

Optimization

Outline

- Recap minibatch SGD
- Discuss semester timeline
- Some prerequisites (SGD, contour plots, EMA)
- Momentum
- Adaptive learning rate
 - AdaGrad (Adaptive Gradients)
 - RMSProp (Root-Mean-Square Backpropagation)
 - Adam (Adaptive Moment Estimation)

Recap: Minibatch SGD

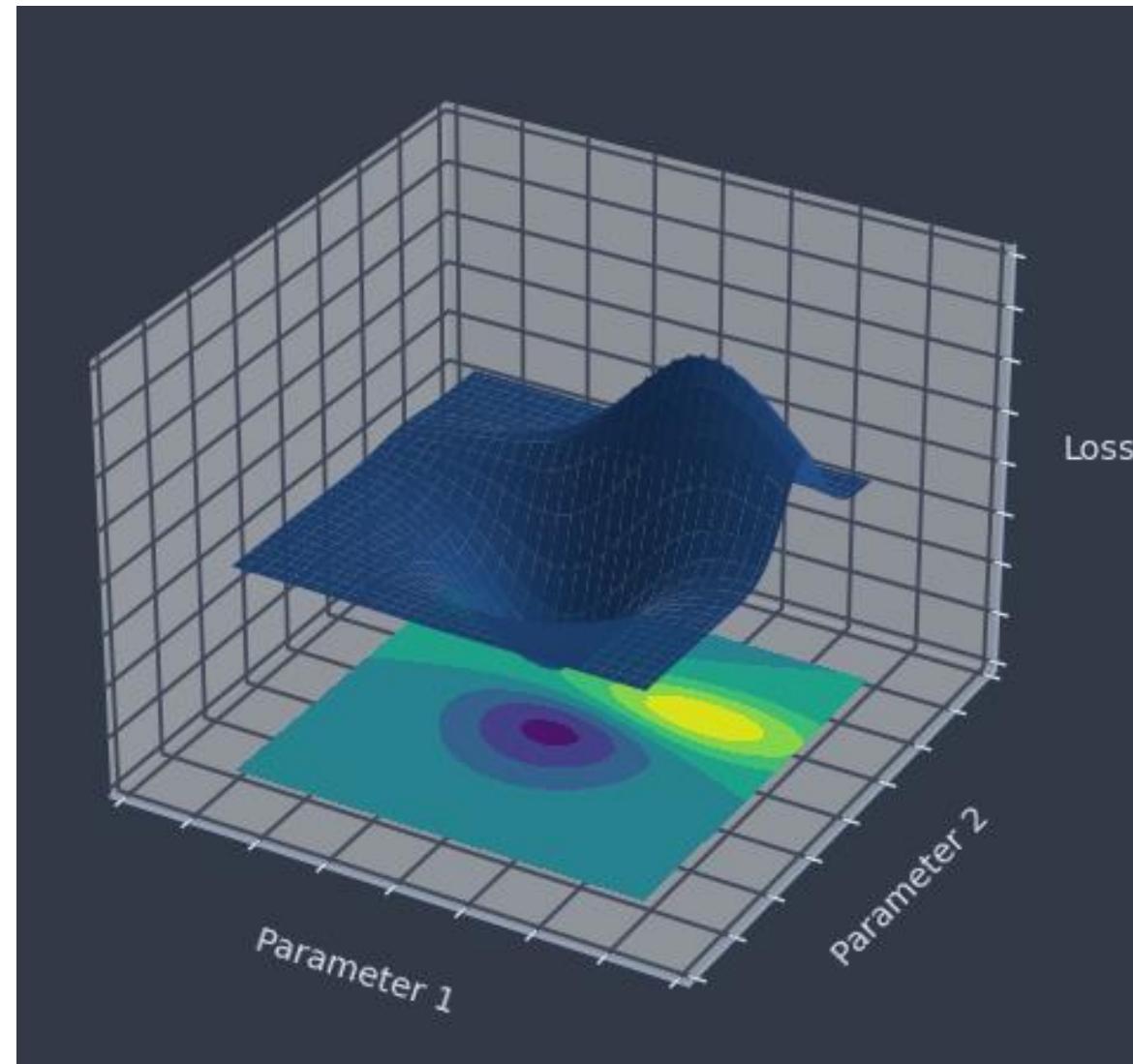
- Take five minutes to draw
 - Whatever will help you remember (no correct or incorrect drawings)
 - You'll keep a running drawing log the rest of the semester

