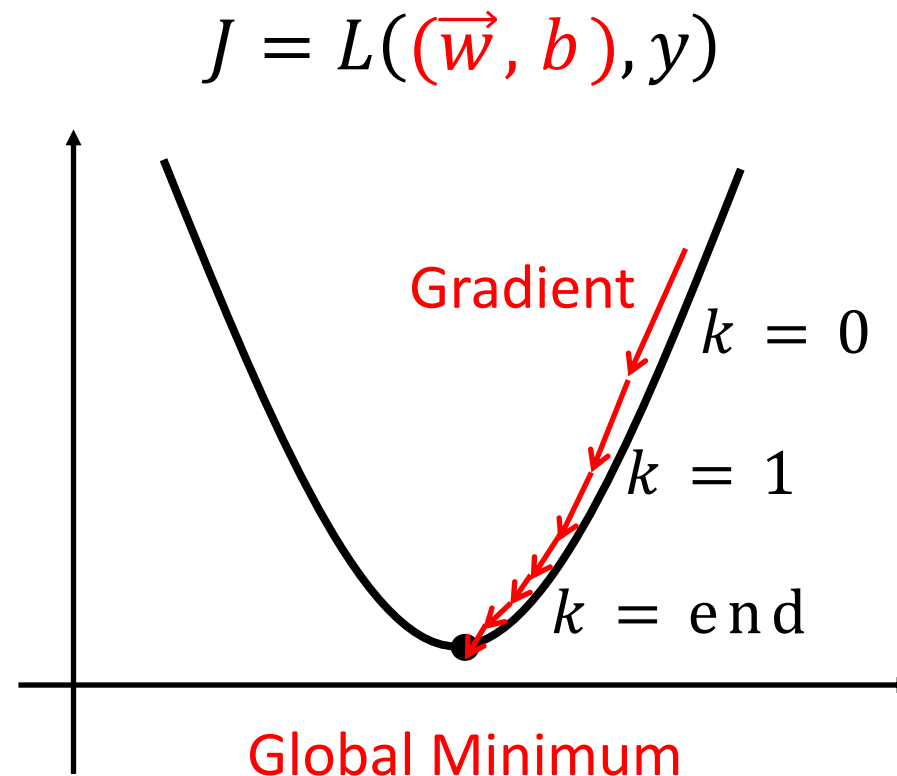# Variable Learning Rates

$$\vec{w}_{k+1} = \vec{w}_k - \alpha \nabla_{\vec{w}} J(\vec{w}_k; b)$$

$$\alpha = \frac{1}{1 + decr \cdot epnum} \alpha_0$$

$$b_{k+1} = b_k - \alpha \frac{\partial}{\partial b} J(\vec{w}; b_k)$$

$$\alpha = d^{epnum} \cdot \alpha_0$$

$$\alpha = f(hp_1, hp_2, \dots)$$

$$\alpha = \frac{d}{\sqrt{epnum}} \cdot \alpha_0$$

## Momentum

"Accelerate" gradients vectors in the right directions, to lead to faster converging.

## AdaGrad

Adagrad uses a different learning rate for every parameter $w_j$ at every step k. It eliminates the need to manually tune the learning rate.
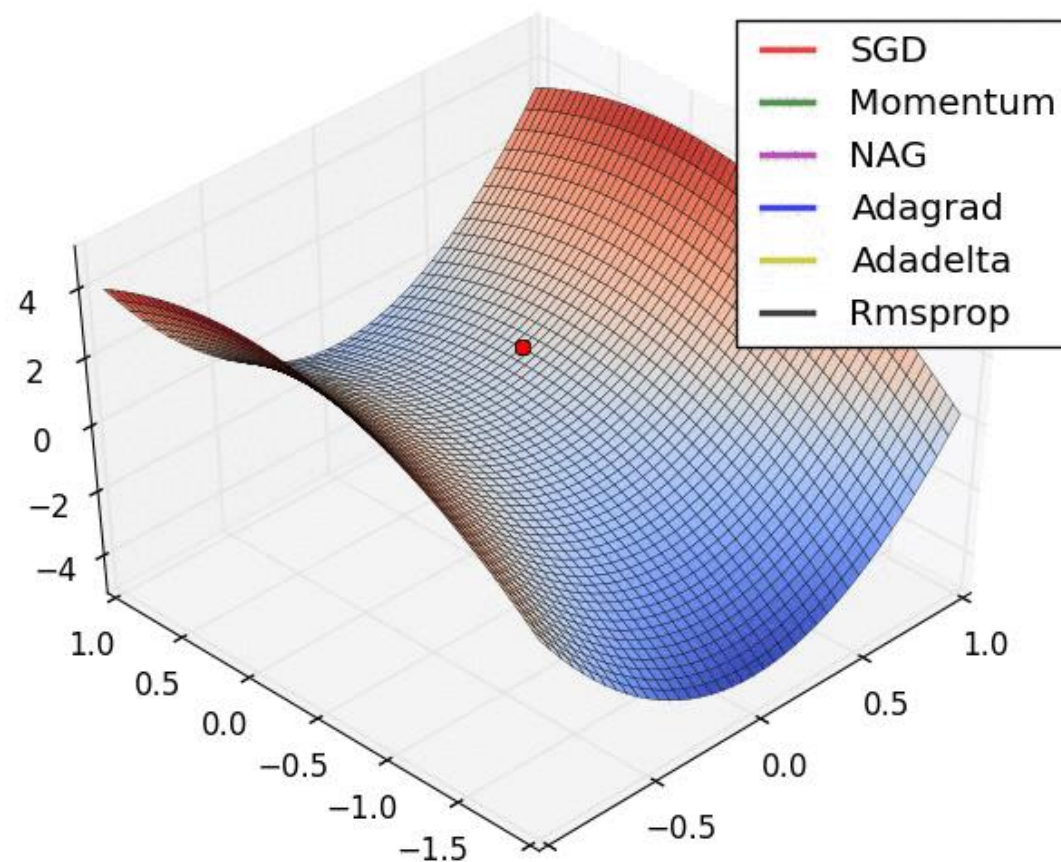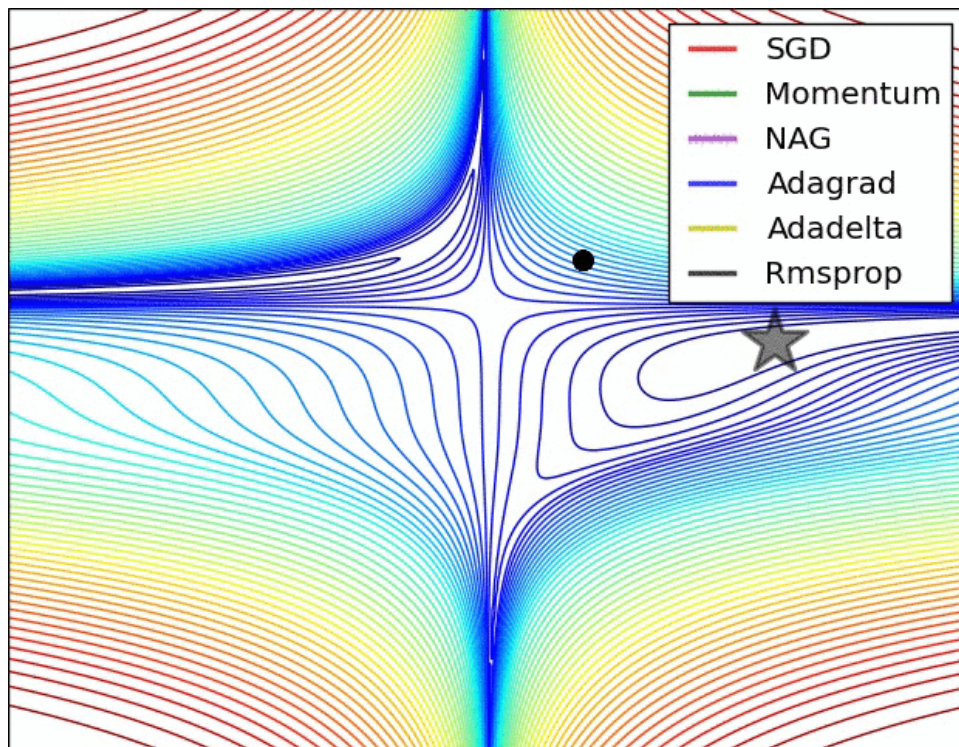
## RMSProp

