# Convergence Analysis of Optimizers
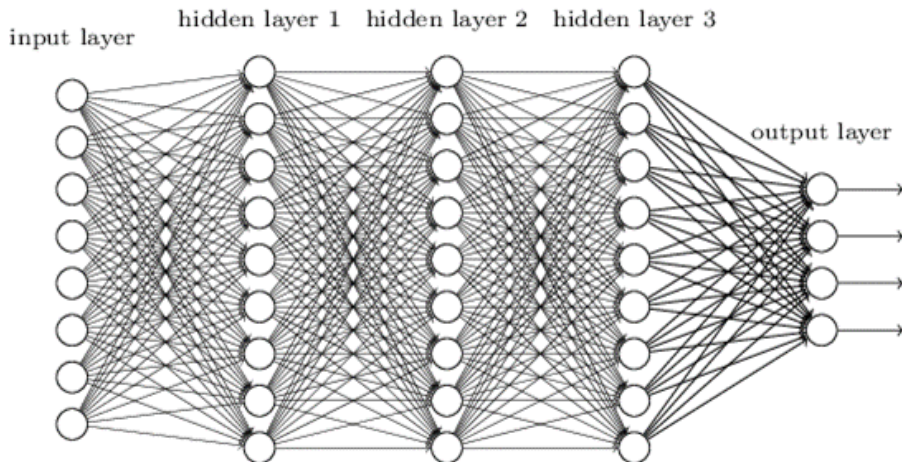
Anthony Garcia

# Contents (Contribution)

# 1. Deep Learning Optimization Overview

▶ The optimization problem in deep learning is non-convex.
▶ Through initialization, (dynamic) learning rate, and stochastic elements, we can make gradient descent-based methods converge to the global minimum.

# 2. Unified Perspective for Optimizers (GD)

$$\mathbf{x_{t+1}} = \mathbf{x_t} - \alpha \nabla \mathbf{f(x)}$$

- ▶ Step size (learning rate) : $\alpha$
- ▶ Step direction : $\nabla f(x)$

모든 자료를 다 검토해서
내 위치의 산기울기를 계산해서
갈 방향을 찾겠다.

**GD**

**Momentum**
스텝 계산해서 움직인 후,
아까 내려 오던 관성 방향 또 가자

스텝방향

**SGD**
전부 다봐야 한걸음은
너무 오래 걸리니까
조금만 보고 빨리 판단한다
같은 시간에 더 많이 간다

스텝사이즈

**Adam**
RMSProp + Momentum
방향도 스텝사이즈도 적절하게!

**RMSProp**
보폭을 줄이는 건 좋은데
이전 맥락 상황봐가며 하자.

**Adagrad**
안가본곳은 성큼 빠르게 걸어 훑고
많이 가본 곳은 잘아니까
갈수록 보폭을 줄여 세밀히 탐색

# 2. Unified Perspective for Optimizers (SGD)

**Stochastic Gradient Descent (Mini-batch)**

- $g_t = \nabla \tilde{f}_t(x) = \dfrac{1}{t} \sum_{i=1}^{t} \nabla f_i(x_t)$

- $x_{t+1} = x_t - \alpha g_t$

Mini-batch: A training method dividing the entire dataset into small subsets

- Noise Injection $\rightarrow$ escape local minima, find global minima
- Stochastic Learning
- Computational Efficiency
- Memory Efficiency

# 2. Unified Perspective for Optimizers (SGDM)

**Stochastic gradient descent with momentum**

- $g_t = \dfrac{1}{t} \sum_{i=1}^{t} \nabla f_i(x_t)$
- $m_t = \beta m_{t-1} + (1 - \beta)g_t$
- $x_{t+1} = x_t - \alpha m_t$

Momentum: Inertia from previous gradients + current gradient

- $m_t = (0.081 g_{t-2} + 0.09 g_{t-1}) + 0.1 g_t$, $(\beta = 0.9)$
- SGDM is SGD with inertia (SGD is special case of SGDM in $\beta = 1$)
- Reduction of oscillations

