

# Optimization Algorithms in Deep Learning

---

by Dr. Rishikesh Yadav

September 29, 2025

Assistant Professor, School of Mathematical and Statistical Sciences, IIT Mandi, India

# Table of Contents

1. Recap: Regression and Classification
2. Role of Gradients and Hessians in Optimization
3. Famous Optimization Algorithms Used in Deep Learning
  - 3.1 Newton-Raphson Algorithm
  - 3.2 Gradient Descent and its Variants
    - Batch Gradient Descent: The Standard One
    - Stochastic Gradient Descent (SGD)
    - Mini-batch Gradient Descent
  - 3.3 Momentum Algorithm: Accelerating Gradient Descent
  - 3.4 Adaptive Learning Rates-Based Algorithms
    - Adagrad: Adaptive Gradient Algorithm
    - RMSProp: Root Mean Square Propagation
    - Adam: Adaptive Moment Estimation
  - 3.5 Theoretical Convergence of Optimizers
4. Practicals: Training Regression and Classification

## Recap: Regression and Classification

---

# Supervised Learning: General Framework

- **Goal:** Learn (estimate) a function  $f(\mathbf{x}; \mathbf{w})$  that maps several input features  $\mathbf{x}$  to an output variable  $y$ .

$$y \approx f(\mathbf{x}; \mathbf{w}) \quad (\text{there will be some error terms as well})$$

- **Data:** A collection of input–output pairs

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$$

- Inputs (features):  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})^\top$
- Outputs (targets):  $y_i$  (numeric or categorical)
- **Learning:** Estimate the unknown parameters (weights in machine learning)  $\mathbf{w}$  by minimizing a suitable loss function  $\mathcal{L}(\mathbf{w})$ .
- **Prediction:** Once  $\hat{\mathbf{w}}$  is obtained, use

$$\hat{y} = f(\mathbf{x}; \hat{\mathbf{w}})$$

for new input  $\mathbf{x}$ .

# Regression Models

- **Task:** Predict a **continuous output**  $y \in \mathbb{R}$ .
- **One simpler model (Assume  $f$  is linear in  $\mathbf{w}$ ):** Multiple Linear Regression (MLR)

$$y_i = w_0 + w_1 x_{i1} + w_2 x_{i2} + \cdots + w_p x_{ip} + \varepsilon_i$$

- **Loss function (Mean Squared Error):**

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

- **Ordinary Least Squares (OLS) solution:**

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

where

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

- **Important Note:** For all regression types of problems, we might not get the closed-form expressions of  $\hat{\mathbf{w}}$ .

# Classification Models

- **Task:** Predict a **categorical output**  $y \in \{1, 2, \dots, K\}$ .
- **Model:** Logistic regression (binary or multinomial)

$$P(Y = k \mid \mathbf{x}, \mathbf{w}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x})}{\sum_{j=1}^K \exp(\mathbf{w}_j^\top \mathbf{x})}$$

- **Loss function (Cross-Entropy / Log-Loss):**

$$\mathcal{L}(\mathbf{w}) = - \sum_{i=1}^n \sum_{k=1}^K \mathbf{1}\{y_i = k\} \log P(Y = k \mid \mathbf{x}_i, \mathbf{w})$$

- **Estimation:** No closed-form solution; parameters  $\mathbf{w}$  are found using iterative optimization methods (e.g., gradient descent, Adam, Newton–Raphson etc.)
- **Prediction:** Assign the class with the highest predicted probability

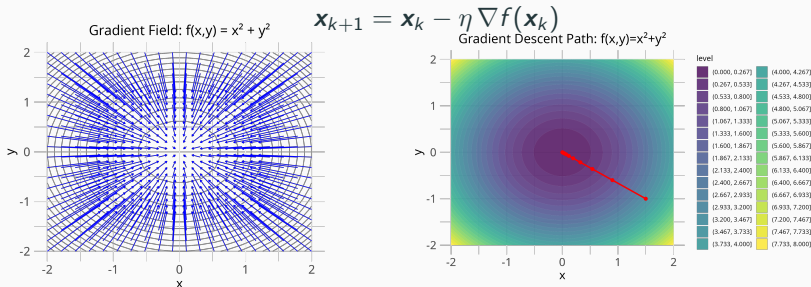
$$\hat{y}_i = \arg \max_{k \in \{1, \dots, K\}} P(Y = k \mid \mathbf{x}_i, \mathbf{w})$$