"Extended" and weighted version of AdaGrad

## AdaM

Adaptive learning rate + Momentum

# Advanced Methods

# Advanced Methods

**SGD**

Fixed $\alpha$    First pass

# Advanced Methods

**SGD**
Fixed $\alpha$    First pass

RSMProp & AdaDelta
**adaptive**

Adam
**adaptive + momentum**

Worth a try if SGD
fails to converge

Standard optimizer in
DL community

# PART 3:

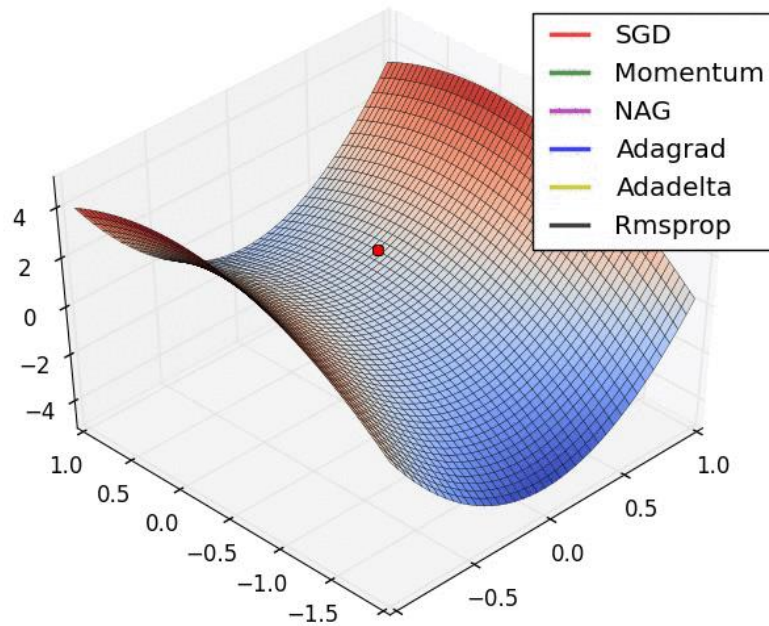## Optimization Techniques in Deep Learning

# Optimizer vs Optimization Techniques

**Optimizers**

**Optimization Techniques**

**Optimizers**

**Optimization Techniques**

# Optimizer vs Optimization Techniques

**Optimizers**

**Optimization Techniques**

- Vanilla SGD
- Momentum
- AdaGrad
- RMSProp
- Adam

# Optimizer vs Optimization Techniques

**Optimizers**

- Vanilla SGD
- Momentum
- AdaGrad
- RMSProp
- Adam

**Optimization Techniques**

Everything else that contributes to optimization

# Optimizer vs Optimization Techniques

## Optimizers

- Vanilla SGD
- Momentum
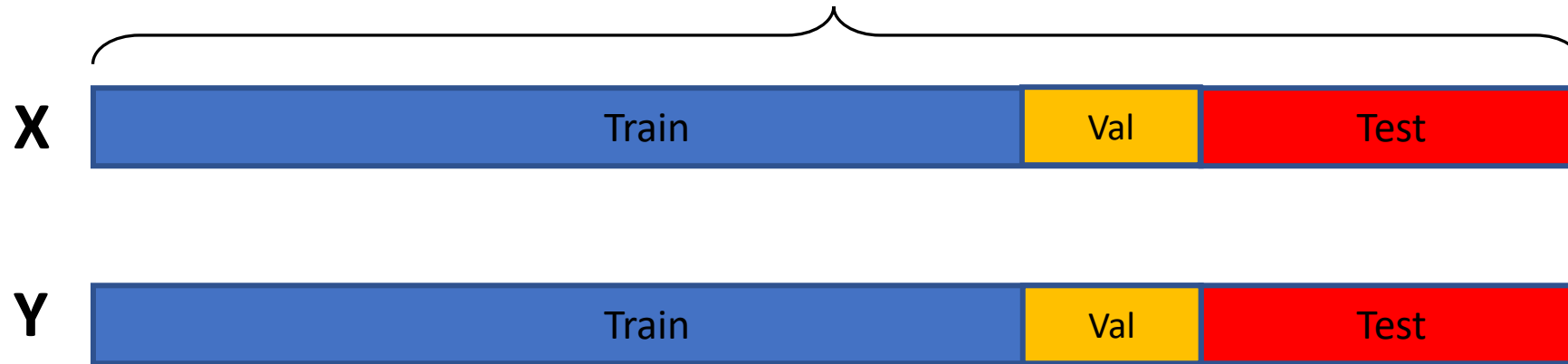- AdaGrad
- RMSProp
- Adam