# 2. Unified Perspective for Optimizers (Adagrad)

**Adagrad : dynamic learning rate based on the accumulated square sum of gradients**

- $g_i^{(k)} = \dfrac{1}{t} \sum_{j=1}^{k} \nabla f_i(x_j)$

- $v_i^{(k)} = \sum_{j=1}^{k} (g_i^{(j)})^2$

- $x_i^{(k+1)} = x_i^{(k)} - \dfrac{\alpha}{\sigma + \sqrt{v_i^{(k)}}} g_i^{(k)}$

  ($\sigma$ is a small value to prevent division by zero)

Dynamic learning rate: less sensitive to learning rate

- Adagrad is SGD with dynamic learning rate
  (SGD is special case of Adagrad in $v_i^{(k)} = 0$)
- large step size initially, small step size near optimum

## 2. Unified Perspective for Optimizers (RMSProp)

**RMSProp : exponentially weighted moving average**

- $g_i^{(k)} = \dfrac{1}{t} \sum_{j=1}^{k} \nabla f_i(x_j)$

- $v_i^{(k+1)} = \beta_2 v_i^{(k)} + (1 - \beta_2)(g_i^{(k)})^2$

- $\hat{v}_i^{(k+1)} = \dfrac{v_i^{(k+1)}}{1 - \beta_2^{k+1}}$

- $x_i^{(k+1)} = x_i^{(k)} - \dfrac{\alpha}{\sigma + \sqrt{\hat{v}_i^{(k+1)}}} g_i^{(k)} = x_i^{(k)} - \dfrac{\alpha}{\sigma + \mathsf{RMS}(g_i)} g_i^{(k)}$

Exponentially weighted moving average: not monotonically decreasing

- to prevent the learning rate from decreasing too rapidly in Adagrad
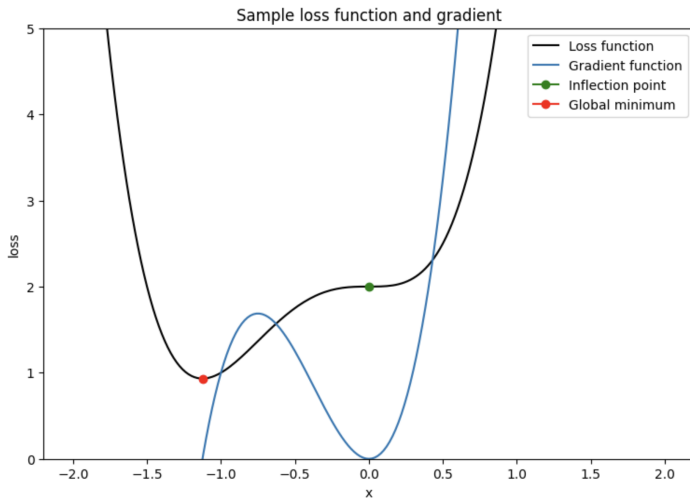
# 2. Unified Perspective for Optimizers (Adam)