## Role of Gradients and Hessians in Optimization

---

# Role of Gradients in Optimization

- The **gradient**  $\nabla f(x)$  points in the direction of the steepest increase of the function.
- Optimization algorithms use the **negative gradient** to move towards a minimum.
- **Intuition:**
  - If  $\nabla f(x) = 0$ , we are at a **stationary point**.
  - Gradient magnitude  $\|\nabla f(x)\|$  indicates how steep the surface is.
- **Example:** Gradient Descent updates

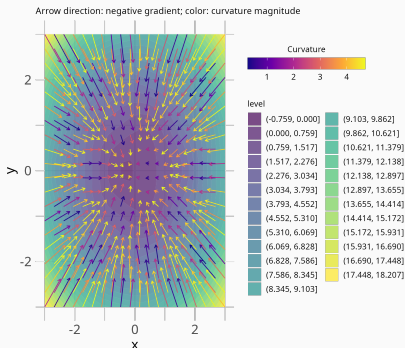




# Role of Hessians in Optimization

- The **Hessian matrix**  $H(\mathbf{x}) = \nabla^2 f(\mathbf{x})$  captures the **curvature** of the function.
- **Key roles:**
  - Determines if a stationary point is a minimum, maximum, or saddle.
  - Guides second-order methods (e.g., Newton's method).

- Positive definite Hessian  $\Rightarrow$  local minimum.
- Negative definite Hessian  $\Rightarrow$  local maximum.
- Indefinite Hessian  $\Rightarrow$  saddle point.



- **Example:** Newton-Raphson update

$$\mathbf{x}_{k+1} = \mathbf{x}_k - H(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)$$

# **Famous Optimization Algorithms Used in Deep Learning**

---

# Agenda

## 3. Famous Optimization Algorithms Used in Deep Learning

### 3.1 Newton-Raphson Algorithm

### 3.2 Gradient Descent and its Variants

- Batch Gradient Descent: The Standard One

- Stochastic Gradient Descent (SGD)

- Mini-batch Gradient Descent

### 3.3 Momentum Algorithm: Accelerating Gradient Descent

### 3.4 Adaptive Learning Rates-Based Algorithms

- Adagrad: Adaptive Gradient Algorithm

- RMSProp: Root Mean Square Propagation

- Adam: Adaptive Moment Estimation

### 3.5 Theoretical Convergence of Optimizers

# Newton-Raphson Method

- **What is Newton–Raphson?** An iterative optimization method using both the **gradient** and the **Hessian**:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \left[ \nabla_{\mathbf{w}}^2 L(\mathbf{w}_t) \right]^{-1} \nabla_{\mathbf{w}} L(\mathbf{w}_t)$$

- Uses **curvature information** (Hessian) for optimal step size and direction
- **Application Example:** classical **Logistic regression** with small, well-conditioned datasets

## Advantages over Gradient Descent:

- ✓ **Quadratic convergence** near optimum (vs. linear)
- ✓ **No learning rate**  $\eta$  to tune - step size is adaptive
- ✓ **Accounts for curvature** - follows natural shape of loss landscape

# Geometric Intuition of Hessian in Newton's Method

- Hessian ( $H$ ) captures the **curvature of the loss surface**:

$$H = \nabla^2 \mathcal{L}(\mathbf{w})$$

- Newton update:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - H^{-1} \nabla \mathcal{L}(\mathbf{w}_t)$$

- Geometric intuition:

- Directions with **steep curvature**  $\rightarrow$  move **smaller steps**
- Directions with **flat curvature**  $\rightarrow$  move **larger steps**
- Hessian inverse rescales gradient according to curvature.

When Hessian is singular or non-invertible:

- Occurs at **saddle points**, flat regions, or redundant parameters.

# Geometric Intuition of Hessian in Newton's Method

- Hessian ( $H$ ) captures the **curvature of the loss surface**:

$$H = \nabla^2 \mathcal{L}(\mathbf{w})$$

- Newton update:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - H^{-1} \nabla \mathcal{L}(\mathbf{w}_t)$$

- Geometric intuition:

- Directions with **steep curvature**  $\rightarrow$  move **smaller steps**
- Directions with **flat curvature**  $\rightarrow$  move **larger steps**
- Hessian inverse rescales gradient according to curvature.

When Hessian is singular or non-invertible:

- Occurs at **saddle points**, flat regions, or redundant parameters.
- Using a **generalized inverse (pseudoinverse)**  $H^+$  is possible:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - H^+ \nabla \mathcal{L}(\mathbf{w}_t)$$

# Geometric Intuition of Hessian in Newton's Method

- Hessian ( $H$ ) captures the **curvature of the loss surface**:

$$H = \nabla^2 \mathcal{L}(\mathbf{w})$$

- Newton update:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - H^{-1} \nabla \mathcal{L}(\mathbf{w}_t)$$

- Geometric intuition:

- Directions with **steep curvature**  $\rightarrow$  move **smaller steps**
- Directions with **flat curvature**  $\rightarrow$  move **larger steps**
- Hessian inverse rescales gradient according to curvature.