## Optimization Techniques

- Data splitting (Train/Val/Test)
- Regularization
- Data normalization
- Batch-normalization
- Network initialization
- Hyperparameter tunings

# Cross Validation in Supervised Learning

Dataset (X inputs, Y targets)

**X** | Train | Val | Test |

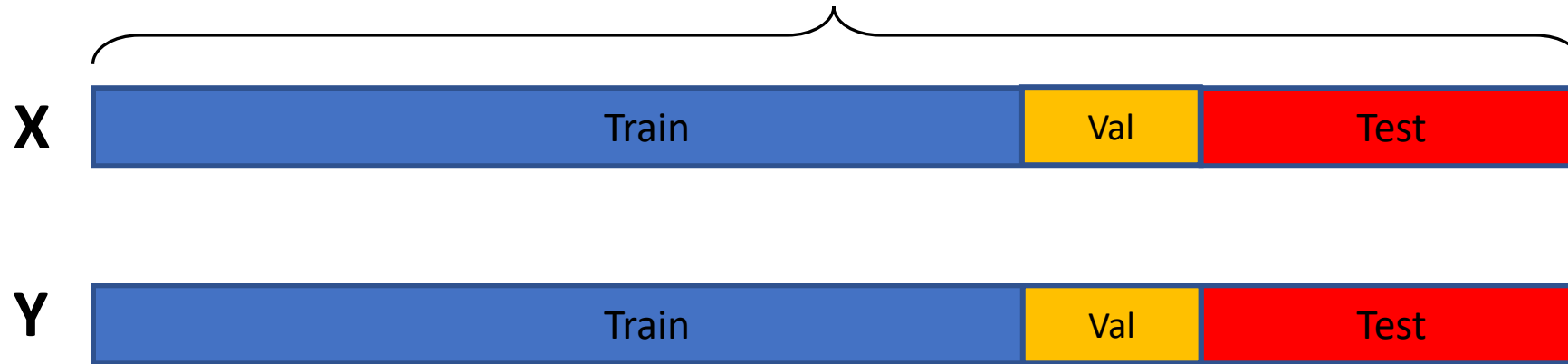**Y** | Train | Val | Test |

During training

Training loss ↓

Validation performance ↑

After training

Testing performance ↑

# Cross Validation in Supervised Learning

Dataset (X inputs, Y targets)