**Adam(Momentum + RMSProp)**

- $g_i^{(k)} = \dfrac{1}{t} \sum_{j=1}^{k} \nabla f_i(x_j)$
- $m^{(k+1)} = \beta_1 m^{(k)} + (1 - \beta_1) g_i^{(k)}$
- $v_i^{(k+1)} = \beta_2 v_i^{(k)} + (1 - \beta_2)(g_i^{(k)})^2$
- $\hat{m}^{(k+1)} = \dfrac{m^{(k+1)}}{1 - \beta_1^{k+1}}$
- $\hat{v}_i^{(k+1)} = \dfrac{v_i^{(k+1)}}{1 - \beta_2^{k+1}}$
- $x_i^{(k+1)} = x_i^{(k)} - \dfrac{\alpha}{\sigma + \sqrt{\hat{v}_i^{(k+1)}}} m^{(k+1)} = x_i^{(k)} - \dfrac{\alpha}{\sigma + \mathsf{RMS}(g_i)} m^{(k+1)}$
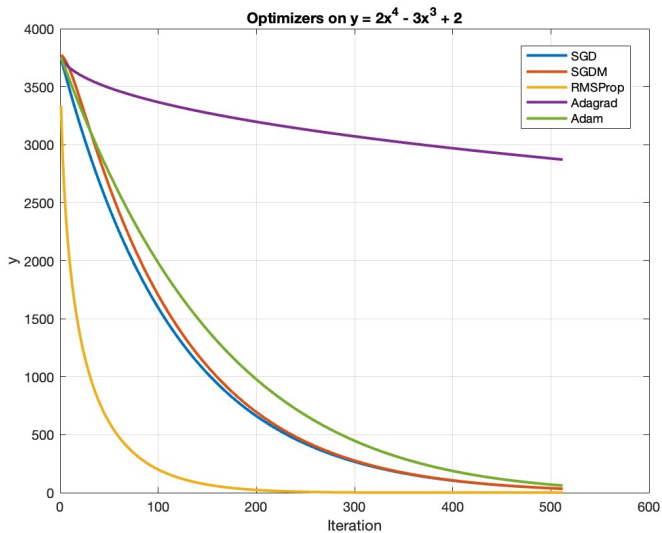
Momentum + RMSProp:
- Adam is the most generalized form among all optimizers (SGDM, Adagrad, RMSProp are special cases of Adam)

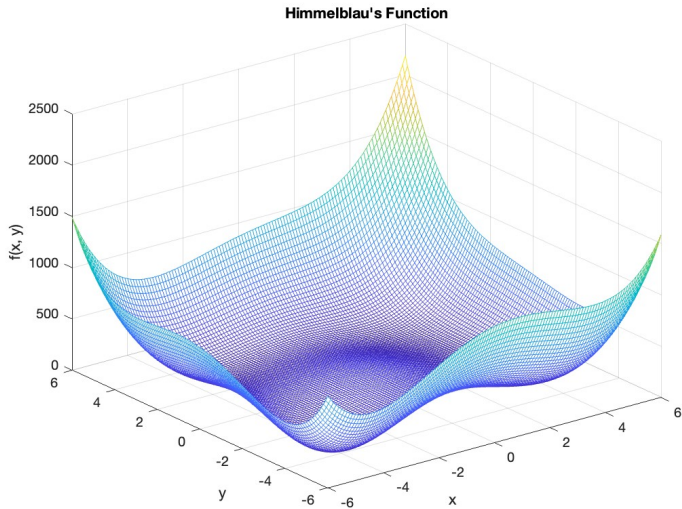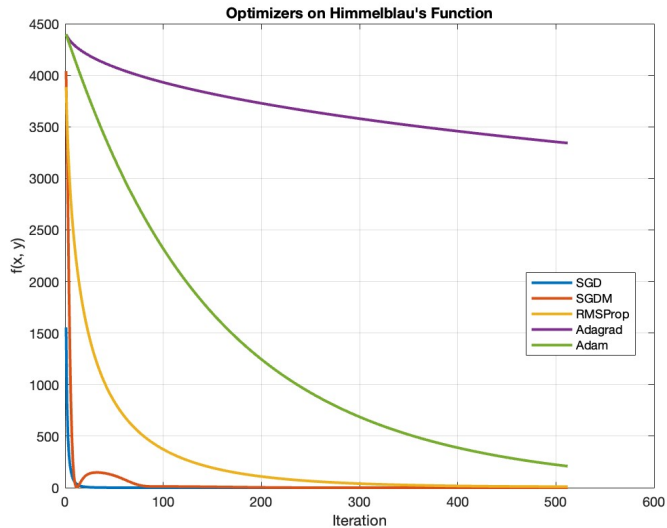# 3. Polynomial Example



