

Optimization

Outline

- Recap minibatch SGD

- Discuss semester timeline

1 Some prerequisites (SGD, contour plots, EMA)

2 Momentum

3 Adaptive learning rate

- 1 • AdaGrad (Adaptive Gradients)

- 2 • RMSProp (Root-Mean-Square Backpropagation)

- 3 • Adam (Adaptive Moment Estimation)

Recap: Minibatch SGD

- Take five minutes to draw
 - Whatever will help you remember (no correct or incorrect drawings)
 - You'll keep a running drawing log the rest of the semester

Semester vs. Topic Timelines

↓ Foundations

- Math
 - Calculus
 - Linear algebra
- Programming
 - Python
 - Notebooks
 - Libraries
- Computing
 - HPC
 - CLI

NN Basics

- Terminology
- History
- Ethics
- Neurons
- Networks
- Backpropagation
- Autodiff
- Gradient descent

NN Intermediate

- Optimization
- Initialization
- Normalization
- Overfitting
- Convolutions
- Recurrent

NN Advanced

- Transfer learning
- Inference
- GANs
- Reinforcement
- Attention
- Transformers
- Stable diffusion

Done!

MB-SGD Equations

W, b
parameters
↓
 Θ

updated params
↓
 Θ_{t+1} := $\Theta_t - \alpha \frac{\nabla L(\hat{y}_b, y_b)}{\nabla \Theta_t}$