Semester vs. Topic Timelines

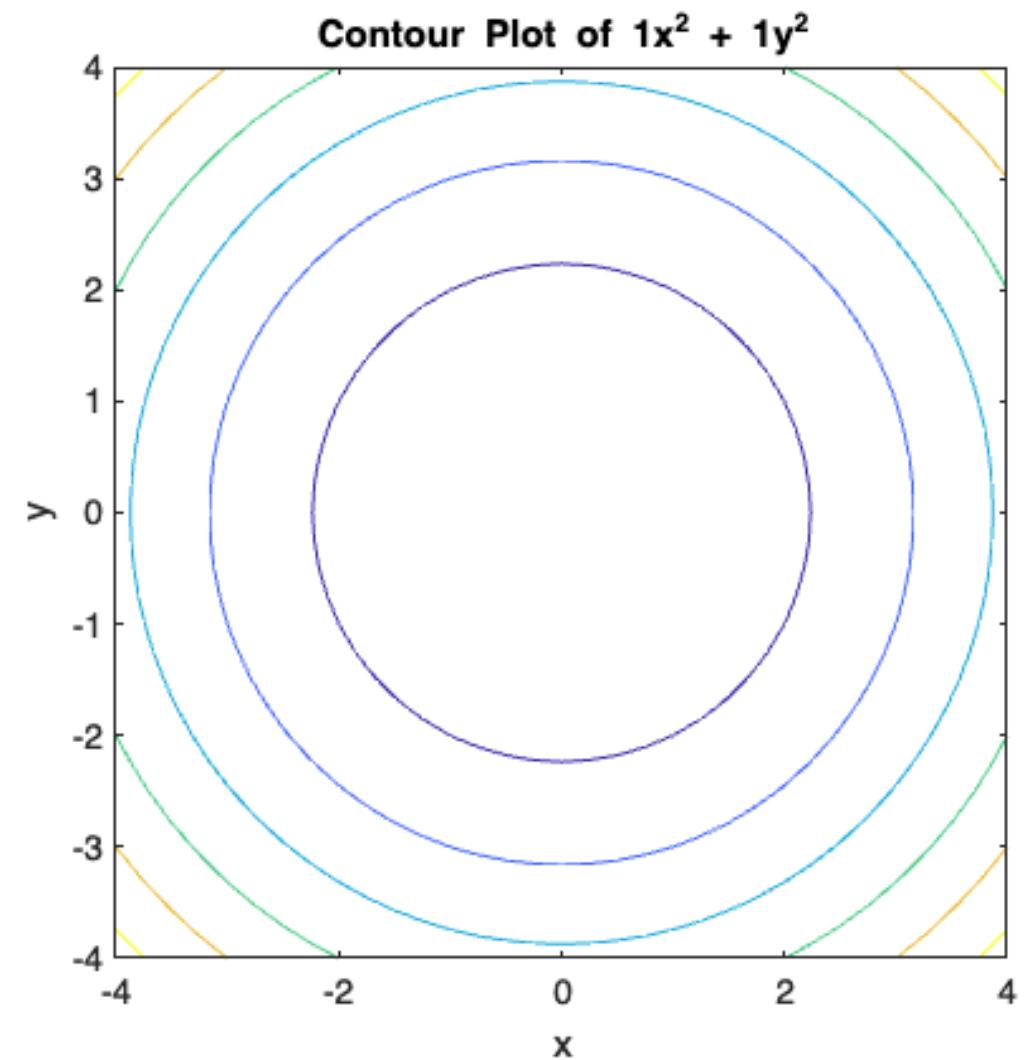
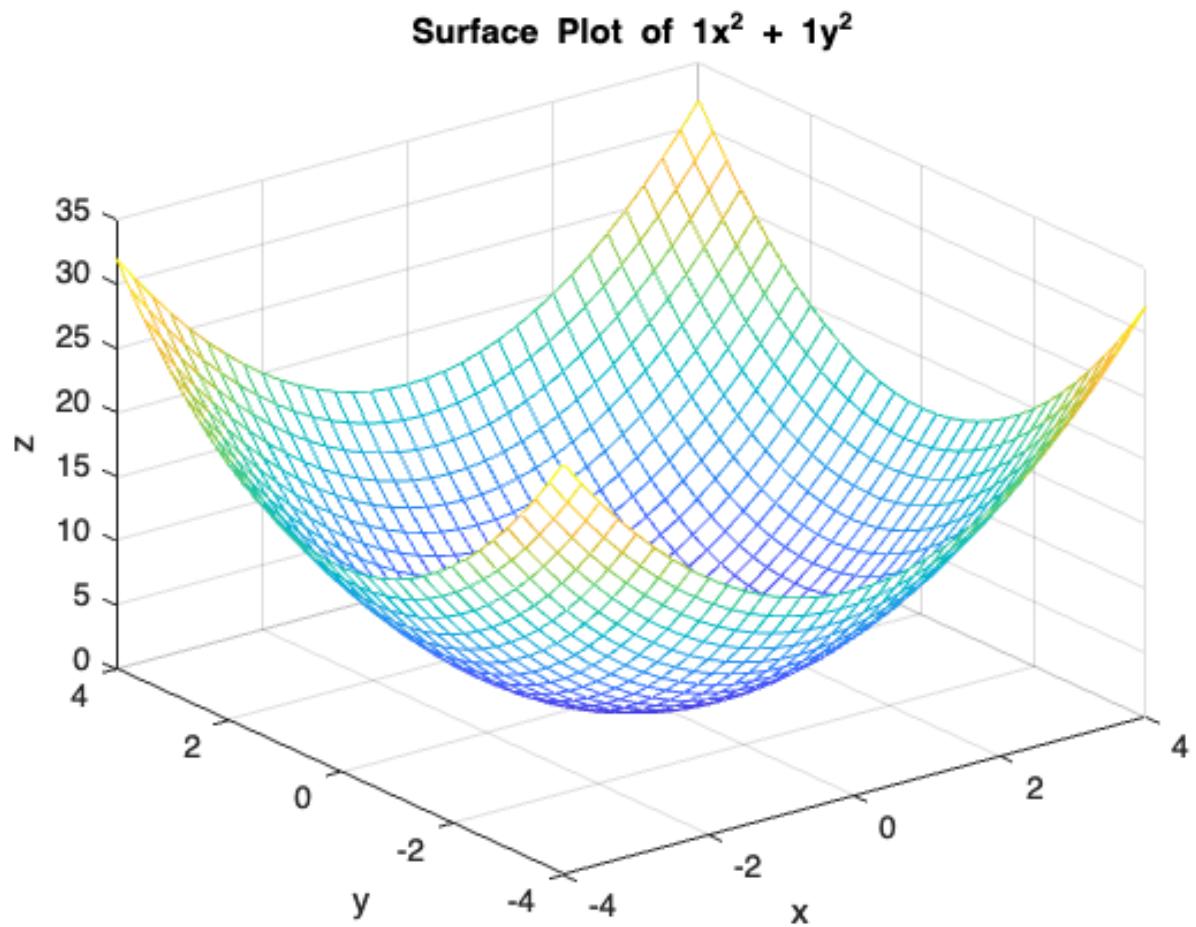
<u>Foundations</u>	<u>NN Basics</u>	<u>NN Intermediate</u>	<u>NN Advanced</u>
<ul style="list-style-type: none">• Math<ul style="list-style-type: none">• Calculus• Linear algebra• Programming<ul style="list-style-type: none">• Python• Notebooks• Libraries• Computing<ul style="list-style-type: none">• HPC• CLI	<ul style="list-style-type: none">• Terminology• History• Ethics• Neurons• Networks• Backpropagation• Autodiff• Gradient descent	<ul style="list-style-type: none">• Optimization• Initialization• Normalization• Overfitting• Convolutions• Recurrent	<ul style="list-style-type: none">• Transfer learning• Inference• GANs• Reinforcement• Attention• Transformers• Stable diffusion