Common ratios

**X** | Train | Val | Test

70, 15, 15

80, 10, 10

**Y** | Train | Val | Test

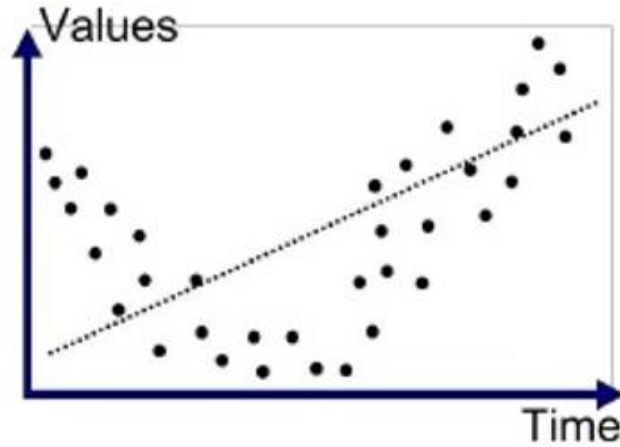60, 20, 20

**During training**

Training loss ↓

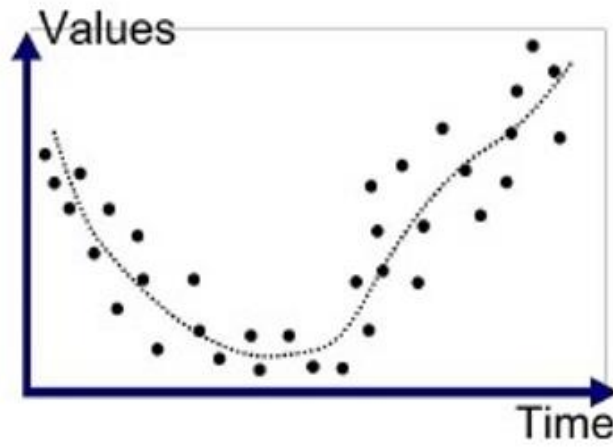Validation performance ↑

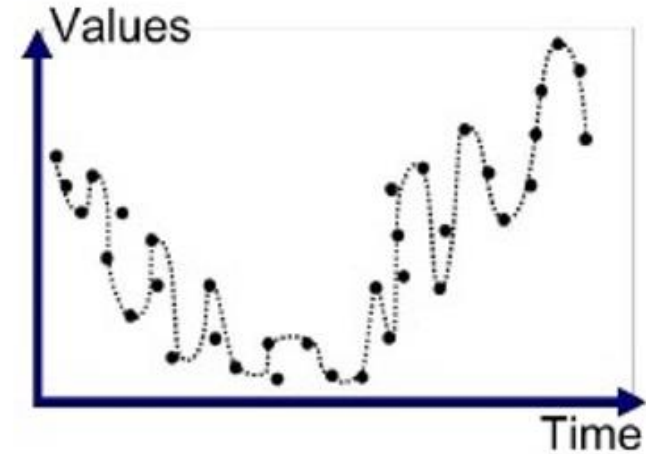**After training**

Testing performance ↑

# Overfitting vs Underfitting



Underfitted      Good Fit/Robust      Overfitted

# Overfitting vs Underfitting



Underfitted

Good Fit/Robust

Overfitted

Bad training accuracy
Bad testing accuracy

Good training accuracy
Good testing accuracy

Great training accuracy
Bad testing accuracy

# Overfitting vs Underfitting



Underfitted

Good Fit/Robust

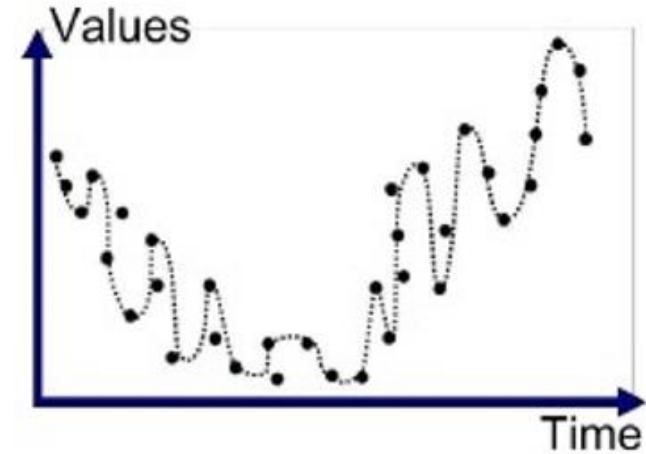Overfitted

Bad training accuracy
Bad testing accuracy

Good training accuracy
Good testing accuracy

Great training accuracy
Bad testing accuracy

**High Bias**

**High Variance**
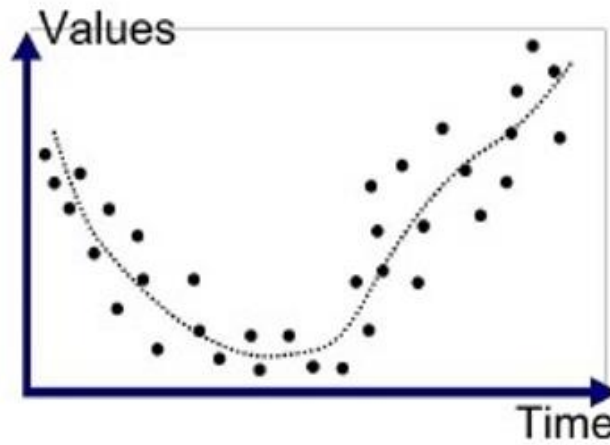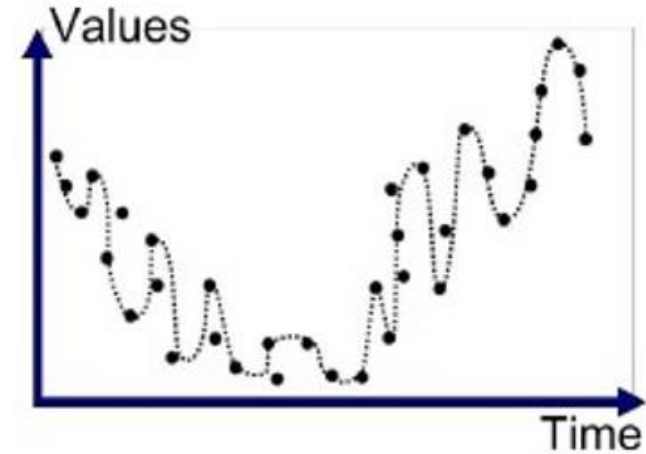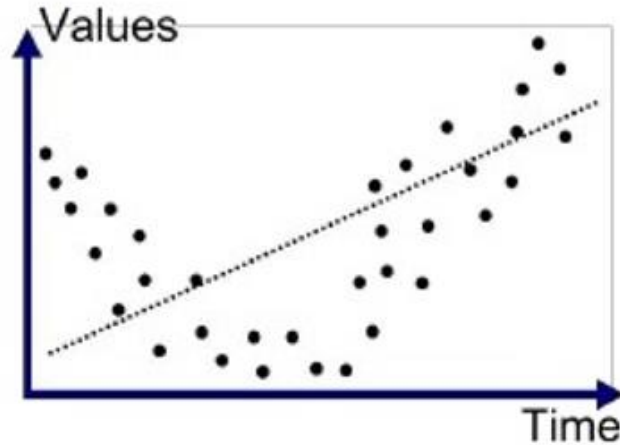
# Remedies for Overfitting/Underfitting



