

Computational Intelligence

Samaneh Hosseini

Isfahan University of Technology

Outline

- Optimization Algorithms
 - Understanding Exponentially Weighted Averages
 - Bias Correction in Exponentially Weighted Averages
 - Gradient Descent with Momentum
 - RMSprop
 - Adam Optimization Algorithm
 - Learning Rate Decay
 - The Problem of Local Optima
- Hyperparameter Tuning
 - Tuning Process
 - Using an Appropriate Scale to pick Hyperparameters
 - Hyperparameters Tuning in Practice

Optimization Algorithms:

Adam Optimization Algorithm

Adam optimization algorithm

$$v_{dw}=0, \quad s_{dw}=0, \quad v_{db}=0, \quad s_{db}=0$$

on iteration t :

compute dw, db using current mini-batch

$$v_{dw} = \beta_1 v_{dw} + (1 - \beta_1) dw \quad v_{db} = \beta_1 v_{db} + (1 - \beta_1) db \leftarrow \text{momentum } \beta_1$$

$$s_{dw} = \beta_2 s_{dw} + (1 - \beta_2) dw^2 \quad s_{db} = \beta_2 s_{db} + (1 - \beta_2) db^2 \leftarrow \text{RMS}_{\text{prob}}$$

$$v_{dw}^{\text{corrected}} = v_{dw} / (1 - \beta_1^t), \quad v_{db}^{\text{corrected}} = v_{db} / (1 - \beta_1^t)$$

$$s_{dw}^{\text{corrected}} = s_{dw} / (1 - \beta_2^t), \quad s_{db}^{\text{corrected}} = s_{db} / (1 - \beta_2^t)$$

$$w := w - \alpha \frac{v_{dw}^{\text{corrected}}}{\sqrt{s_{dw}^{\text{corrected}} + \epsilon}}$$

$$b := b - \alpha \frac{v_{db}^{\text{corrected}}}{\sqrt{s_{db}^{\text{corrected}} + \epsilon}}$$

Hyperparameters Choice

→ α : learning rate needs to be tune ←

β_1 : 0.9 → (\underline{dw})