When Hessian is singular or non-invertible:

- Occurs at **saddle points**, flat regions, or redundant parameters.
- Using a **generalized inverse (pseudoinverse)**  $H^+$  is possible:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - H^+ \nabla \mathcal{L}(\mathbf{w}_t)$$

- **Intuition:** move along **curved directions only**, ignore flat directions.

# Geometric Intuition of Hessian in Newton's Method

- Hessian ( $H$ ) captures the **curvature of the loss surface**:

$$H = \nabla^2 \mathcal{L}(\mathbf{w})$$

- Newton update:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - H^{-1} \nabla \mathcal{L}(\mathbf{w}_t)$$

- Geometric intuition:

- Directions with **steep curvature**  $\rightarrow$  move **smaller steps**
- Directions with **flat curvature**  $\rightarrow$  move **larger steps**
- Hessian inverse rescales gradient according to curvature.

## When Hessian is singular or non-invertible:

- Occurs at **saddle points**, flat regions, or redundant parameters.
- Using a **generalized inverse (pseudoinverse)**  $H^+$  is possible:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - H^+ \nabla \mathcal{L}(\mathbf{w}_t)$$

- **Intuition:** move along **curved directions only**, ignore flat directions.
- **Effect:** slower or less precise convergence, but still usable.



# Geometric Intuition of Hessian in Newton's Method

- Hessian ( $H$ ) captures the **curvature of the loss surface**:

$$H = \nabla^2 \mathcal{L}(\mathbf{w})$$

- Newton update:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - H^{-1} \nabla \mathcal{L}(\mathbf{w}_t)$$

- Geometric intuition:

- Directions with **steep curvature**  $\rightarrow$  move **smaller steps**
- Directions with **flat curvature**  $\rightarrow$  move **larger steps**
- Hessian inverse rescales gradient according to curvature.

When Hessian is singular or non-invertible:

- Occurs at **saddle points**, flat regions, or redundant parameters.
- Using a **generalized inverse (pseudoinverse)**  $H^+$  is possible:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - H^+ \nabla \mathcal{L}(\mathbf{w}_t)$$

- **Intuition:** move along **curved directions only**, ignore flat directions.
- **Effect:** slower or less precise convergence, but still usable.
- **Alternative practical fix: Damped Newton:**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - (H + \lambda I)^{-1} \nabla \mathcal{L}(\mathbf{w}_t)$$

ensures positive definiteness and stable updates.

# From Newton-Raphson to Gradient Descent

- **Hessian Computation:**  $\nabla_w^2 \mathcal{L}(\mathbf{w})$  requires  $O(p^2)$  memory
  - Modern DL: millions of parameters  $\rightarrow$  infeasible storage
- **Matrix Inversion:**  $[\nabla_w^2 \mathcal{L}(\mathbf{w})]^{-1}$  costs  $O(p^3)$  computationally
- **Non-Convexity:** Hessian may not be positive definite  $\rightarrow$  convergence issues
- **Sensitivity:** Requires careful initialization and well-conditioned problems

## Gradient Descent: The Practical Choice

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla_w \mathcal{L}(\mathbf{w}_t)$$

- **Memory:** Stores only gradient ( $O(p)$ ) vs. Hessian ( $O(p^2)$ )
- **Speed:**  $O(p)$  per update vs. Newton's  $O(p^3)$
- **Robustness:** Works reliably even in non-convex landscapes
- **Scalability:** Mini-batch version enables training on massive datasets

## 3. Famous Optimization Algorithms Used in Deep Learning

### 3.1 Newton-Raphson Algorithm

### 3.2 Gradient Descent and its Variants

Batch Gradient Descent: The Standard One

Stochastic Gradient Descent (SGD)

Mini-batch Gradient Descent

### 3.3 Momentum Algorithm: Accelerating Gradient Descent

### 3.4 Adaptive Learning Rates-Based Algorithms

Adagrad: Adaptive Gradient Algorithm

RMSProp: Root Mean Square Propagation

Adam: Adaptive Moment Estimation

### 3.5 Theoretical Convergence of Optimizers

# Gradient Descent: The Core Idea

- **Goal:** Minimize a loss function  $L(w)$  to find the best model parameters  $w$  for regression/classification.
- **Intuition:** Find the lowest point in a valley by always walking downhill.