SGD Equations

Contour Plots

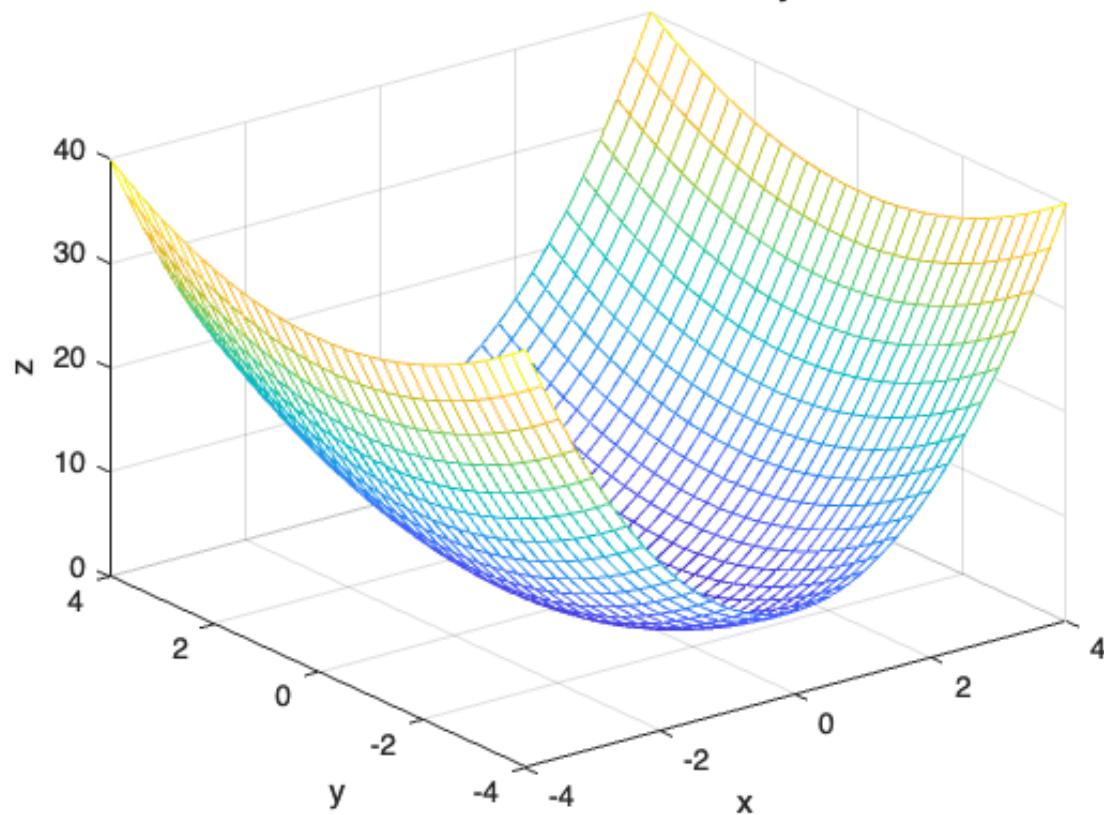


Uniform Gradients

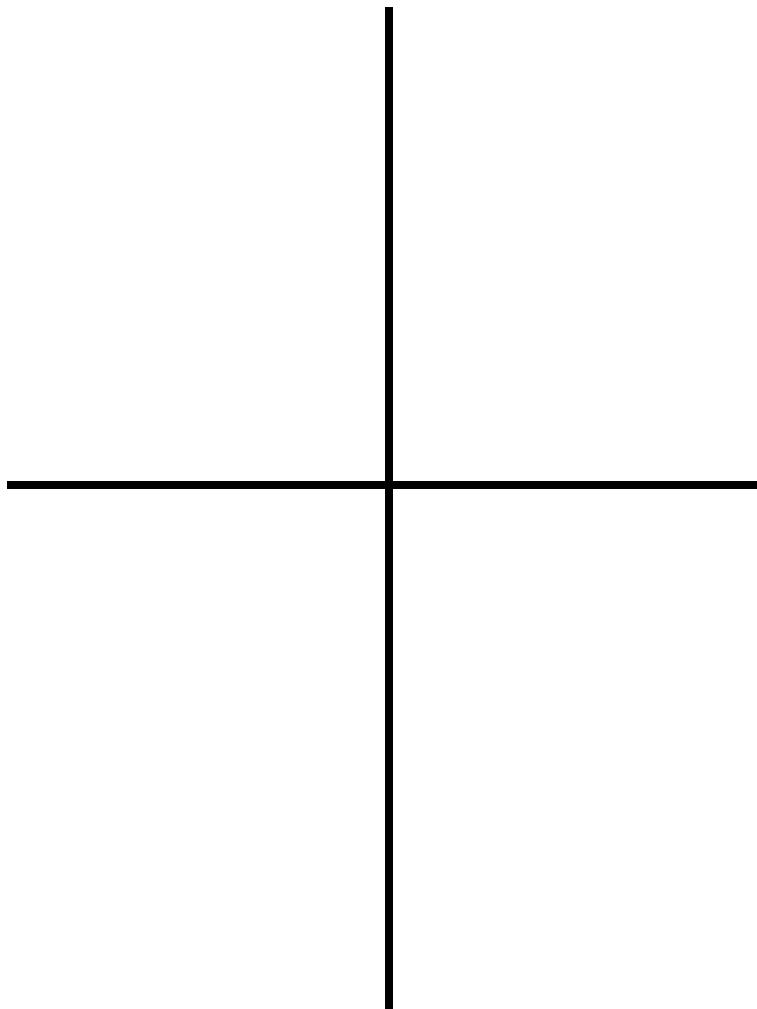


Contour Plots

Surface Plot of $2x^2 + 0.5y^2$



Contour Plot



