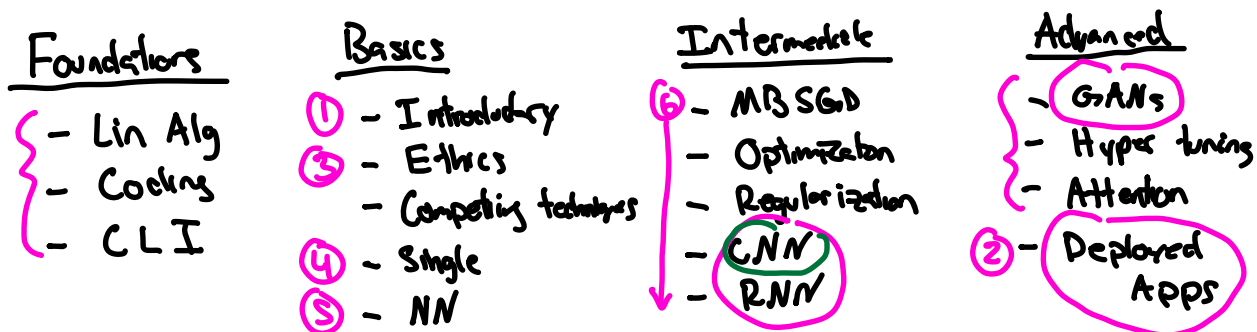


Optimization

$$\nabla \mathcal{L}(\hat{\psi}, \psi)$$

- Timeline
- Optimization
 - Momentum
 - Adaptive Learning Rates (Adagrad, RMSprop, Adam)

Timeline



Don't wait until we cover something in class to use it.

Mini-Batch SGD

learned parameters

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

Matrix of all parameters

Learning rate

loss function

$$\frac{\nabla \mathcal{L}(\hat{y}_b, y_b)}{\nabla \Theta_t}$$

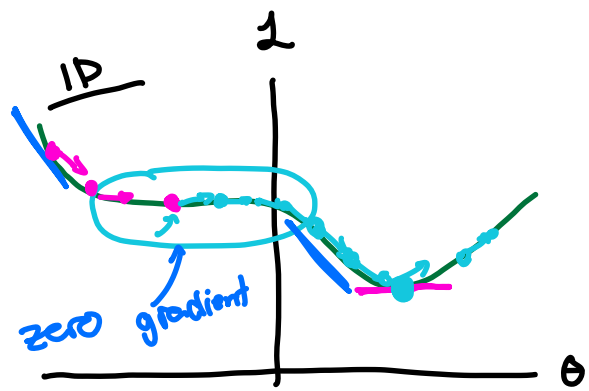
$\nabla_{\Theta} \mathcal{L}_b$

gradient of loss w.r.t. each parameter

- How do we pick η
- Should it be the same for all parameters.
- How do we escape saddle points?