## The Update Rule:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla_w \mathcal{L}(\mathbf{w}_t)$$

- $\mathbf{w}$ : Model parameters (e.g., weights)
- $\mathcal{L}(\mathbf{w})$ : Loss function (e.g., Mean Squared Error, Log Loss)
- $\eta$ : Learning rate (step size)
- $\nabla_w \mathcal{L}(\mathbf{w})$ : Gradient (direction of steepest ascent)

**Motivation:** Fundamental algorithm for training models like **Linear Regression** and **Logistic Regression**.

# Batch Gradient Descent: The Standard Version

## Algorithm:

1. Initialize parameters  $\mathbf{w}$  randomly.
2. **Compute Gradient** over the **entire dataset**:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{w}} \text{Loss}(f_{\mathbf{w}}(\mathbf{x}^{(i)}), y^{(i)})$$

3. **Update:**  $\mathbf{w} = \mathbf{w} - \eta \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})$
4. Repeat until convergence.

## Characteristics:

- ✓ **Pros:** Stable convergence. Guaranteed for convex functions.
- ✗ **Cons:** **Very slow for large datasets.** One update requires a full data pass.

# Stochastic Gradient Descent (SGD)

**Core Idea:** Use a **single, random** training example  $(\mathbf{x}^{(i)}, y^{(i)})$  to compute a **noisy** gradient.

**Algorithm per Epoch (one complete pass through data):**

1. Shuffle the entire dataset.
2. For each example (input)  $i$  in the dataset:
  - 2.1 Compute gradient for one example:  $\nabla_w \mathcal{L}(\mathbf{w}; \mathbf{x}^{(i)}, y^{(i)})$
  - 2.2 Update immediately:  $\mathbf{w} = \mathbf{w} - \eta \nabla_w \mathcal{L}(\mathbf{w}; \mathbf{x}^{(i)}, y^{(i)})$

**Key Properties:**

- ✓ **Pros:** Extremely fast per update. Can escape local minima due to noise.
- ✗ **Cons:** **Very noisy path.** Loss may fluctuate heavily. Harder to converge precisely.

# Mini-batch Gradient Descent

**Core Idea:** The **best compromise**. Use a small random subset (a **mini-batch**) of size  $b$  to compute the gradient.

## Algorithm per Epoch:

1. Shuffle the dataset.
2. For each batch of  $b$  examples:
  - 2.1 Compute gradient for the batch:
$$\nabla_w \mathcal{L}(\mathbf{w}) = \frac{1}{b} \sum_{k=1}^b \nabla_w \text{Loss}(f_w(\mathbf{x}^{(k)}), y^{(k)})$$
  - 2.2 Update parameters:  $\mathbf{w} = \mathbf{w} - \eta \nabla_w \mathcal{L}(\mathbf{w})$

## Key Properties:

- ✓ **Pros:** **Efficient** and leverages GPU parallelism. **More stable** than SGD.
- ✗ **Cons:** Introduces the batch size  $b$  as a new hyperparameter to tune.

## Comparison: GD, SGD, and Mini-batch GD

Criterion	Batch GD	Stochastic GD	Mini-batch GD
Gradient	Full dataset	Single example	Small batch ( $b$ )
Speed/Update	Slow	Very Fast	Fast
Stability	Smooth	Noisy	Moderate
Memory	High	Low	Medium
Parallelization	Difficult	No	Excellent
Use Case	Small datasets	Large datasets	Deep Learning

**Conclusion:** For most modern machine learning tasks, especially deep learning, [Mini-batch Gradient Descent](#) is the preferred algorithm.



# Epochs and Iterations

- **Iteration:** One parameter update using a single batch of data.

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t)$$

- **Epoch:** One complete pass over the entire dataset.
- Relationship between them:

$$\text{Iterations per epoch} = \frac{N}{\text{Batch Size}}$$

- Each epoch revisits all training samples once; multiple epochs refine model parameters progressively.

## Interpretation:

- **Iterations** = micro updates within an epoch
- **Epoch** = one full cycle through the dataset.