Exponential Moving Average (EMA)

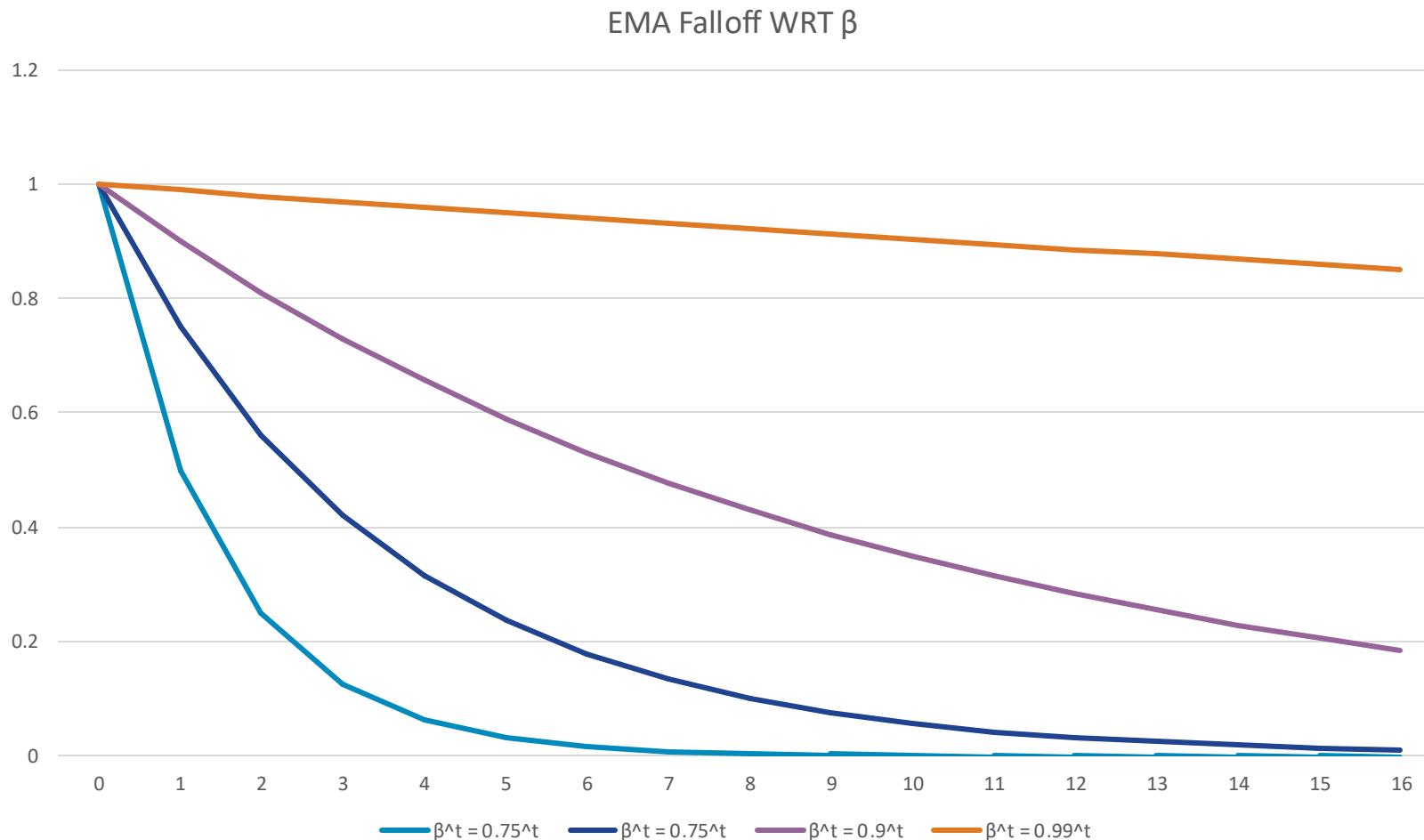
- A simple moving average is the unweighted mean of the previous k values
- EMA is a rule-of-thumb technique for smoothing time series data