β_2 : 0.999 → $(\underline{\frac{dw^2}{\alpha}})$

ϵ : 10^{-8}

Adam : Adaptive momentum estimation

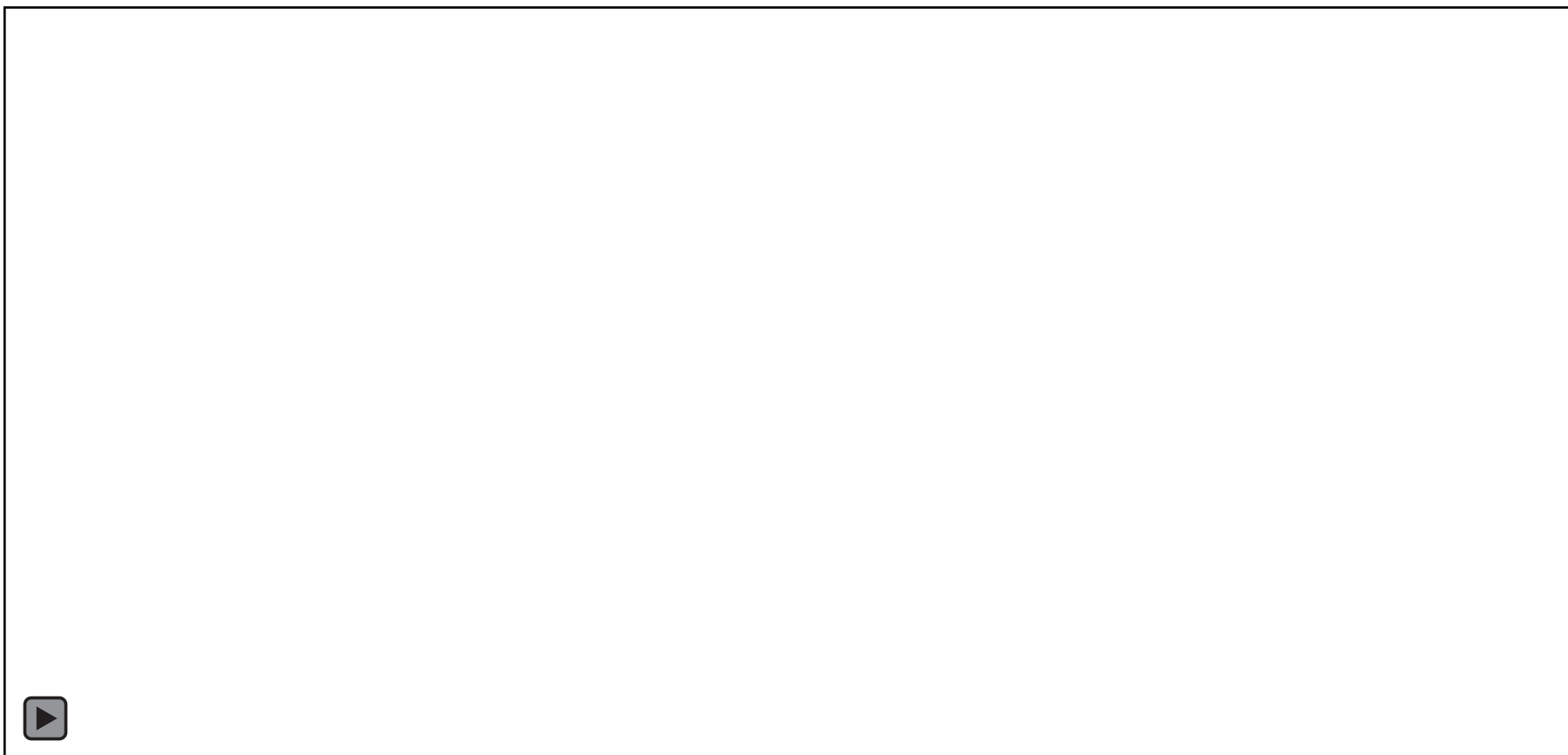
Setting the Learning Rate

Small learning rate converges slowly and gets stuck in false local minima



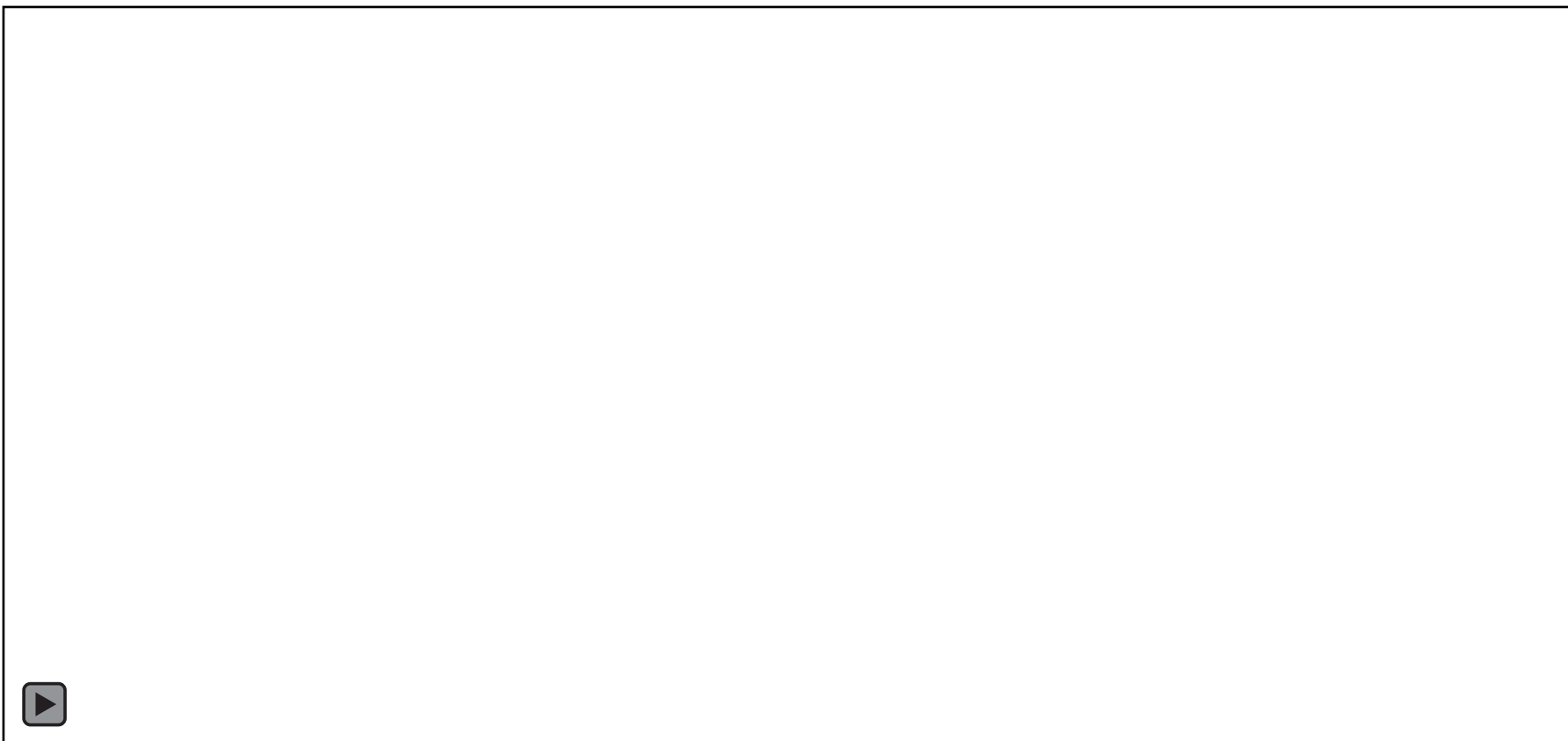
Setting the Learning Rate

Large learning rates overshoot, become unstable and diverge



Setting the Learning Rate

Stable learning rates converge smoothly and avoid local minima



How to deal with this?

Idea 1:

Try lots of different learning rates and see what works “just right”

Adaptive Learning Rates

Idea 1:

Try lots of different learning rates and see what works “just right”

Idea 2:

Do something smarter!