# Data Division in Gradient Descent Variants

Method	Batch Size	Description	Key Property
Batch GD	$N$	Uses all data per update	Stable but slow
Stochastic GD (SGD)	1	One sample per update	High variance, fast convergence
Mini-Batch GD	$b$ (e.g. 32–512)	Uses small random batches per update	Balanced trade-off between noise and stability

## Important:

- At the start of each *epoch*, the dataset is **randomly shuffled** and then divided into batches. Each batch is used once per epoch — i.e., **sampling without replacement**. (This is widely used now)
- In **Online SGD**, samples are drawn **randomly with replacement**:
  - Each iteration picks a random data point.
  - Some samples may be seen multiple times before others.
  - There is no strict notion of “epoch”.
  - Commonly used in **streaming** or **very large-scale** settings. Barely used

# Choosing the Learning Rate $\eta$ : A Practical Guide

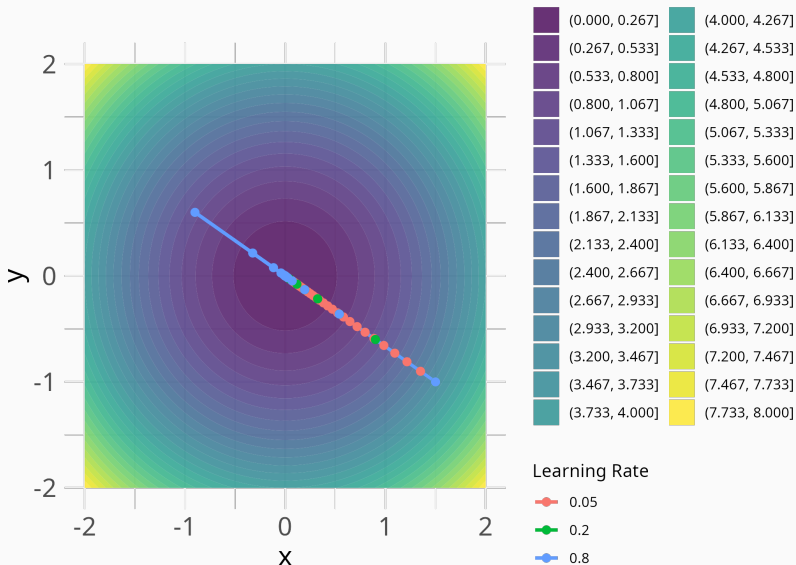
- **The Trade-off:** Newton-Raphson had no  $\eta$  (Hessian adaptive), but GD requires careful tuning

## Effects of Learning Rate:

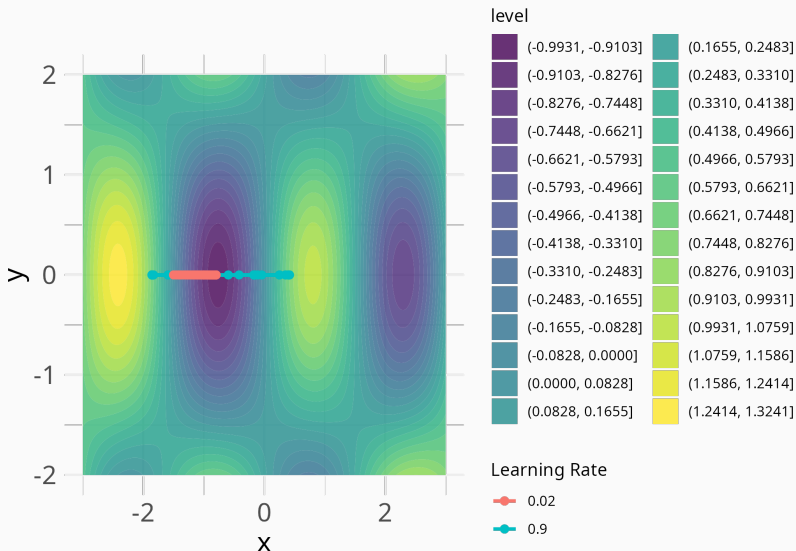
- ✗ **Too Large ( $\eta$  big):** Overshooting, divergence, unstable loss
- ✗ **Too Small ( $\eta$  small):** Extremely slow convergence, stuck in local minima
- ✓ **Just Right:** Stable, efficient convergence to good minimum

**Rule of Thumb:** Start with  $\eta = 0.01$  or  $0.001$  and adjust based on loss curve behavior.

# Gradient Descent with Different Learning Rates: Convex



# Gradient Descent with Different Learning Rates: Non-Convex



Small learning rate **stays**; large learning rate **jumps to next basin**.