CURRENT params
↓
 Θ_t

parameter gradients
↑
 $\nabla L(\hat{y}_b, y_b)$

learning rate
↑
 α

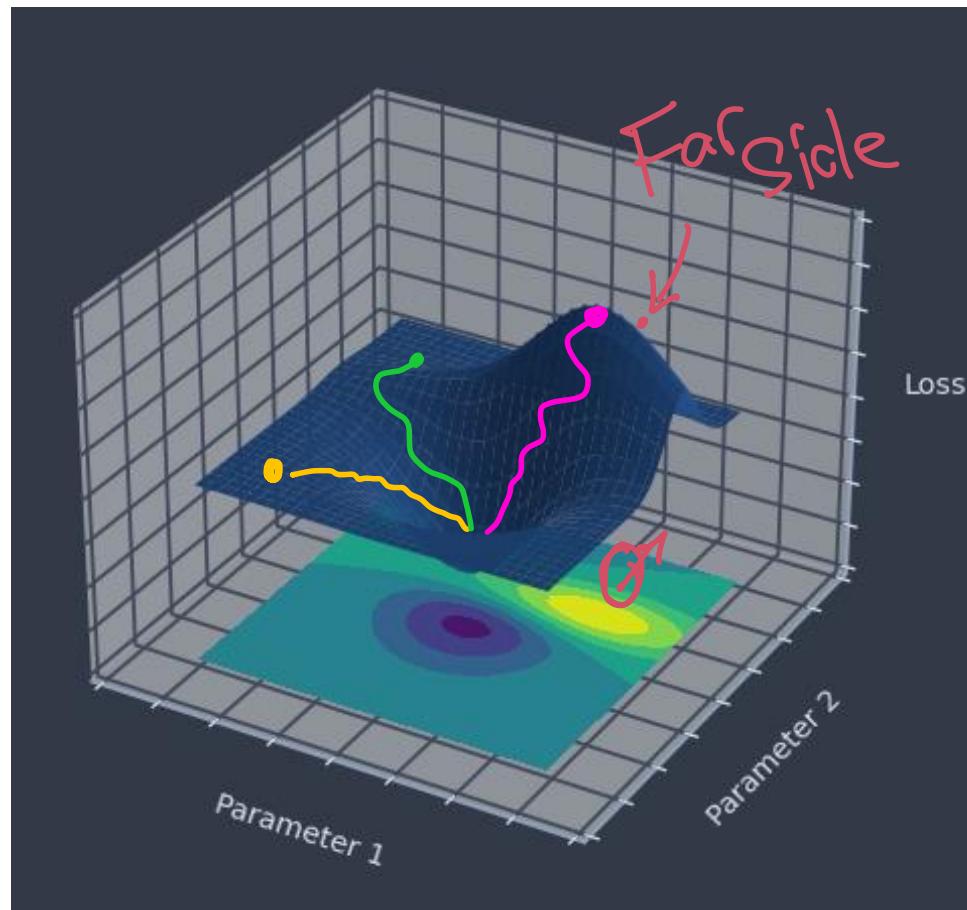
assignment
↑
 Θ_{t+1}

param. - learning_rate · param. grad

g_t

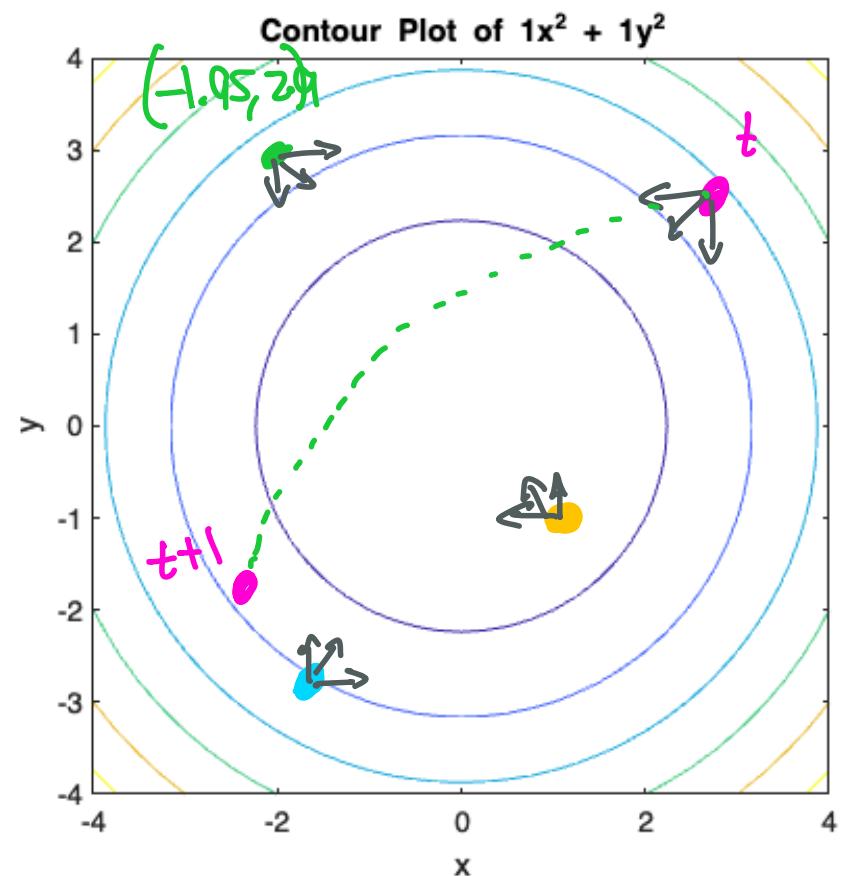
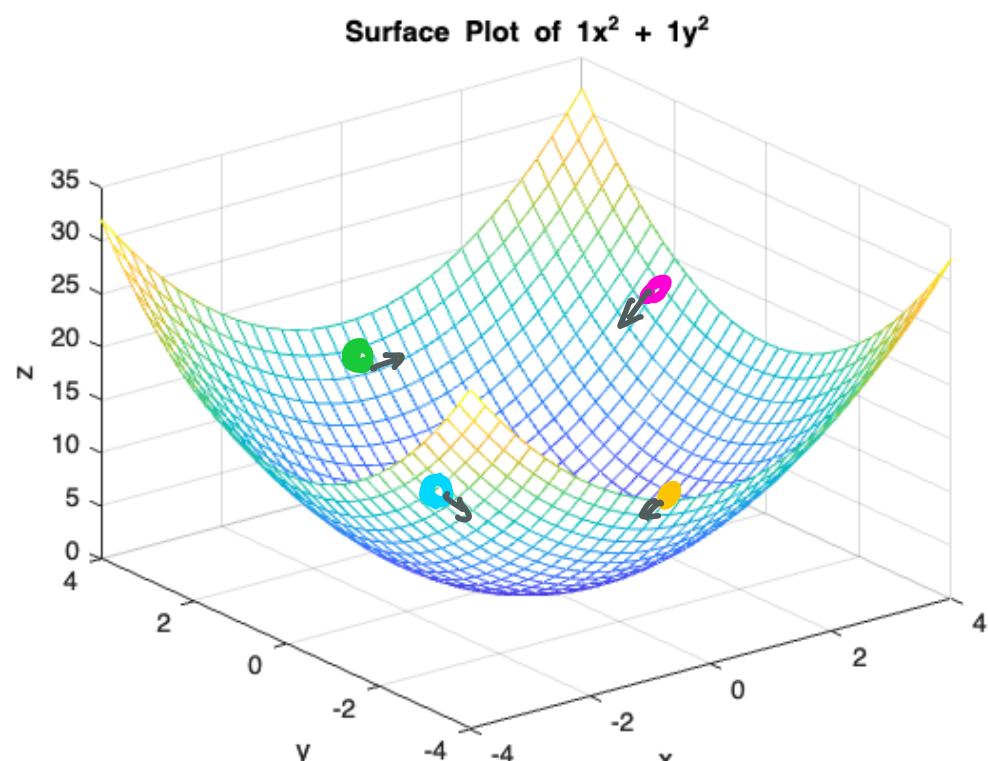
How do we choose a learning rate?