Exponential Moving Average (EMA)

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

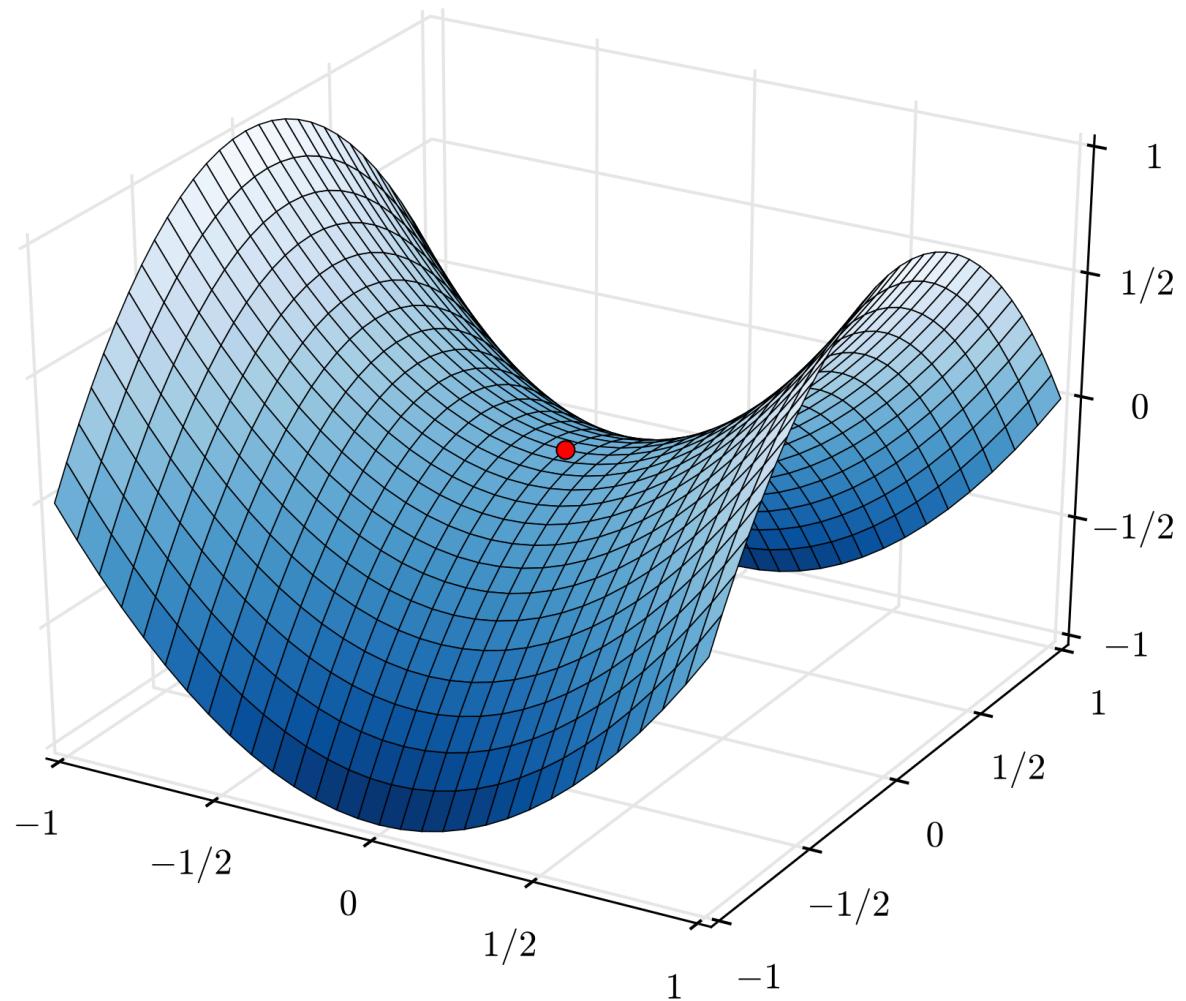
Exponential Moving Average (EMA)