Underfitted          Good Fit/Robust          Overfitted

- More Layers/Neurons
- Longer Training
- Architecture
- Hyperparameter tunings

# Remedies for Overfitting/Underfitting



Underfitted      Good Fit/Robust      Overfitted

- More Layers/Neurons
- Longer Training
- Architecture
- Hyperparameter tunings

- More training data
- Regularization
- Dropout
- Initialization
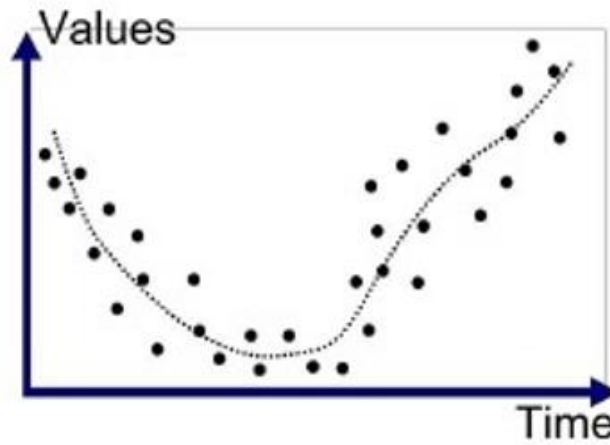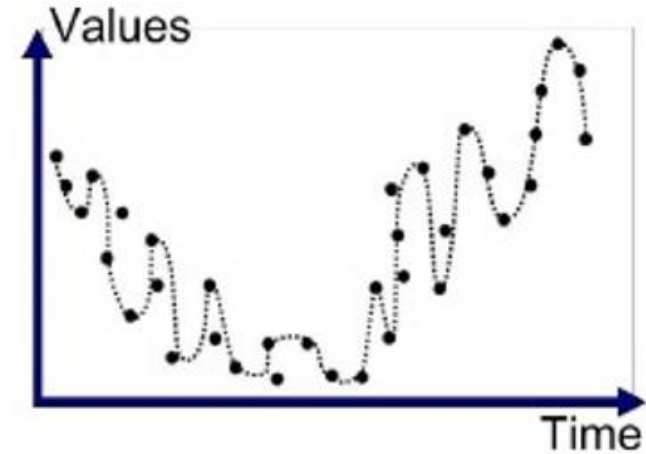
# Remedies for Overfitting/Underfitting



Underfitted

Good Fit/Robust

Overfitted

- More Layers/Neurons
- Longer Training
- Architecture
- Hyperparameter tunings

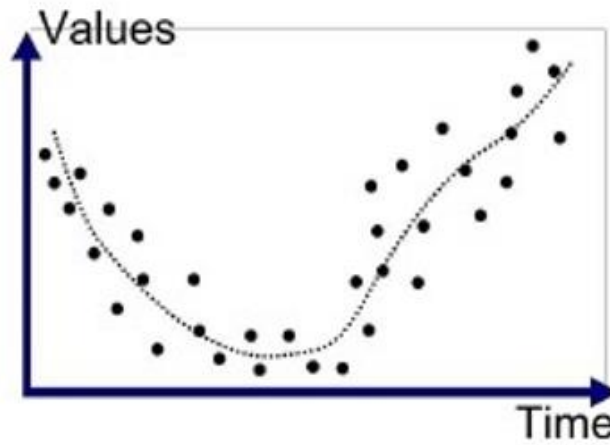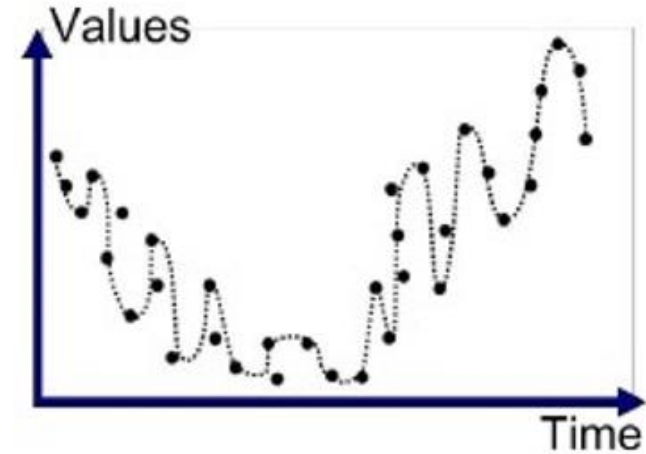- More training data
- Regularization
- Dropout
- Initialization

# L1, L2 Regularizations

## L1 Regularization

$$Loss = Error(y, \hat{y}) + \lambda \sum_{i=1}^{N} |w_i|$$

## L2 Regularization

$$Loss = Error(y, \hat{y}) + \lambda \sum_{i=1}^{N} w_i^2$$

# L1, L2 Regularizations

## L1 Regularization

$$Loss = Error(y, \hat{y}) + \lambda \sum_{i=1}^{N} |w_i|$$



w

## L2 Regularization

$$Loss = Error(y, \hat{y}) + \lambda \sum_{i=1}^{N} w_i^2$$



w

# L1, L2 Regularizations

## L1 Regularization

$$Loss = Error(y, \hat{y}) + \lambda \sum_{i=1}^{N} |w_i|$$

Penalizes sum of absolute values of weights

Results in a sparse model

Not suitable for learning complex patterns

Robust to outliers

## L2 Regularization

$$Loss = Error(y, \hat{y}) + \lambda \sum_{i=1}^{N} w_i^2$$

Penalizes sum of squared values of weights

Results in a dense model

Learns complex patterns

Sensitive to outliers

# Single-layer Regularization

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^{m} L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|\vec{w}\|_2^2$$

$$+ \frac{\lambda}{m} \|\vec{w}\|_1$$

**Cost function**

**Weight regularization terms**

$\vec{w}$

Neuron

# Multi-layer Regularization

$$J(W^{[1]}, b^{[1]}, ..., W^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^{m} L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^{L} ||W^{[l]}||_F^2$$

**Cost function**



$$||W^{[l]}||_F^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l]})^2$$

**Weight regularization term over multiple layer**

# Dropout Regularization



Standard Neural Network

Network with Dropout

# Dropout Regularization



Standard Neural Network      Network with Dropout

Dropout forces the network to learn more robust features + different random subsets of other neurons

# Dropout Regularization



Standard Neural Network

- Effectively spreading the weights
- Similar to L2 reg
- Testing with dropout $p_d=0$

Network with Dropout

- Can depend on weights (W)
- J could not be well defined in each pass

# Data Augmentation



(a)     (b)     (c)     (d)

(e)     (f)     (g)     (h)

# Data Augmentation



- Add noise

- Distortions

- Synthetic images

- Resize resolutions

- Rotation

- Add symmetries

# Early Stopping

# Exploding/Vanishing Gradients

**Very deep** neural network



$x$

$w_1$      $w_l$      $w_L$

$y$

# Exploding/Vanishing Gradients

**Very deep** neural network



$$\hat{y} = w_L \cdot \ldots w_l \cdot \ldots w_2 \cdot w_1 \cdot x$$

$$w_l = w > 1; \qquad \hat{y} = w^L x \to \infty$$

$$w_l = w < 1; \qquad \hat{y} = w^L x \to 0$$

# Exploding/Vanishing Gradients

With **activation**:

$$...w_3\sigma_3(w_2\sigma_2(\sigma_1'(w_1 x))$$

For **gradients**:

$$\frac{\partial J}{\partial w_1} = \sigma_3'(z_3)w_3\sigma_2'(z_2)w_2\sigma_1'(z_1)x$$