Contour Plots

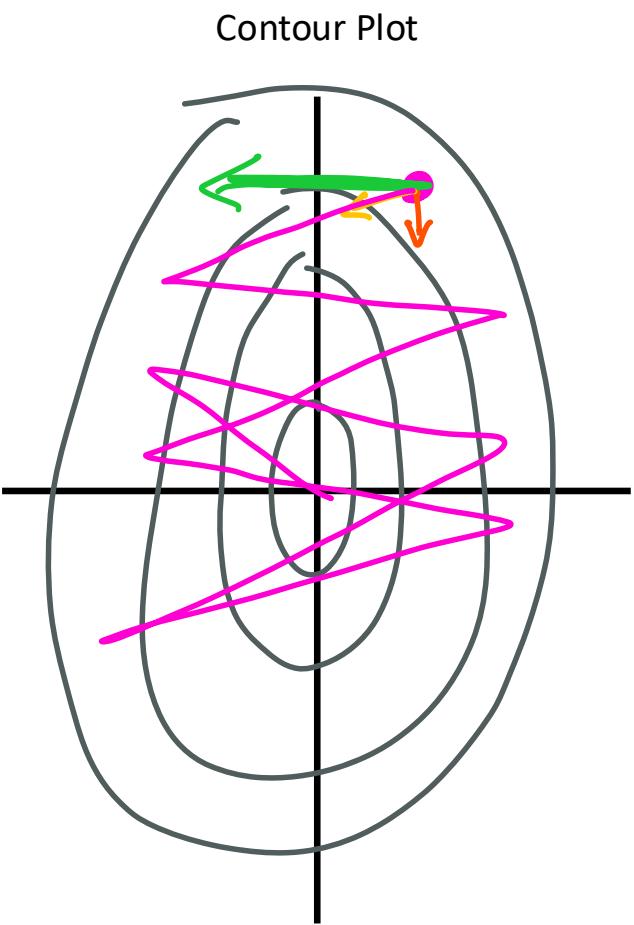
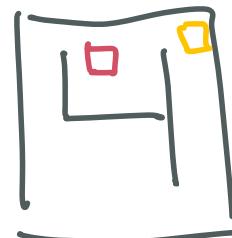
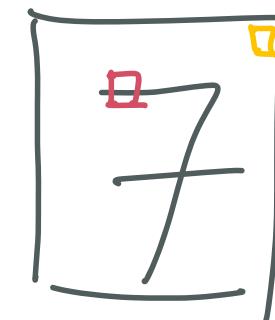
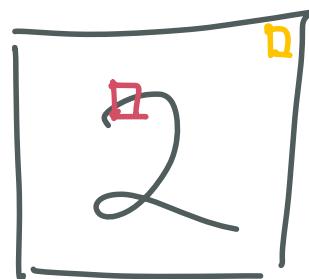
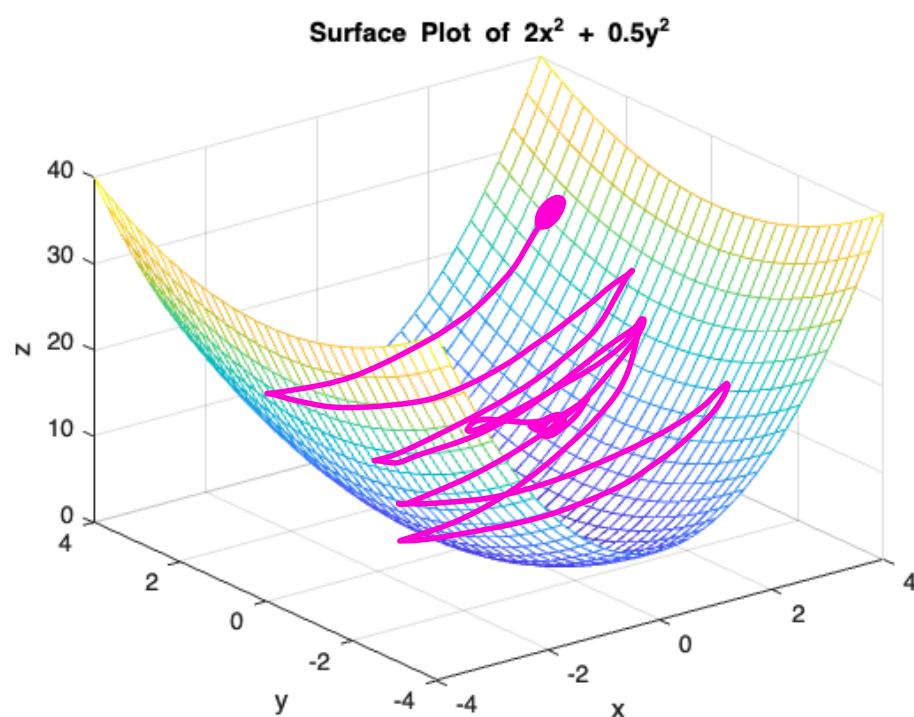


$$\Theta_{t+1} := \Theta_t - \alpha g_t$$

Uniform Gradients



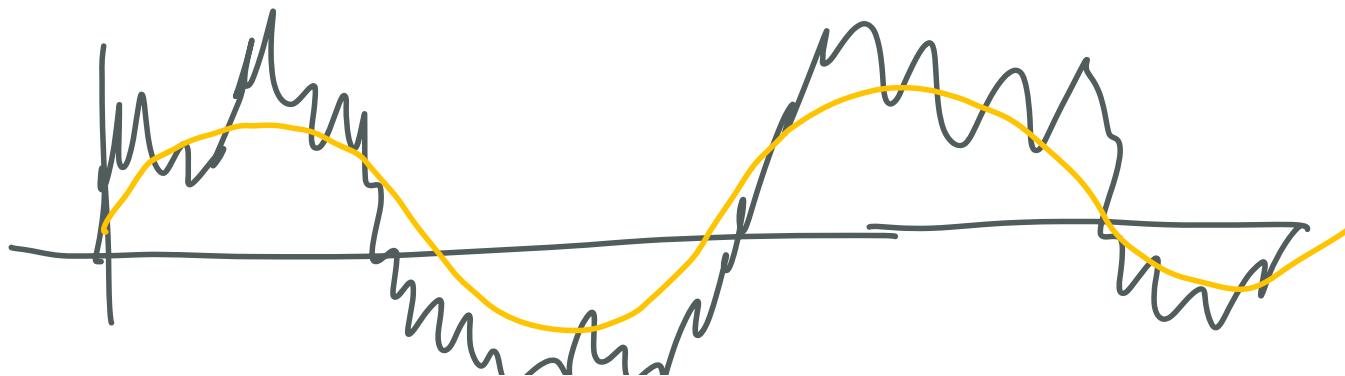
Contour Plots



Exponential Moving Average (EMA)