- 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

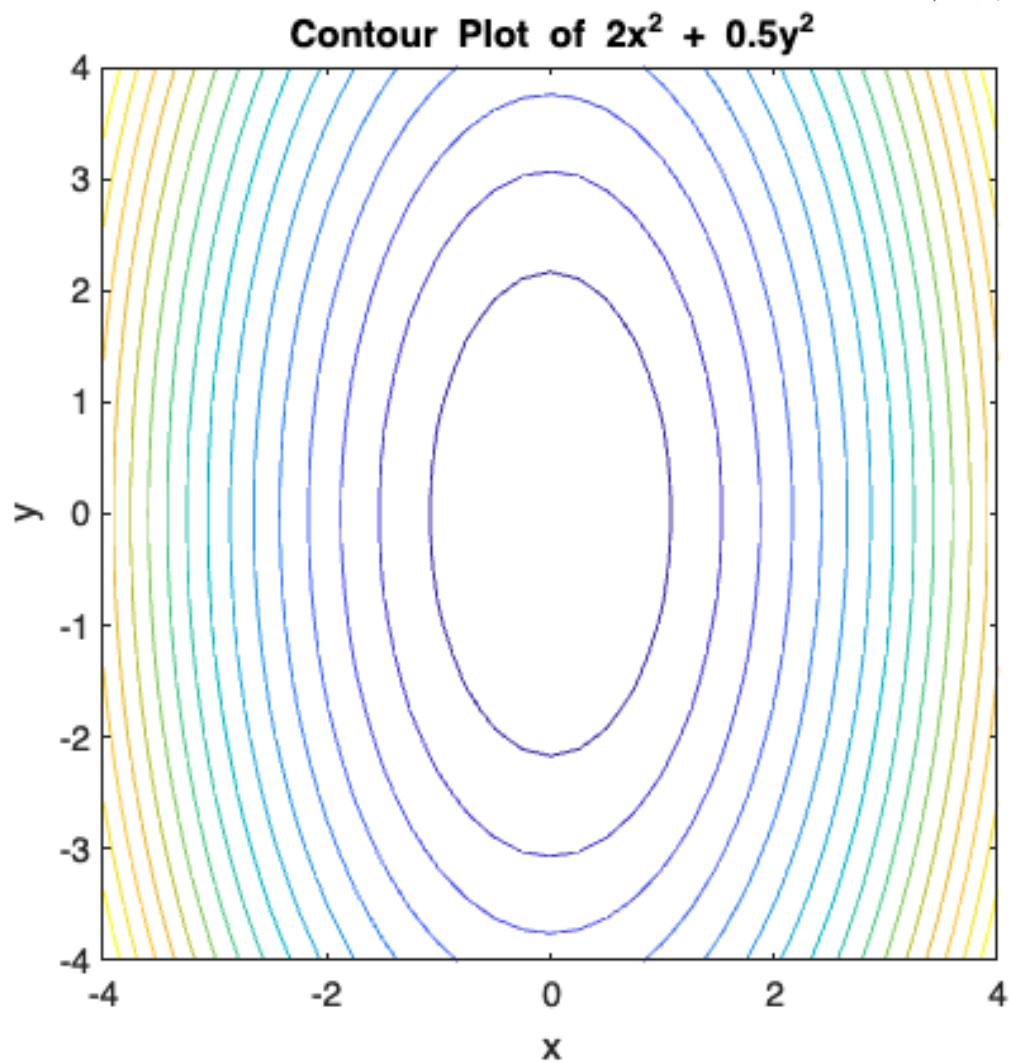
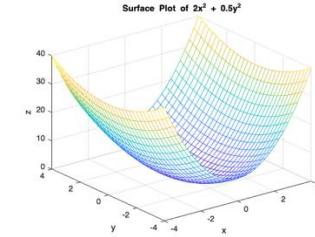