**Zero mean:**

$$\mu = \frac{1}{m} \sum_{i=1}^{m} x^{(i)}$$

$$x^{(i)\mu} = x^{(i)} - \mu$$

**Normalized Variances**

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^{m} x^{(i)^2}$$

$$x^{(i)\mu,\sigma^2} = x^{(i)\mu}./\sigma^2$$

# Intuition for data normalization

If **inputs have different scales**, the **cost function** will also have to include different scales → increased likelihood of instability



Remember to **normalize all sets**: training, validation, testing

First training sample

Mini-batch of 3 training samples

$$\begin{array}{ccc} 2 & 5 & 7 \\ 4 & 3 & 2 \\ 7 & 8 & 9 \end{array}$$

Input

Layer 1

First training sample

Mini-batch of 3 training samples

$$\begin{bmatrix} 2 & 5 & 7 \\ 4 & 3 & 2 \\ 7 & 8 & 9 \end{bmatrix}$$

Input

Layer 1

Layer 1 output for first sample

$$\begin{bmatrix} 3 & 5 & 7 & 5 & 7 \\ 1 & 5 & 3 & 7 & 8 \\ 7 & 8 & 4 & 2 & 4 \end{bmatrix}$$

$mean = 3.7$
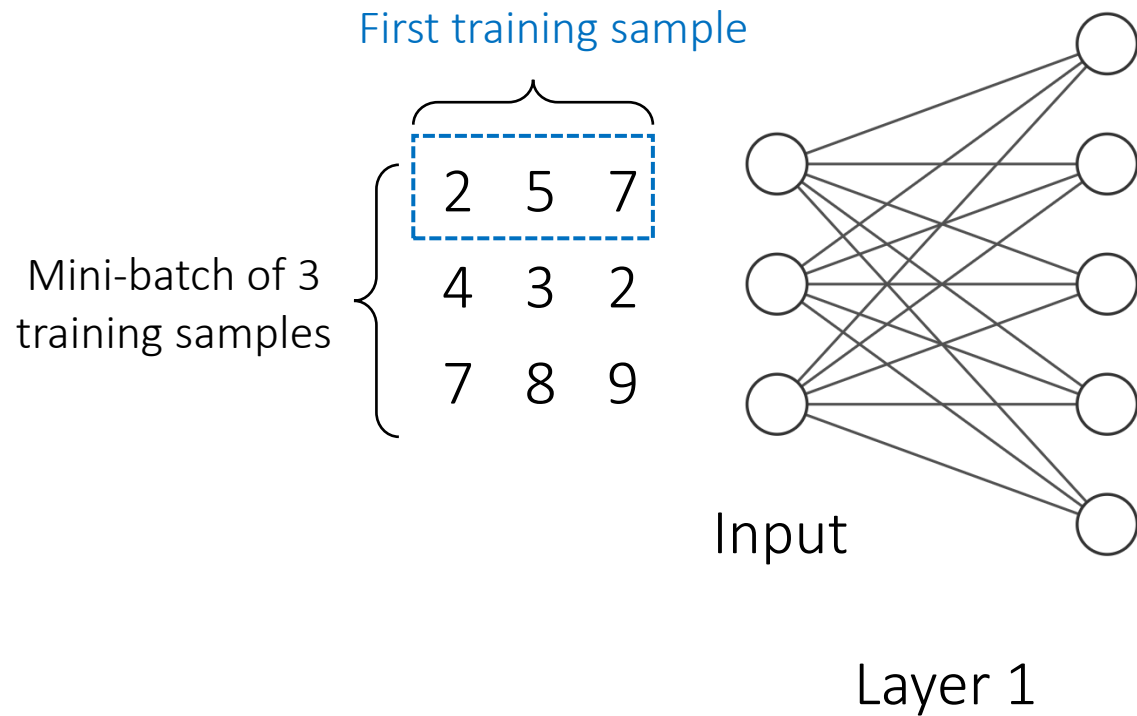$std = 2.5$

Normalized layer 1 output

$$\begin{bmatrix} -0.3 & -0.7 & 1.4 & 0.2 & 0.4 \\ -1.1 & -0.7 & -1.0 & 1.1 & 1.0 \\ 1.3 & 1.4 & -0.4 & -1.3 & -1.4 \end{bmatrix}$$

$mean \approx 0$
$std \approx 1$

Batch normalization

64

# Remedies for Vanishing/Exploding Gradients: Weight Initialization



Proper weight initialization plays essential roles in preventing exploding/vanishing gradients

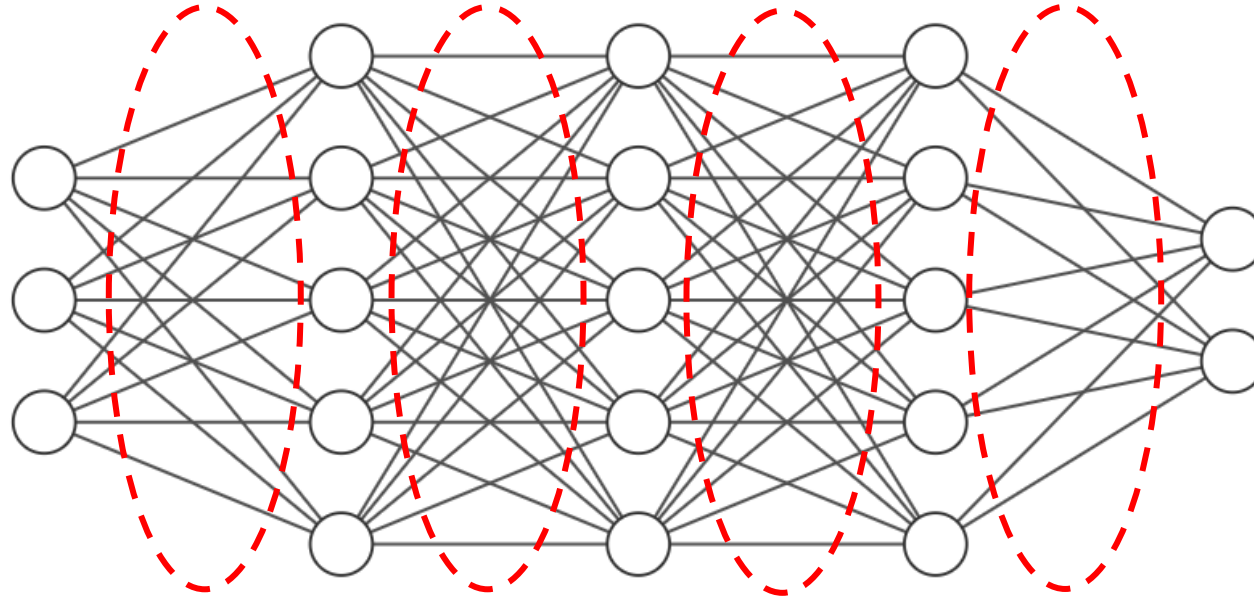# Remedies for Vanishing/Exploding Gradients: Weight Initialization