- A simple moving average is the unweighted mean of the previous k values
- EMA is a rule-of-thumb technique for smoothing time series data

smoothed value $\rightarrow \bar{y}_{t+1} := \beta \bar{y}_t + (1-\beta) y_t$ *next value in the sequence*

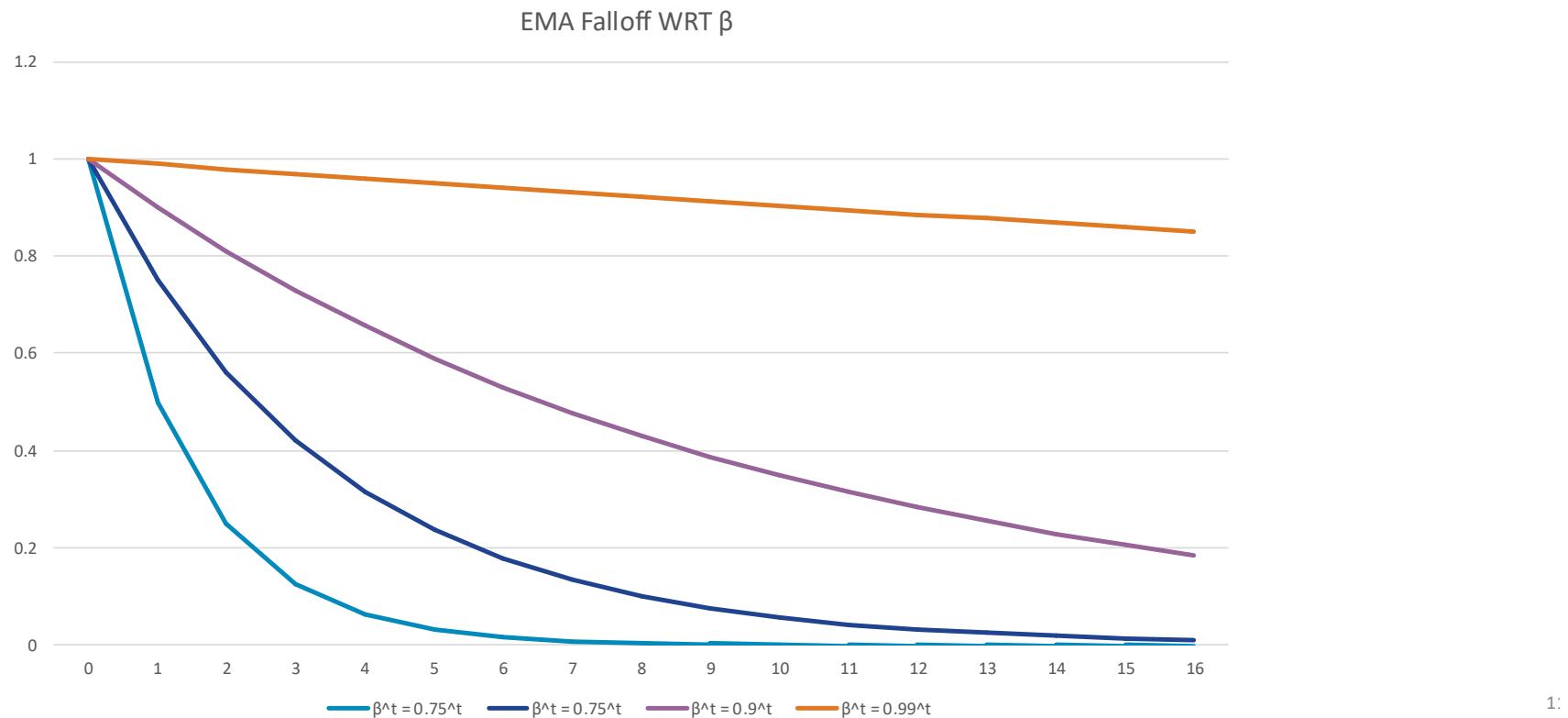


Exponential Moving Average (EMA)

- <https://cs.pomona.edu/classes/cs181r/book/19-SensorFusion.html>

Exponential Moving Average (EMA)