## Practical Strategies:

- **Grid Search:** Try values like 0.1, 0.01, 0.001, 0.0001
- **Learning Rate Schedule:** Start large, decrease over time (e.g.,  $\eta_t = \frac{\eta_0}{1+t}$ )
- **Adaptive Methods:** Use algorithms like the momentum algorithm, Adam, Adagrad that auto-tune  $\eta$

## 3. Famous Optimization Algorithms Used in Deep Learning

### 3.1 Newton-Raphson Algorithm

### 3.2 Gradient Descent and its Variants

Batch Gradient Descent: The Standard One

Stochastic Gradient Descent (SGD)

Mini-batch Gradient Descent

### 3.3 Momentum Algorithm: Accelerating Gradient Descent

### 3.4 Adaptive Learning Rates-Based Algorithms

Adagrad: Adaptive Gradient Algorithm

RMSProp: Root Mean Square Propagation

Adam: Adaptive Moment Estimation

### 3.5 Theoretical Convergence of Optimizers

# Momentum: Inertia in Optimization

## Motivation:

- Plain (Vanilla) SGD oscillates heavily in ravines (narrow valleys).
- Wastes time moving back-and-forth in directions of steep curvature.
- **Idea:** build inertia — keep moving in the average descent direction.

## Update Rule:

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

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{m}_t$$

- $\mathbf{g}_t$ : current gradient
- $\mathbf{m}_t$ : velocity (exponentially averaged gradients).
- $\beta$ : momentum coefficient (0.9 typical).

## Characteristics:

- ✓ Faster convergence in ravines.
- ✓ Reduces oscillations.
- ✗ May overshoot if  $\beta$  too high.

Look at the YouTube lecture of Andrew NG: [watch?v=k8fTYJPd3\\_I](https://www.youtube.com/watch?v=k8fTYJPd3_I)



# Agenda

## 3. Famous Optimization Algorithms Used in Deep Learning

### 3.1 Newton-Raphson Algorithm

### 3.2 Gradient Descent and its Variants

Batch Gradient Descent: The Standard One

Stochastic Gradient Descent (SGD)

Mini-batch Gradient Descent

### 3.3 Momentum Algorithm: Accelerating Gradient Descent

### 3.4 Adaptive Learning Rates-Based Algorithms

Adagrad: Adaptive Gradient Algorithm

RMSProp: Root Mean Square Propagation

Adam: Adaptive Moment Estimation

### 3.5 Theoretical Convergence of Optimizers

# The Need for Adaptive Learning Rates

- **Problem with Fixed  $\eta$ :** Same learning rate for all parameters, all time
- **Real-World Data:** Features have different scales and frequencies
- **Sparse Features:** Some features appear rarely but are important

## Why Adaptation Helps:

- **Sparse Gradients:** Rare features need larger updates when they appear
- **Ill-Conditioned Problems:** Loss landscape has different curvatures in different directions
- **Training Dynamics:** Need larger steps initially, smaller steps near convergence

**Goal:** Automate the per-parameter learning rate adjustment during training

# Adagrad: Adaptive per-Parameter Learning

- **Core Idea:** Scale learning rates by historical gradient magnitudes
- Parameters with large gradients get smaller learning rates, and vice versa

## Update Rule:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{\sqrt{\mathbf{G}_t + \epsilon}} \odot \mathbf{g}_t$$

where:

- $\mathbf{G}_t$ : sum of squares of all past gradients (per parameter)
- $\epsilon$ : small constant for numerical stability
- $\mathbf{g}_t = \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t)$ : gradient vector at time step  $t$
- $\odot$ : element-wise multiplication

## Characteristics:

- ✓ Excellent for sparse data and features
- ✗ Learning rate decreases too aggressively over time

# RMSPProp: Fixing Adagrad's Aggressive Decay

- **Improvement:** Use exponentially weighted moving average of gradients
- Prevents learning rate from vanishing too quickly

## Update Rule:

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) \mathbf{g}_t^2$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{\sqrt{\mathbf{v}_t + \epsilon}} \odot \mathbf{g}_t$$

## Key Innovation:

- $\beta$ : decay rate (typically 0.9) - controls how much history to keep
- Recent gradients matter more than ancient ones
- Learning rates adapt to recent gradient behavior

**Result:** More stable and practical than Adagrad

# Adam: Combining Momentum and Adaptive Rates

- **Best of Both Worlds:** Momentum (like physical inertia) + Adaptive learning rates
- **Most Popular:** Default choice for many deep learning applications

## Algorithm <sup>1</sup>:

1. Compute momentum:  $\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$
2. Compute gradient magnitudes:  $\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2$
3. Bias correction:  $\hat{\mathbf{m}}_t = \frac{\mathbf{m}_t}{1 - \beta_1^t}$ ,  $\hat{\mathbf{v}}_t = \frac{\mathbf{v}_t}{1 - \beta_2^t}$
4. Update:  $\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{\sqrt{\hat{\mathbf{v}}_t} + \epsilon} \hat{\mathbf{m}}_t$

**Default Parameters:**  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$

---

<sup>1</sup>Kingma, Diederik P., and Jimmy Ba (2014). “Adam”: A Method for Stochastic Optimization.