Design an adaptive learning rate that “adapts” to the landscape

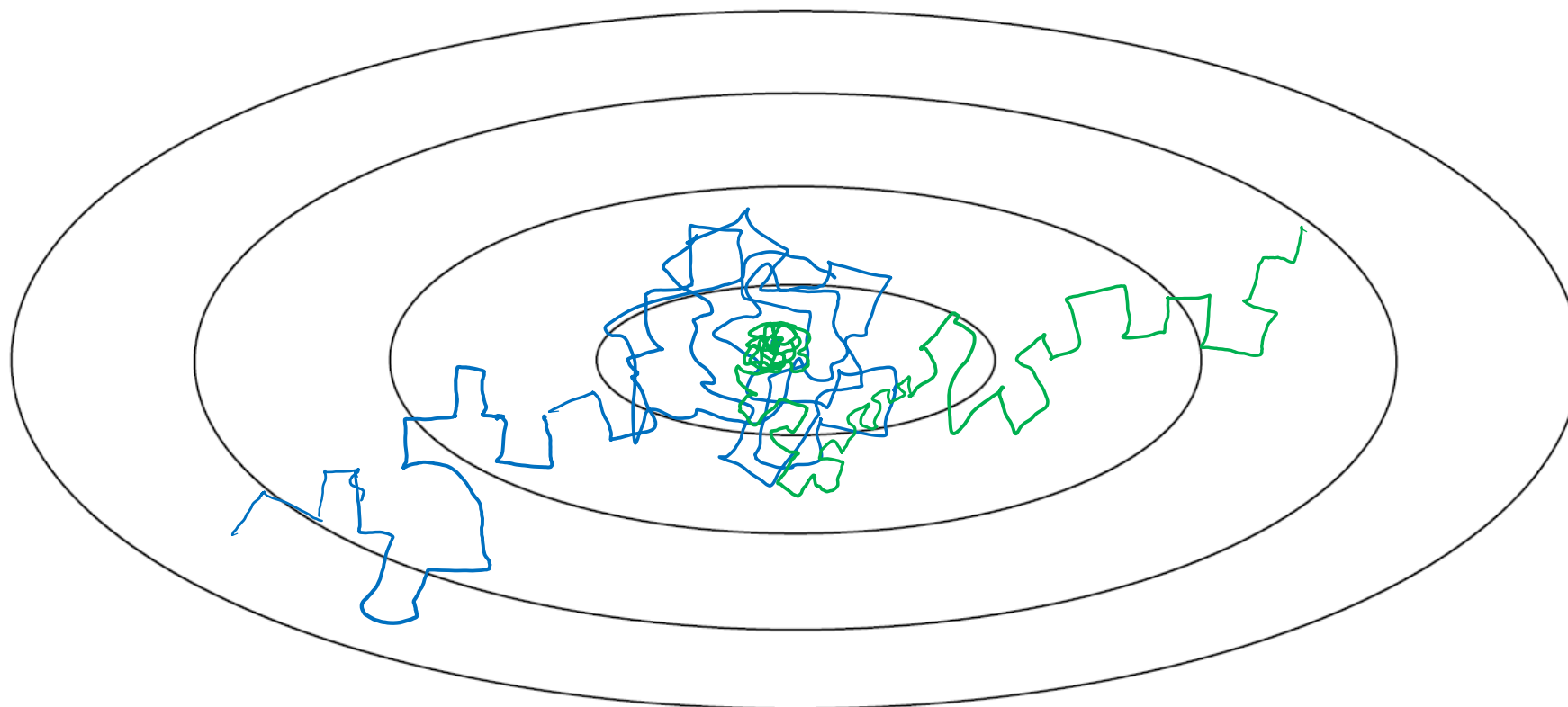
Adaptive Learning Rates

- Learning rates are no longer fixed
- Can be made larger or smaller depending on:
 - how large gradient is
 - how fast learning is happening
 - size of particular weights
 - etc...

Optimization Algorithms: Learning Rate Decay

Learning Rate Decay

slowly reduce α

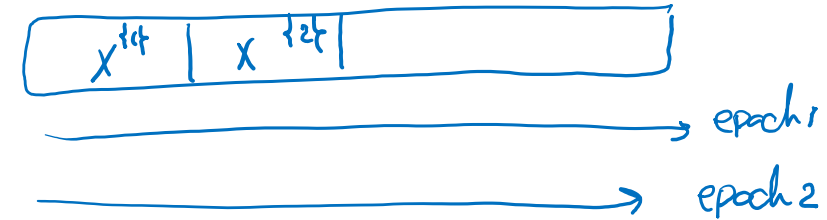


Learning Rate Decay

1 epoch: 1 pass through data

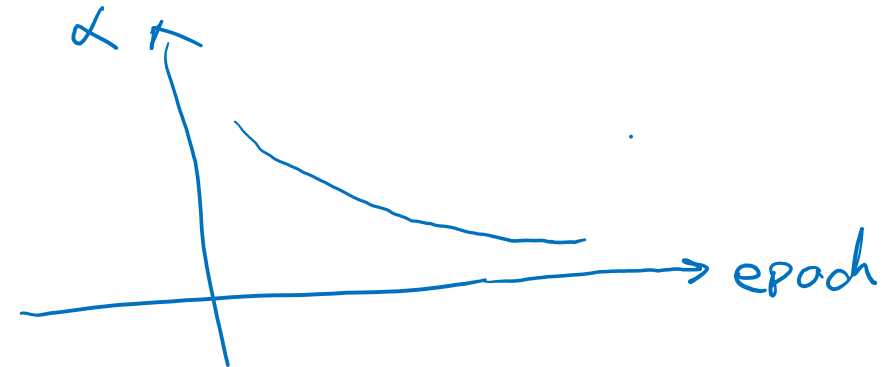
$$\alpha = \frac{\alpha_0}{1 + \text{decay-rate} * \text{epoch-num}}$$

