

## AdaGrad (Adaptive Gradient Algorithm)

Parameter-Specific Learning Rate Adaptation

### Core Concept

Adaptively adjusts learning rate  
for each parameter based on historical gradients

### Benefits

- ✓ No manual learning rate tuning needed
- ✓ Advantageous for sparse data

### Update Rule

$$\theta = \theta - \eta / \sqrt{G + \epsilon} \odot \nabla L(\theta)$$

G: Cumulative sum of squared gradients

⊗: Element-wise multiplication

### Main Drawbacks

- ! Learning rate monotonically decreases over time
- ! Aggressive decay can cause premature training stop



### Epsilon Parameter

Prevents division by zero (typically 1e-8)

### Adaptation Mechanism

#### Infrequent Parameters

Larger updates

#### Frequent Parameters

Smaller updates

### Best Use Cases

Sparse Data Processing

Natural Language Processing