Proper weight initialization plays essential roles in preventing exploding/vanishing gradients

↓

Faster convergence

# Network Initialization

- Zero → Problematic
- Random Normal (0,1) -> Problematic
- Xavier (tanh):

$$Var(w^{[l]}) : 1/n^{[l-1]}$$

$$w^{[l]} = N(0,1) \cdot \sqrt{\frac{1}{n^{[l-1]}}}$$

# Network Initialization

- He (ReLU):

$$Var(w^{[l]}) : 2/n^{[l-1]}$$

- Other:

$$w^{[l]} = N(0, 1) \cdot \sqrt{\frac{2}{n^{[l-1]}}}$$

$$Var(w^{[l]}) : \frac{2}{n^{[l-1]} + n^{[l]}}$$

# Hyperparameters

- Learning rate

- Number of layers

- Neurons in each layer

- Activation function
(ReLU, Tanh, sigmoid)

- Training batch size
(SGD, Mini-batch, Batch Gradient)

- Optimizer
(SGD, Adam, RMS Prop etc)

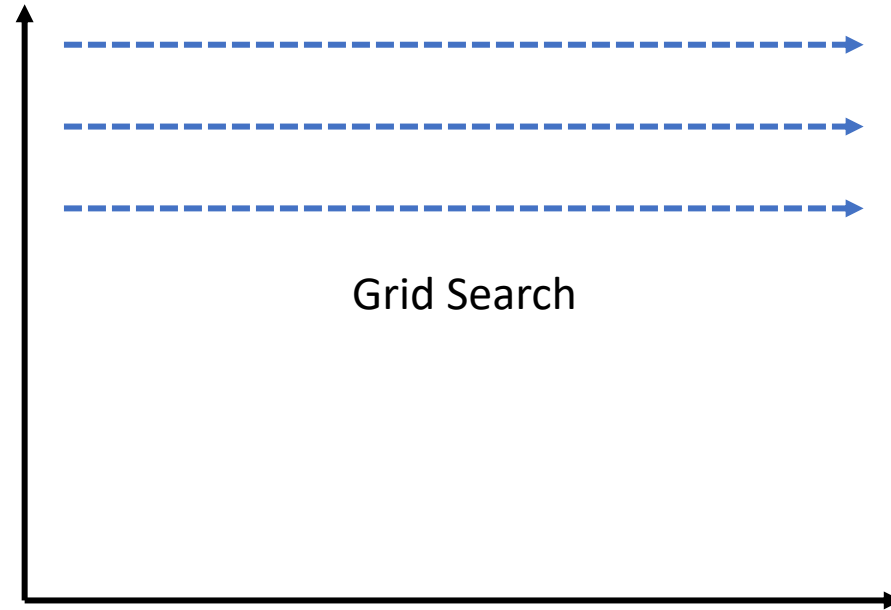- Number of training epochs

# Hyperparameters

- Learning rate

- Number of layers

- Neurons in each layer

- Activation function
  (ReLU, Tanh, sigmoid)

- Training batch size
  (SGD, Mini-batch, Batch Gradient)

- Optimizer
  (SGD, Adam, RMS Prop etc)

- Number of training epochs

Number of layers

Grid Search

Learning rate

# Hyperparameters

- Learning rate

- Number of layers

- Neurons in each layer

- Activation function
  (ReLU, Tanh, sigmoid)

- Training batch size
  (SGD, Mini-batch, Batch Gradient)

- Optimizer
  (SGD, Adam, RMS Prop etc)