Epoch	α
1	0.1
2	0.067
3	0.05
4	0.04
⋮	⋮



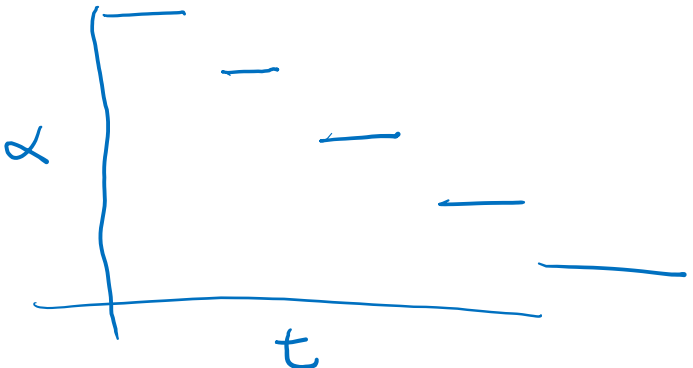
$$\alpha_0 = 0.2$$

$$\text{decay-rate} = 1$$



Other Learning Rate Decay Methods

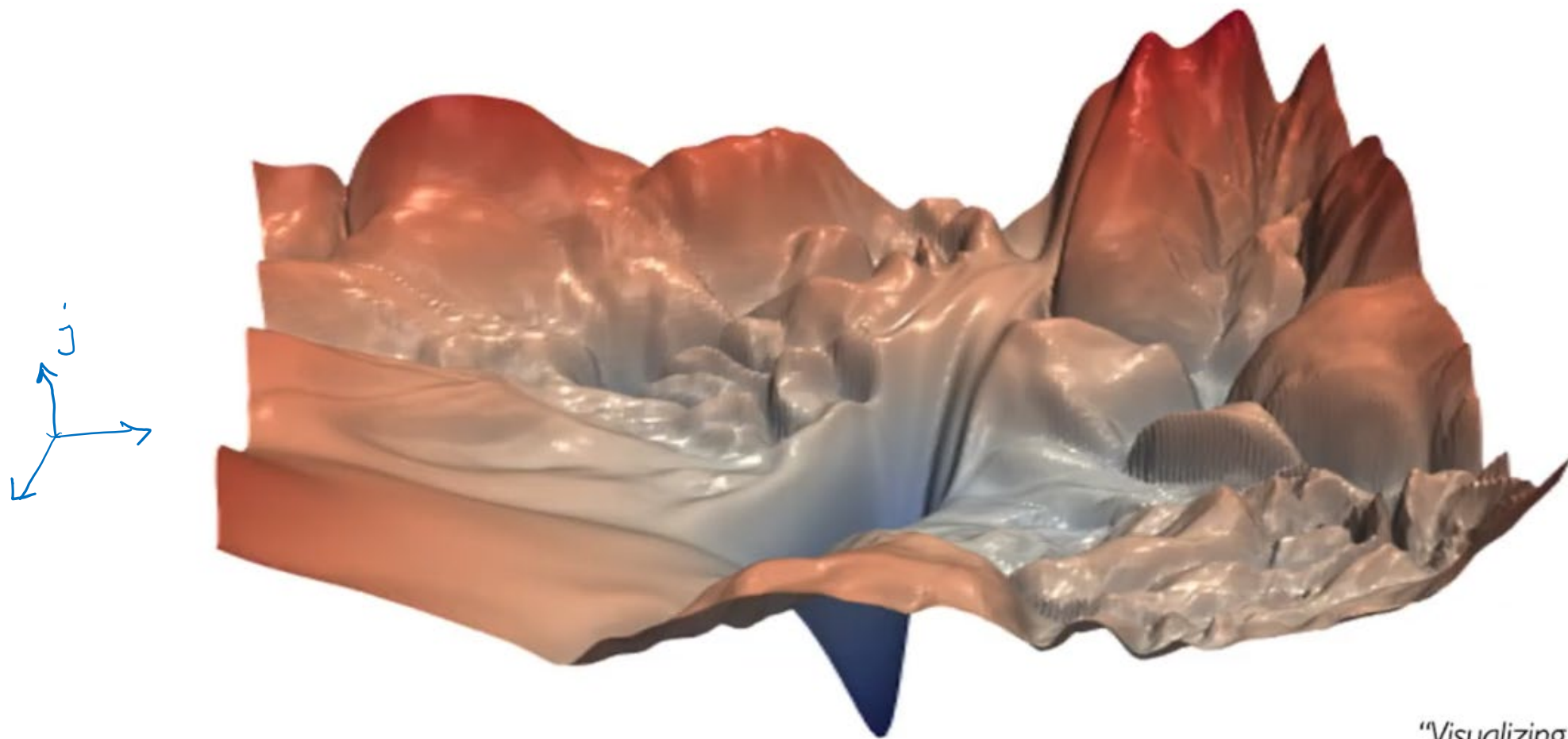
Formula {

$$\alpha = 0.95^{\text{epoch-num}} \alpha_0$$
$$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \alpha_0$$