# Why a Base Learning Rate $\eta$ is Still Needed

## Example: Simple 1D Loss Function

$$\begin{aligned} L(\theta) &= (\theta - 10)^2, & \frac{dL}{d\theta} &= 2(\theta - 10) \\ \tilde{L}(\theta) &= 100(\theta - 10)^2, & \frac{d\tilde{L}}{d\theta} &= 200(\theta - 10) \end{aligned}$$

### Observation:

- Both have the *same shape* and minimum at  $\theta = 10$ .
- But gradients of  $\tilde{L}$  are  $100\times$  larger everywhere.

### Implication for Optimization:

- Using gradient descent:  $\theta_{t+1} = \theta_t - \alpha \frac{dL}{d\theta}$ .
- For  $\tilde{L}$ , the same  $\alpha$  causes  $100\times$  larger parameter jumps.
- Thus, even adaptive methods (AdaGrad, RMSProp, Adam) still need a *base scale*  $\alpha$ .

# Why Adaptive Methods Still Depend on the Base Learning Rate

$\eta$

**Adaptive update form (e.g., AdaGrad, RMSProp, Adam):**

$$\theta_{t+1} = \theta_t - \eta \frac{m_t}{\sqrt{v_t} + \epsilon}$$

where  $m_t$  and  $v_t$  are moving averages of gradients and squared gradients.

**If gradients are scaled:**  $g_t \rightarrow c g_t$

Then

$$m_t \rightarrow c m_t, \quad v_t \rightarrow c^2 v_t$$

Hence

$$\frac{m_t}{\sqrt{v_t}} \rightarrow \frac{c m_t}{|c| \sqrt{v_t}} = \text{same as before.}$$

**So:** Adaptive methods cancel out local gradient magnitude scaling. But the overall step magnitude remains proportional to  $\eta$ .

**Key Point:**  $\eta$  still decides the *global scale of motion*, while the adaptive term adjusts only the *relative step sizes* between parameters.

# Optimizers with Gradient Descent Variants

## Recall: Gradient Descent Variants:

- **Batch GD** — use the full dataset to compute gradient.
- **Mini-batch GD** — compute gradient on small batches.
- **Stochastic GD (SGD)** — compute gradient on a single sample.

## Optimizers (can be applied to any variant):

- Momentum
- Adagrad
- RMSProp
- Adam

## Key idea:

- Each optimizer only requires a gradient to update parameters.
- Gradient can come from **full dataset, mini-batch, or single sample**.
- Therefore, all optimizers are compatible with all gradient descent variants.



# Choosing Optimizers: Adam vs SGD + Momentum

## Adam: Default Choice for Most Problems

- Handles **sparse gradients** and **non-stationary objectives** well
- Adapts learning rate per parameter automatically
- Fast convergence in deep networks, embeddings, RNNs, Transformers

## Cases Where Simpler Optimizers May Be Better (SGD + Momentum)

- Extremely large, sparse, high-dimensional problems (Adagrad sometimes better)
- Memory constraints — fewer moving averages needed
- Desire for more controlled, predictable convergence

# Adaptive Methods: Summary & Comparison

Method	Key Idea	Best For
<b>Adagrad</b>	Scale by sum of all past gradients	Sparse data, NLP
<b>RMSProp</b>	Exponentially weighted average	Non-stationary objectives
<b>Adam</b>	Momentum + adaptive learning rates	General purpose, deep learning
<b>SGD + Momentum</b>	Physical inertia metaphor	Well-tuned convex problems

## Practical Recommendations:

- **Start with Adam:** Good default for most problems
- **Fine-tune with SGD:** For ultimate performance (with careful tuning)
- **Sparse data:** Consider Adagrad or its variants
- **Use default parameters first**, then experiment if needed

## 3. Famous Optimization Algorithms Used in Deep Learning

### 3.1 Newton-Raphson Algorithm

### 3.2 Gradient Descent and its Variants

Batch Gradient Descent: The Standard One

Stochastic Gradient Descent (SGD)

Mini-batch Gradient Descent

### 3.3 Momentum Algorithm: Accelerating Gradient Descent

### 3.4 Adaptive Learning Rates-Based Algorithms

Adagrad: Adaptive Gradient Algorithm

RMSProp: Root Mean Square Propagation

Adam: Adaptive Moment Estimation

### 3.5 Theoretical Convergence of Optimizers

# Theoretical Convergence for Convex Problems

## Convex Optimization: Convergence Behavior of Optimizers

Optimizer	Convergence Behavior
SGD	Converges to global minimum with proper learning rate decay; slow if learning rate too small.
Momentum	Accelerates convergence in ravines; smooths oscillations.
Adagrad	Convergent due to adaptive learning rate; may become too slow for long training.
RMSProp	Stable convergence; handles noisy gradients well.
Adam	Fast convergence by combining momentum + adaptive learning; widely used in practice.