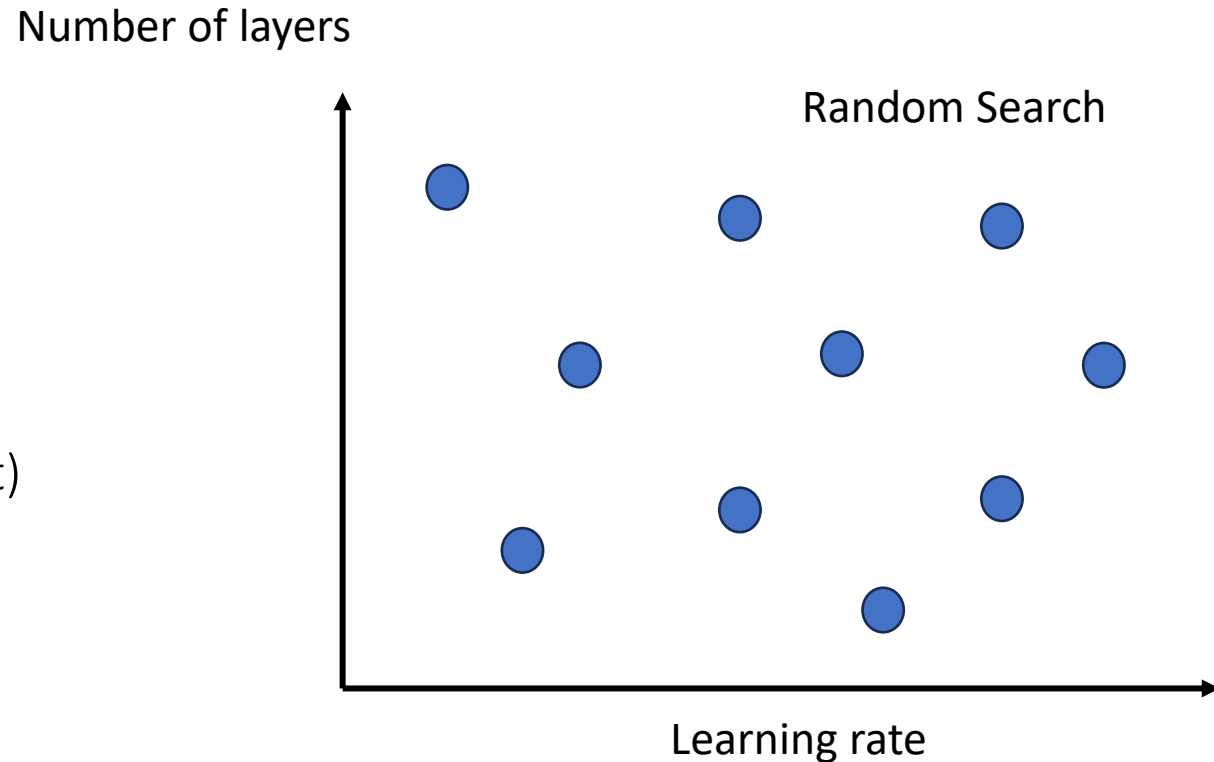- Number of training epochs

Number of layers

Random Search



Learning rate

# Summary

## Optimizers

- Vanilla SGD
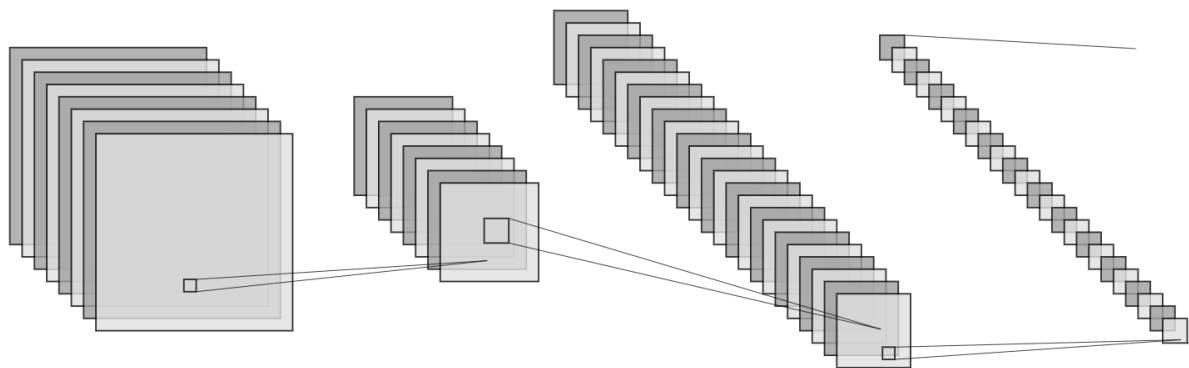- Momentum
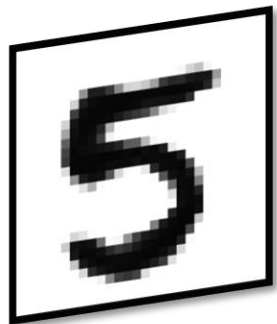- AdaGrad
- RMSProp
- Adam

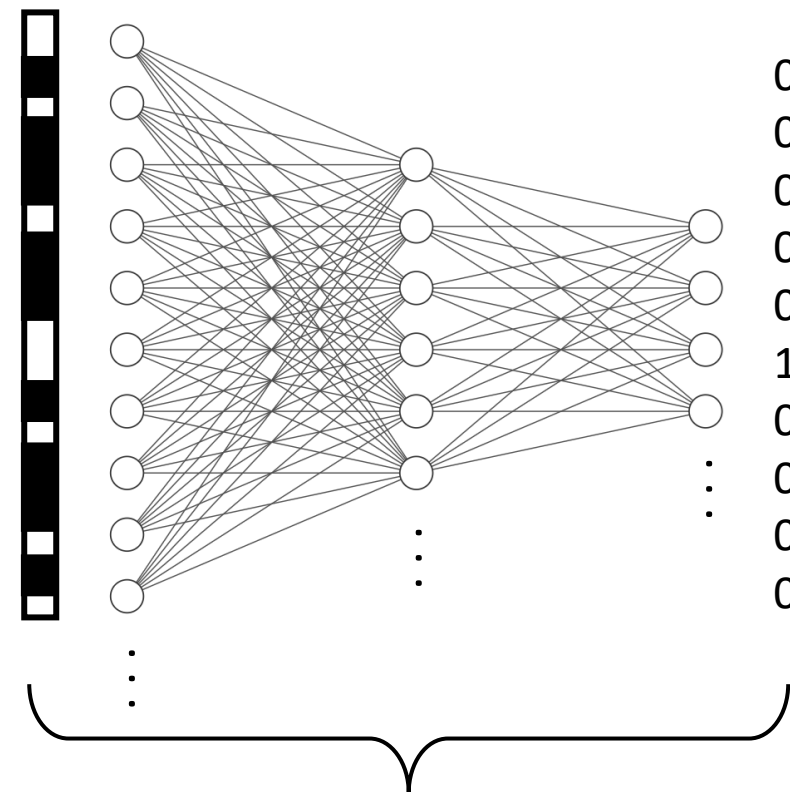Balance

## Optimization Techniques

- Data splitting (Train/Val/Test)
- Regularization
- Data normalization
- Batch-normalization
- Network initialization
- Hyperparameter tunings

# Next episode in EEP 596

Convolution Layers + Pooling Layers
(Image feature extraction)

Fully connected layers
(Classifier)