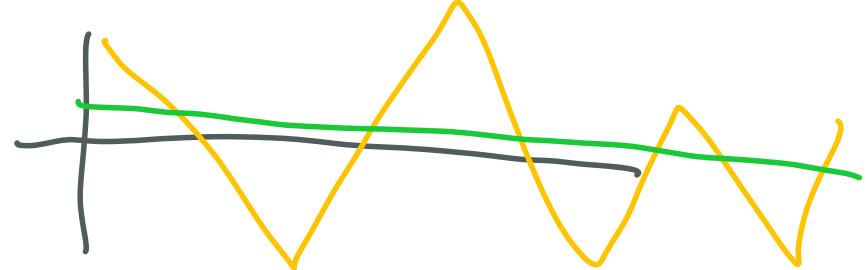
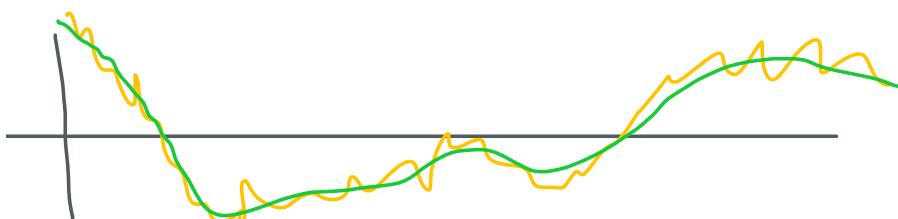
- A simple moving average is the unweighted mean of the previous k values
- EMA is a rule-of-thumb technique for smoothing time series data



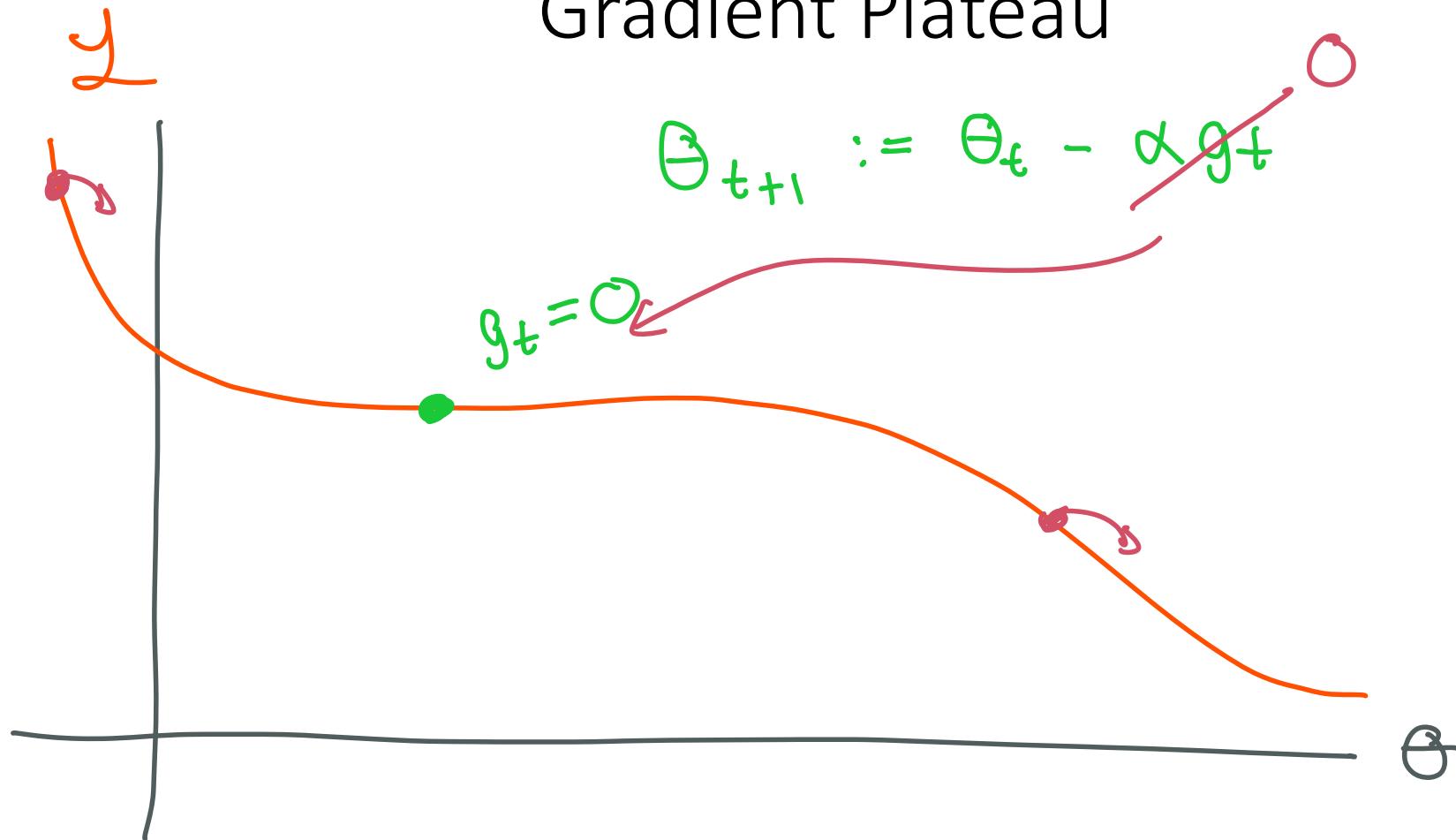
$$\bar{y}_{t+1} := \beta \bar{y}_t + (1-\beta) y_t \quad 0 < \beta < 1$$

$$\bar{y}_{t+1} := \beta (\beta \bar{y}_{t-1} + (1-\beta) y_{t-1}) + (1-\beta) y_t$$