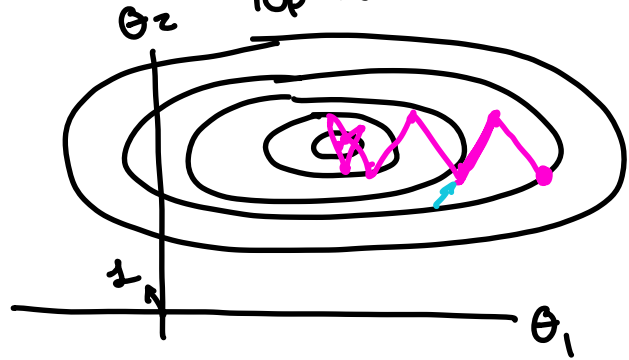


Saddle Points

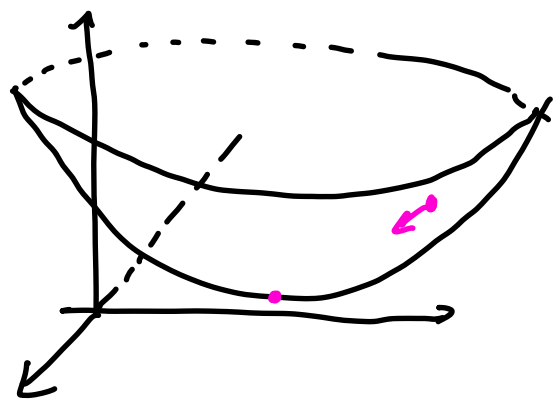


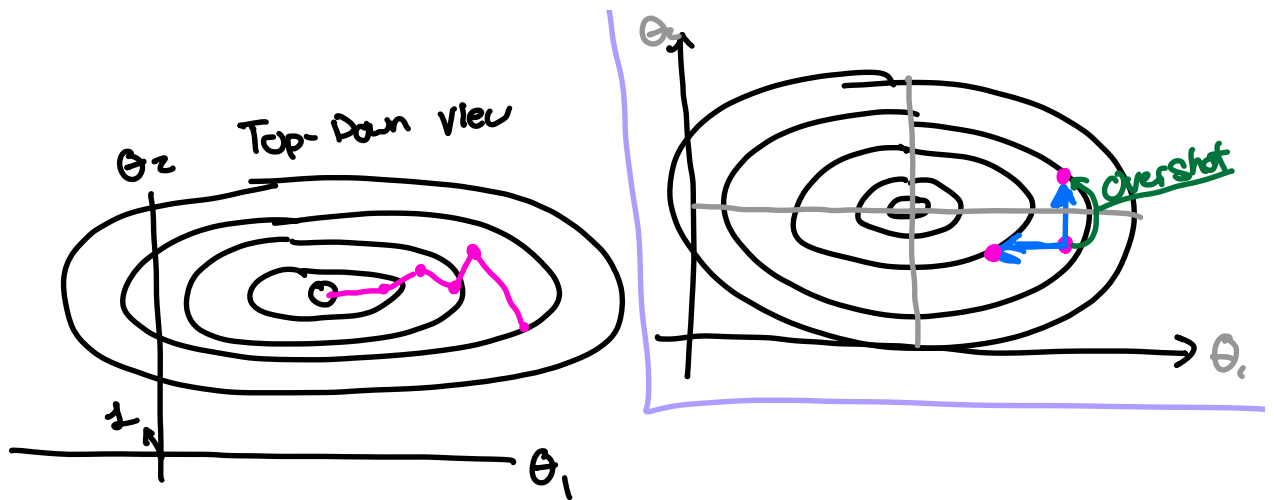
Momentum

Top-Down view



3D-view





"Exponential Moving Average" of ∇

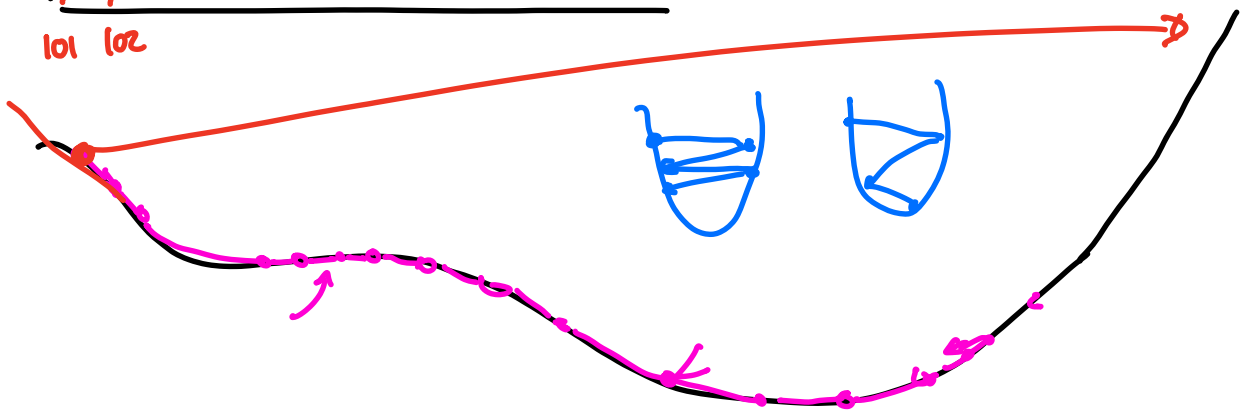
$$\text{averaging} \rightarrow V_{t+1} := \underbrace{\beta}_\text{friction} \underbrace{V_t}_\text{velocity} + (1 - \beta) \underbrace{\nabla_{\theta} L}_\text{acceleration}$$

$$\Theta_{t+1} := \Theta_t - \underbrace{\eta}_\text{rolling average} \underbrace{V_{t+1}}_\text{average gradients}$$