### Key Points:

- All optimizers eventually reach global minimum (theoretical proofs, under strong assumptions).
- Convergence rate and stability depend on learning rate and batch size.

# Theoretical Convergence for Non-Convex Problems

## Non-Convex Optimization: Convergence Behavior of Optimizers

Optimizer	Convergence Behavior
SGD	May oscillate; can get trapped in local minima.
Momentum	Helps escape shallow local minima; can oscillate in complex landscapes.
Adagrad	Learning rate decays too fast; may stop before reaching good minima.
RMSProp	Adapts learning rate; more stable in non-convex landscapes.
Adam	Fast convergence; widely effective, but may overfit or oscillate if not tuned.

### Key Points:

- Global minimum is not guaranteed.
- Optimizer choice affects ability to escape local minima and convergence speed.

# Summary of Algorithm Evolution

- **Newton-Raphson**: second order, converges fast, but memory heavy
- **Gradient Descent**: first-order, scalable
- **SGD and mini batch GD**: mini-batches, stochastic updates
- **Momentum**: smoother and faster convergence
- **AdaGrad/RMSProp/Adam**: adaptive learning rates for deep networks
- **Key Takeaway**: Evolution driven by dataset size, high-dimensional parameters, and non-convex loss landscapes

*References:* Goodfellow et al., 2016; Ruder, 2016; Kingma & Ba, 2014.

## **Practicals: Training Regression and Classification**

---

# Gradient Descent for Training Binary Classification (Logistic Regression)

- **Loss function** (Binary Cross-Entropy / Log-Loss):

$$\mathcal{L}(\mathbf{w}) = -\frac{1}{n} \sum_{i=1}^n [y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)],$$

where  $\hat{y}_i = \sigma(\mathbf{w}^\top \mathbf{x}_i)$  and  $\sigma(z) = \frac{1}{1+e^{-z}}$ .

- **Gradient:**

$$\nabla \mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i) \mathbf{x}_i$$

- Gradient Descent: **Parameter update rule:**

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \cdot \nabla \mathcal{L}(\mathbf{w}^{(t)})$$

where  $\eta$  is the learning rate.



## Lab: Binary Classification Problems

Look at the Jupyter notebook in the

[https://github.com/yadavrishikesh/Deep-Learning-Slides-Code/blob/main/code/DL\\_Optim](https://github.com/yadavrishikesh/Deep-Learning-Slides-Code/blob/main/code/DL_Optim)

- For **linearly separable** classification problem:  
[DL\\_OptimizationAlgorithm\\_Classification\\_LinearSep.ipynb](#)
- For **non-linearly separable** classification problem:  
[DL\\_OptimizationAlgorithm\\_Classification\\_NotLinearSep.ipynb](#)

# Gradient Descent for Training Linear Regression

- **Loss function** (Mean Squared Error, MSE):

$$\mathcal{L}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad \hat{y}_i = w_0 + w_1 x_i$$

- **Gradients:**

$$\frac{\partial \mathcal{L}}{\partial w_1} = -\frac{2}{n} \sum_{i=1}^n x_i (y_i - \hat{y}_i), \quad \frac{\partial \mathcal{L}}{\partial w_0} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i)$$

- **Parameter updates:**

$$w_1 \leftarrow w_1 - \eta \frac{\partial \mathcal{L}}{\partial w_1}, \quad w_0 \leftarrow w_0 - \eta \frac{\partial \mathcal{L}}{\partial w_0}$$

## Exercise: Regression Problems

Look at the Jupyter notebook in the

[https://github.com/yadavrishikesh/Deep-Learning-Slides-Code/blob/main/code/DL\\_Optim](https://github.com/yadavrishikesh/Deep-Learning-Slides-Code/blob/main/code/DL_Optim) and

- Answer the questions from the Notebook:  
[DL\\_OptimizationAlgorithm\\_Regression\\_Exercise.ipynb](#)