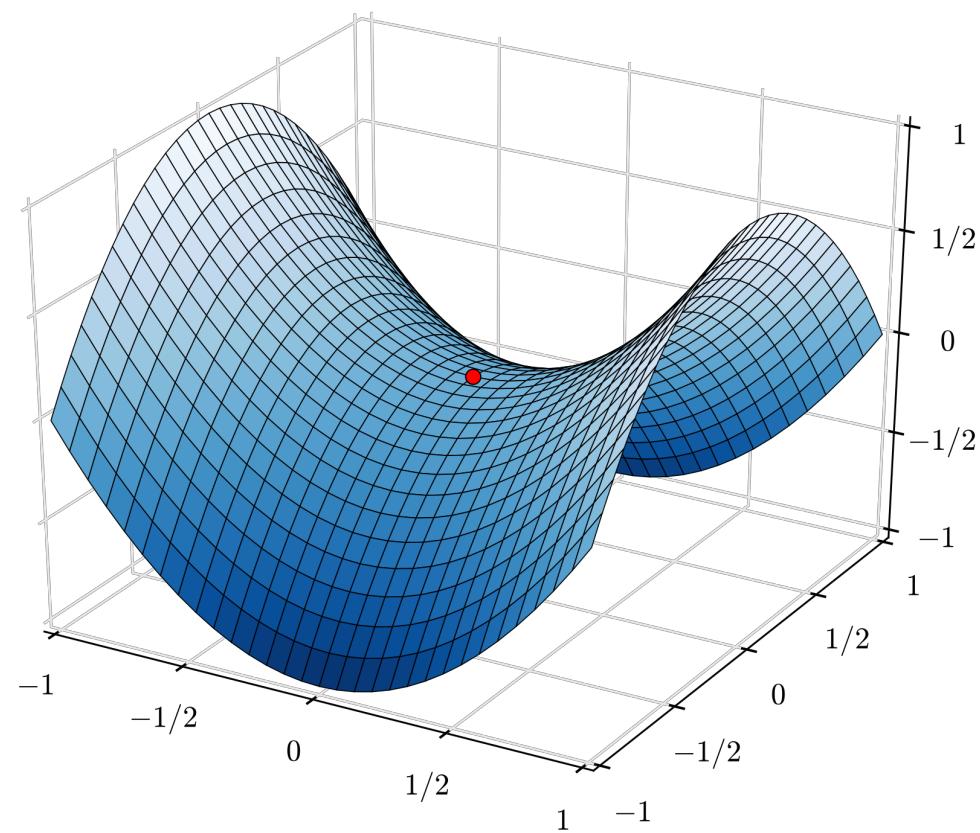
$$\begin{aligned}\bar{y}_{t+1} &:= \beta (\beta (\beta \bar{y}_{t-2} + (1-\beta) y_{t-2}) + (1-\beta) y_{t-1}) + (1-\beta) y_t \\ &:= \beta^3 \bar{y}_{t-2} + \beta^2 (1-\beta) \bar{y}_{t-2} + \beta (1-\beta) \bar{y}_{t-1} + (1-\beta) y_t\end{aligned}$$