Sample loss function and gradient

# 3. Polynomial Example

# 3. Polynomial Example



Himmelblau's Function

# 3. Polynomial Example

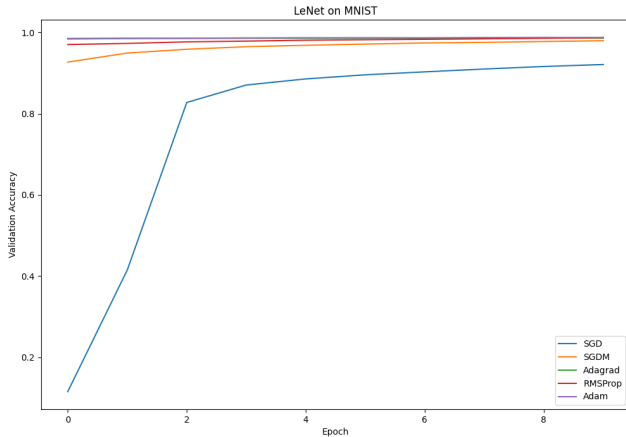# 4. Broad Range of Experiments

**Multiple Datasets, Different Model Architectures, and various optimizer**

| Dataset | Model Architecture |
|---------------|--------------------|
| Mnist | LeNet-5 |
| Fashion-mnist | AlexNet |
| Cifar-10 | VGGNet |

# 4. Broad Range of Experiments

# 4. Broad Range of Experiments



LeNet on MNIST

# 4. Broad Range of Experiments



Alexnet on Fashion MNIST

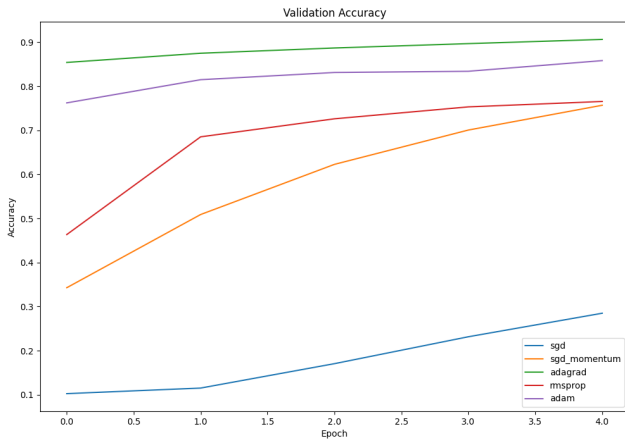# 4. Broad Range of Experiments



Alexnet on Fashion MNIST
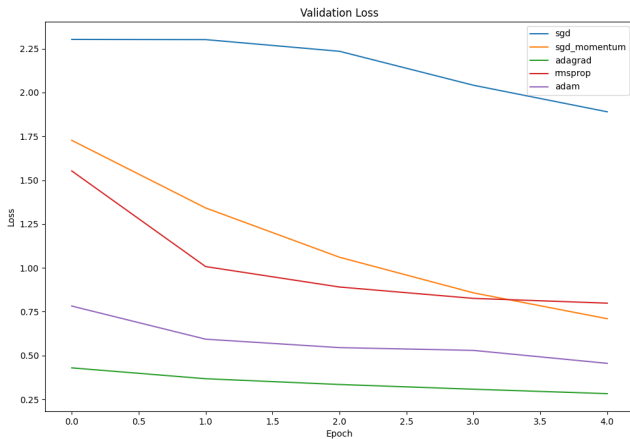
# 4. Broad Range of Experiments

# 4. Broad Range of Experiments

# References

▶ Oh Ilseok, Machine Learning

▶ Mykel J. Kochenderfer, Algorithms for Optimization

▶ Yanli Liu, An Improved Analysis of Stochastic Gradient Descent with Momentum(2020)

▶ John Duchi, Adaptive Subgradient Methods for Online Learning and Stochastic Optimization(2011)

▶ Alexandre Défossez, A Simple Convergence Proof of Adam and Adagrad(2022)

▶ Rachel Ward, AdaGrad stepsizes: Sharp convergence over nonconvex landscapes(2020)

▶ Fangyu Zou, A Sufficient Condition for Convergences of Adam and RMSProp

▶ Diederik P. Kingma, Adam: A Method for Stochastic Optimization(2015)

▶ Pierre Foret, Sharpness-Aware Minimization for Efficiently Improving Generalization(2021)