Gradient Plateau

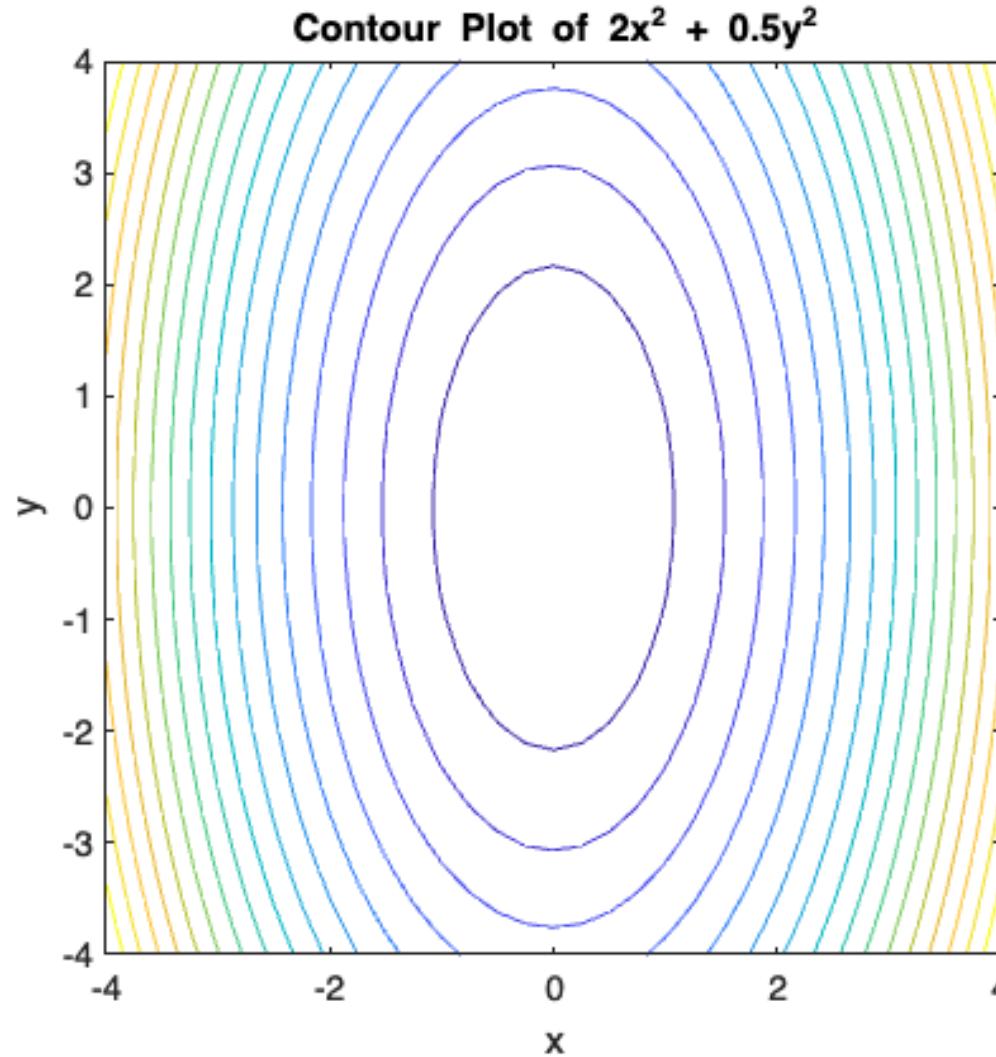
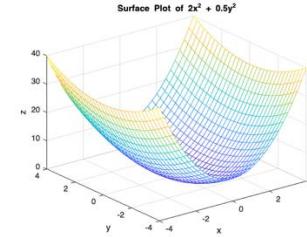
Saddle Points