Gradient Plateau



Saddle Points

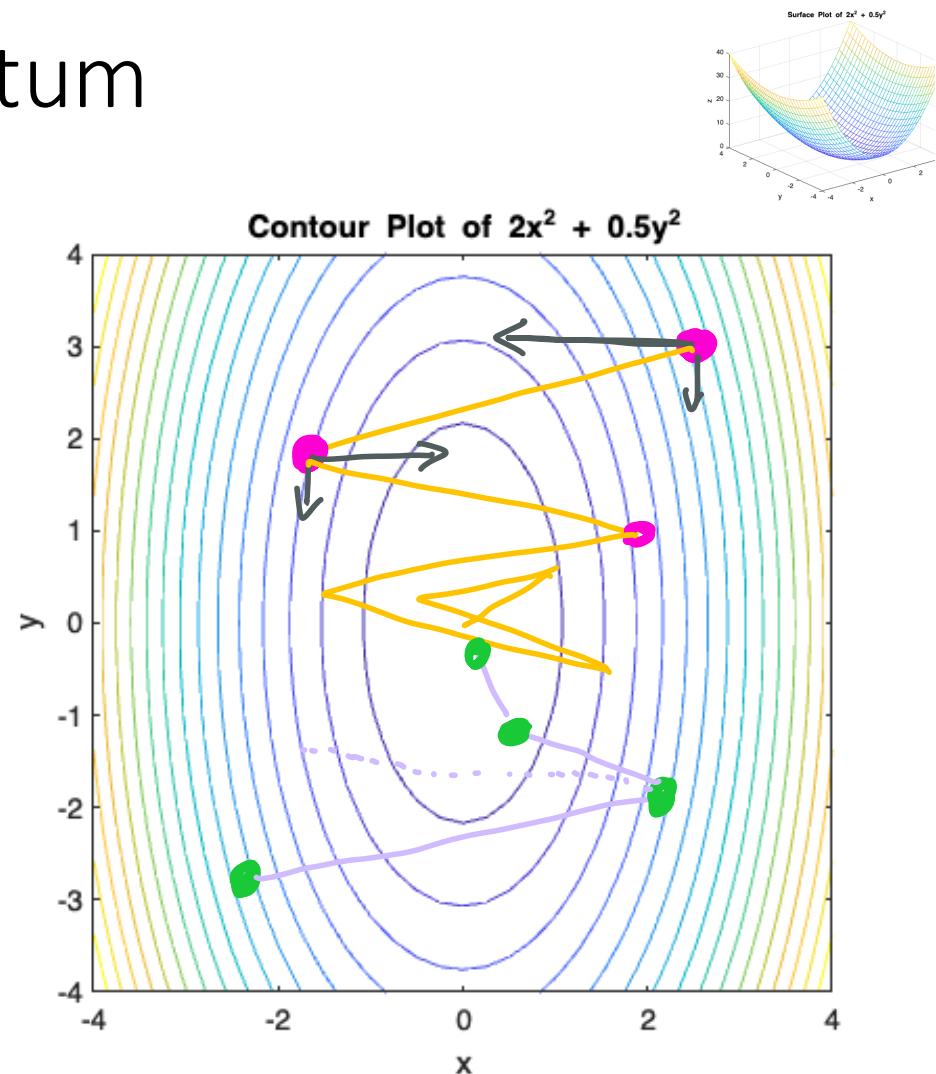


exponential moving average

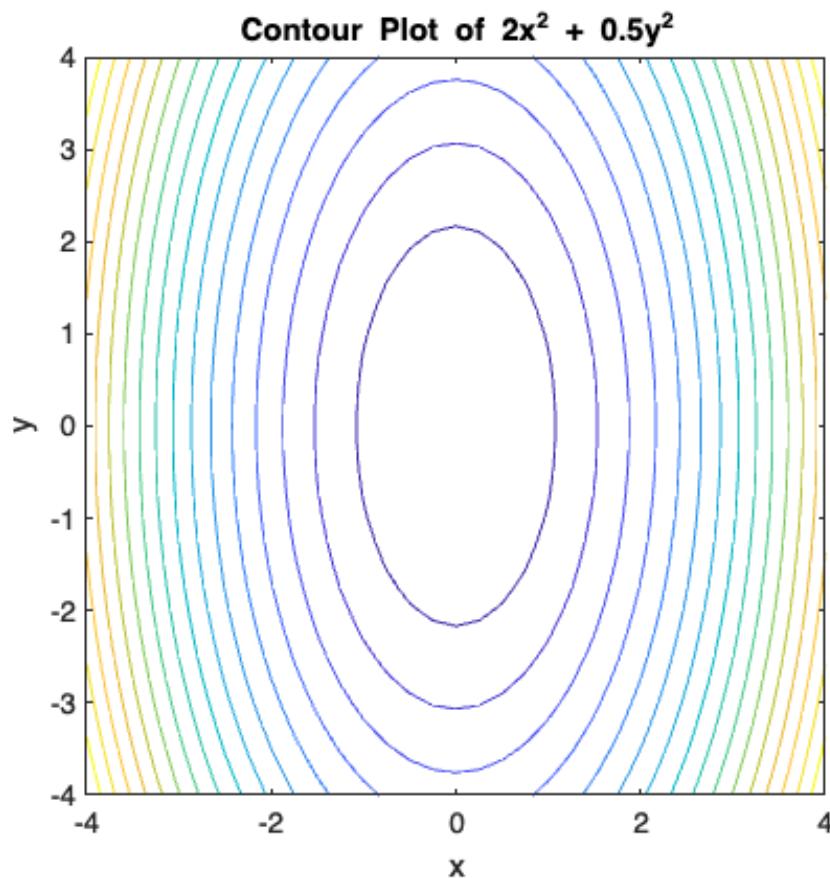
$$\text{Momentum} = \beta_m m_{t-1} + (1 - \beta_m) g_{t-1}$$

$$m_{t+1} := \beta_m m_t + (1 - \beta_m) g_t$$

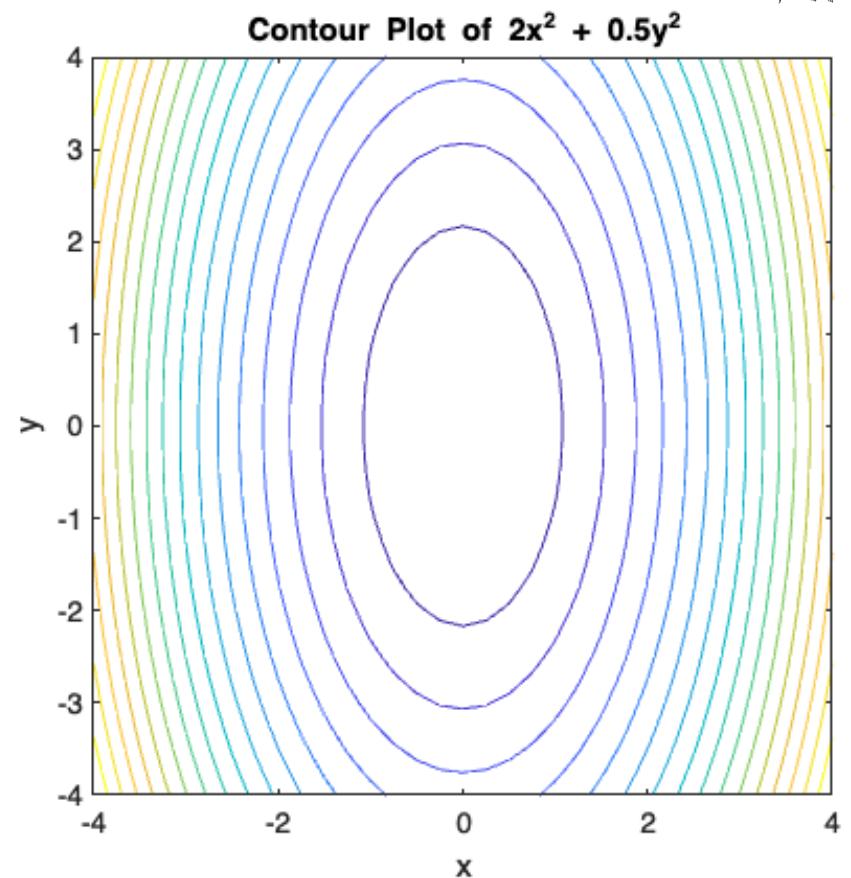
$$\theta_{t+1} := \theta_t - \alpha m_{t+1}$$