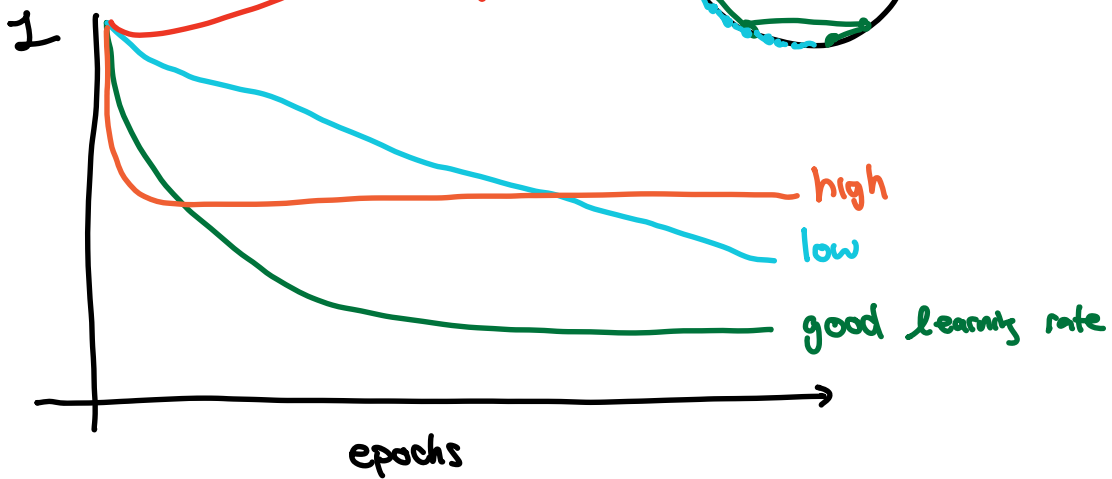
100 time step average



101 102



Adaptive Learning Rates



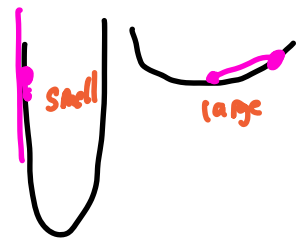
RMSProp "Root-Mean-Square Backprop"

Average magnitude of gradient.

$$g_{t+1}^2 := \beta_R g_t^2 + (1 - \beta_R) \nabla_{\theta} L_b^2$$

$$\Theta_{t+1} := \Theta_t - \underbrace{\frac{\nabla_{\theta} L_b}{\sqrt{g_{t+1}^2} + \epsilon}}_{\text{Effective Learning Rate}}$$

↑ very small value



- What if g^2 is large? → smaller steps
- What if g^2 is small? → larger steps

Adam "Adaptive Moment Estimation"

- A good first choice
- Combines momentum w/ RMSProp

$$v_{t+1} = \beta_m v_t + (1 - \beta_m) \nabla_{\theta} \mathcal{L}_b$$

$$\hat{v}_{t+1} = \frac{v_{t+1}}{1 - \beta_m^t} \quad 0.9^t \quad \frac{1}{1-0.9} = \frac{1}{0.1}, \quad \frac{1}{1-0.9^{100}} = 1$$

$$g_{t+1}^2 = \beta_r g_t^2 + (1 - \beta_r) \nabla_{\theta} \mathcal{L}_b^2$$

$$\hat{g}_{t+1}^2 = \frac{g_{t+1}^2}{1 - \beta_r^t}$$

$$\Theta_{t+1} := \Theta_t - \eta \frac{\hat{v}_{t+1}}{\sqrt{\hat{g}_{t+1}^2 + \epsilon}}$$

← momentum

← adaptive learning rate

Bias
Correction
for
EMA