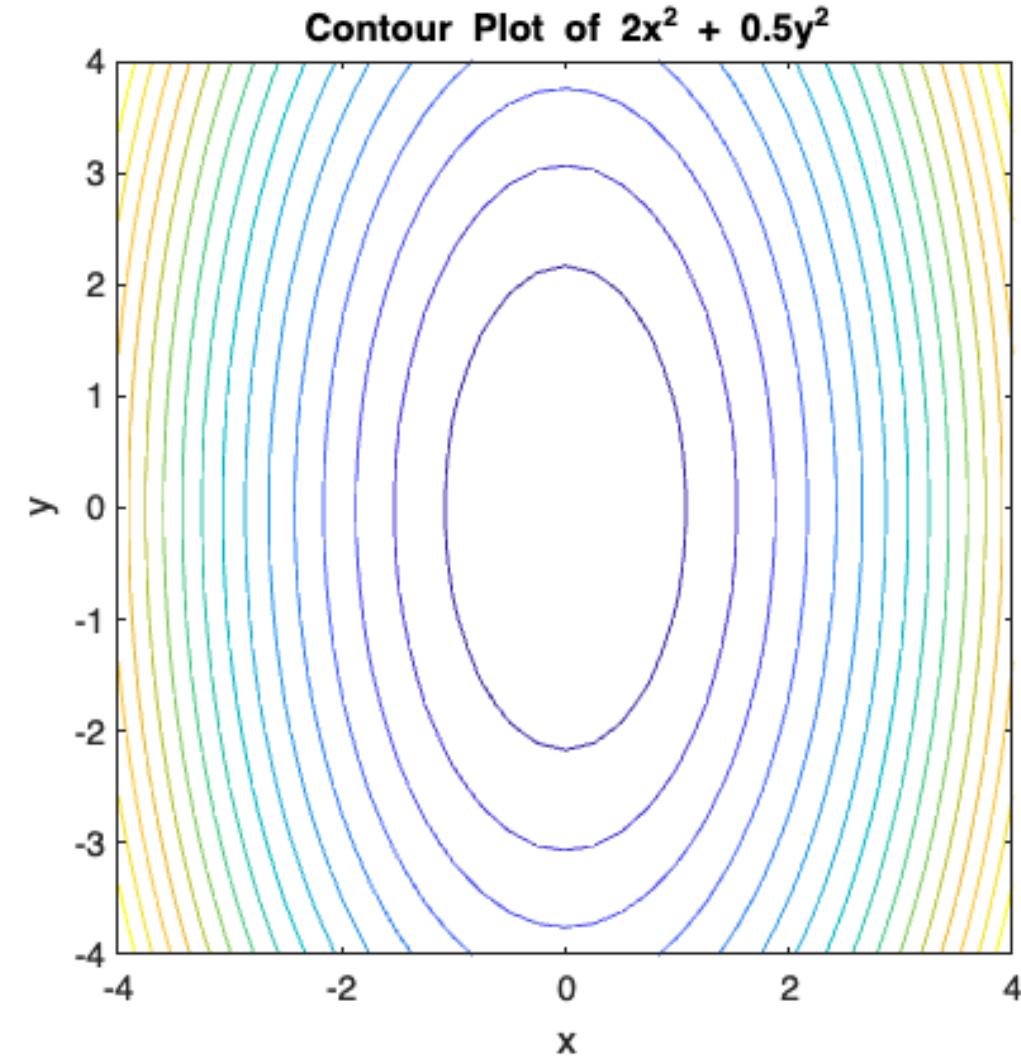
Momentum



Momentum



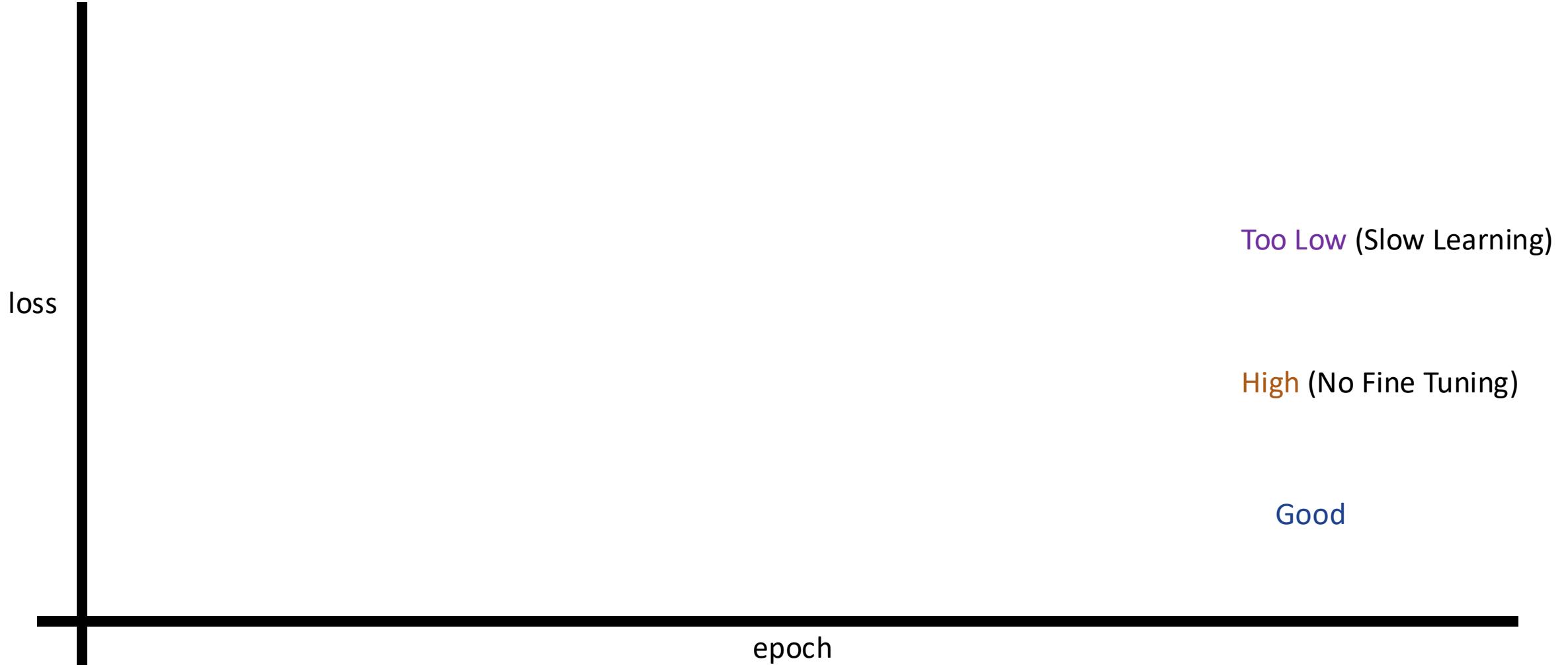
Without Momentum



With Momentum

Learning Rates

Too High (No Learning)



Too Low (Slow Learning)

High (No Fine Tuning)

Good

AdaGrad

RMSProp

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