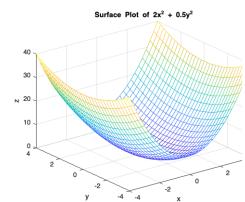
Momentum



Without Momentum

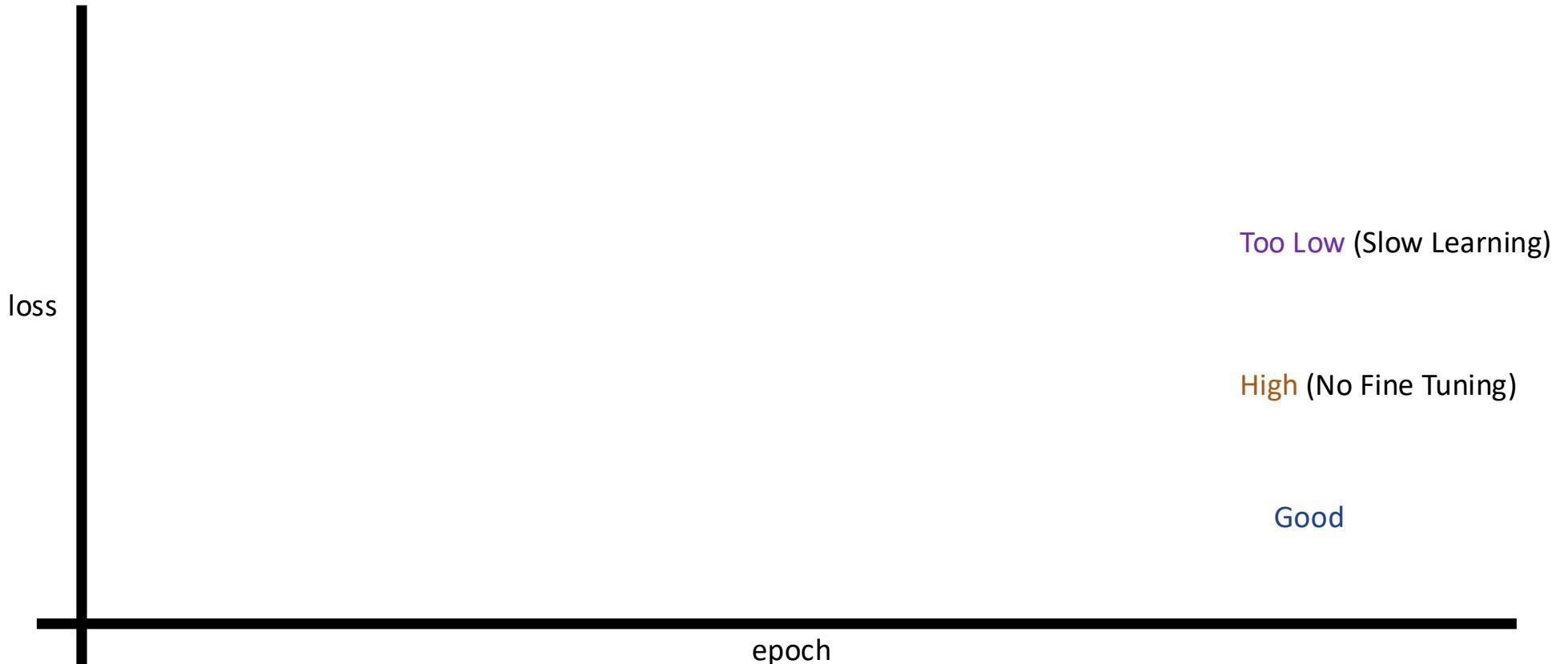


With Momentum



Learning Rates

Too High (No Learning)



AdaGrad

RMSProp

$z = Aw^T + b \rightarrow \text{Linear}$

$A = \sigma(z) \rightarrow \text{Sigmoid}$

Adam

Recent Optimizers

- SGD (1951)
 - SGD+Momentum (1999)
 - AdaGrad (2011)
 - AdaDelta (2012)
 - RMSProp (2013)
 - Adam (2014)
 - NADAM (2015)
 - AdamW (2017)
 - AdaShift (2018)
 - AggMo (2018)
 - LAMB (2019)
 - AMSGrad (2019)
 - Adabelief (2020)
 - MADGRAD (2021)
 - AdaSmooth (2022)
- Just a sample
- See: <https://johnchenresearch.github.io/demon/> for more information