

# Machine Learning

## Optimizers

Curtis Larsen

Utah Tech University—Computing

Spring 2025

# Objectives

---

# Outline

## Objectives:

- ▶ Understand convergence
- ▶ Understand and implement optimizers

# Convergence

---

# Idea

- ▶ Gradient descent
- ▶ When does it stop?
- ▶ How long until it stops?
- ▶ Controlled through gradient descent and learning rate

# Math

Gradient descent step.

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (1)$$

# Momentum

---

# Ideas

- ▶ Track sum of gradient values
- ▶ Use sum to update parameters
- ▶ Gradient is an acceleration now, instead of a speed



## Math

$$\mathbf{m} \leftarrow \beta \mathbf{m} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (2)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{m} \quad (3)$$

$$0 \leq \beta \leq 1$$

# Implementation

```
keras.optimizers.SGD(momentum=0.9)
```

# Nesterov

---

# Ideas

- ▶ Look into future for gradient

## Math

$$\mathbf{m} \leftarrow \beta \mathbf{m} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta} + \beta \mathbf{m}) \quad (4)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{m} \quad (5)$$

$$0 \leq \beta \leq 1$$

# Implementation

```
keras.optimizers.SGD(momentum=0.9,  
                      nesterov=True)
```

# AdaGrad

---

# Ideas

- ▶ Scale updates in each direction based on historical gradient size.



# Math

$$\mathbf{s} \leftarrow \mathbf{s} + \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \otimes \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (6)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \div \sqrt{\mathbf{s} + \epsilon} \quad (7)$$

# Implementation

```
keras.optimizers.Adagrad()
```

# RMSProp

---

# Ideas

- ▶ Scale updates in each direction based on recent history.

## Math

$$\mathbf{s} \leftarrow \rho \mathbf{s} + (1 - \rho) \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \otimes \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (8)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \div \sqrt{\mathbf{s} + \epsilon} \quad (9)$$

# Implementation

```
keras.optimizers.RMSprop(rho=0.9)
```

# Adam

# Ideas

- ▶ Track mean historical gradient (first moment)
- ▶ Track variance in historical gradient (second moment)
- ▶ Scale steps based on them



## Math

$$\mathbf{m} \leftarrow \beta_1 \mathbf{m} + (1 - \beta_1) \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (10)$$

$$\mathbf{s} \leftarrow \beta_2 \mathbf{s} + (1 - \beta_2) \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \otimes \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (11)$$

$$\hat{\mathbf{m}} \leftarrow \frac{\mathbf{m}}{1 - \beta_1^t} \quad (12)$$

$$\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \beta_2^t} \quad (13)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta \hat{\mathbf{m}} \div \sqrt{\hat{\mathbf{s}} + \epsilon} \quad (14)$$

# Implementation

```
keras.optimizers.Adam(beta_1=0.9,  
                        beta_2=0.999)
```

# Summary

---

# Topics

- ▶ Momentum
- ▶ RMSProp
- ▶ Adam