Manual decay

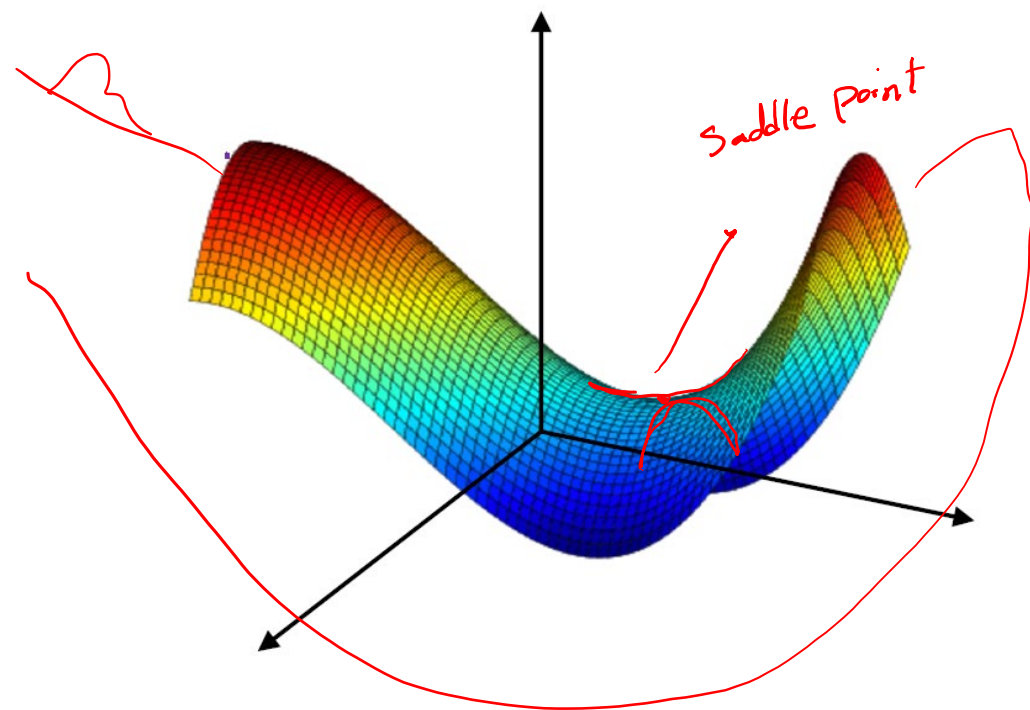
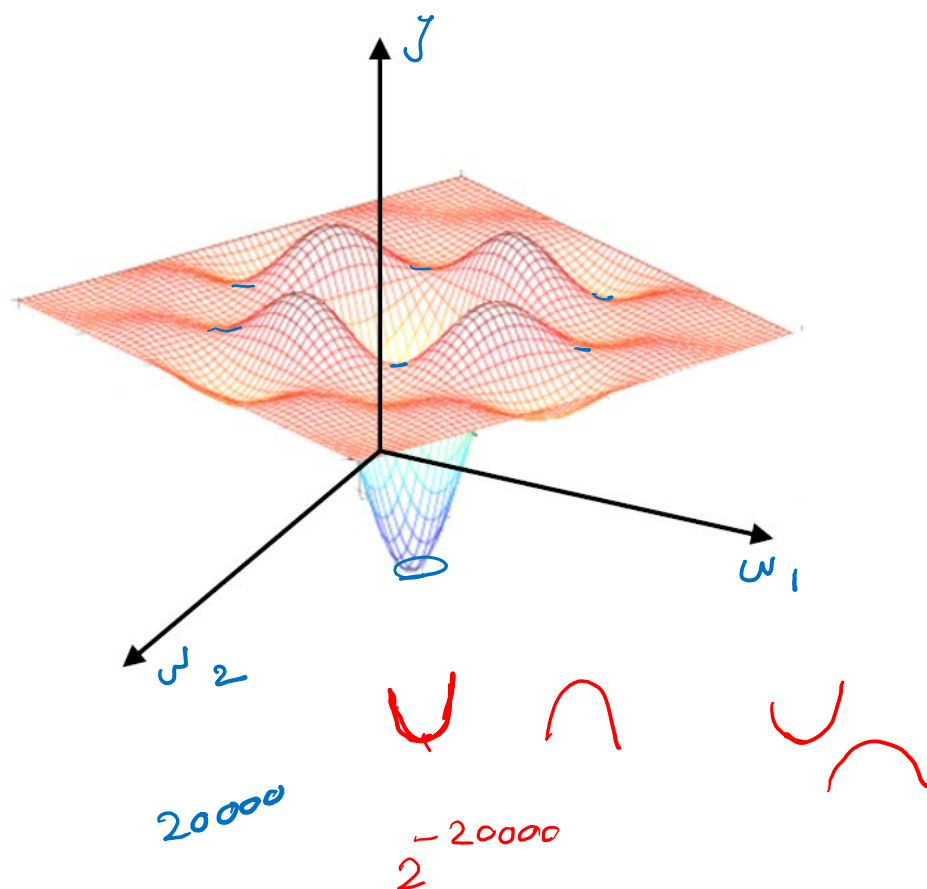
Optimization Algorithms: The Problem of Local Optima

Training Neural Network is Difficult

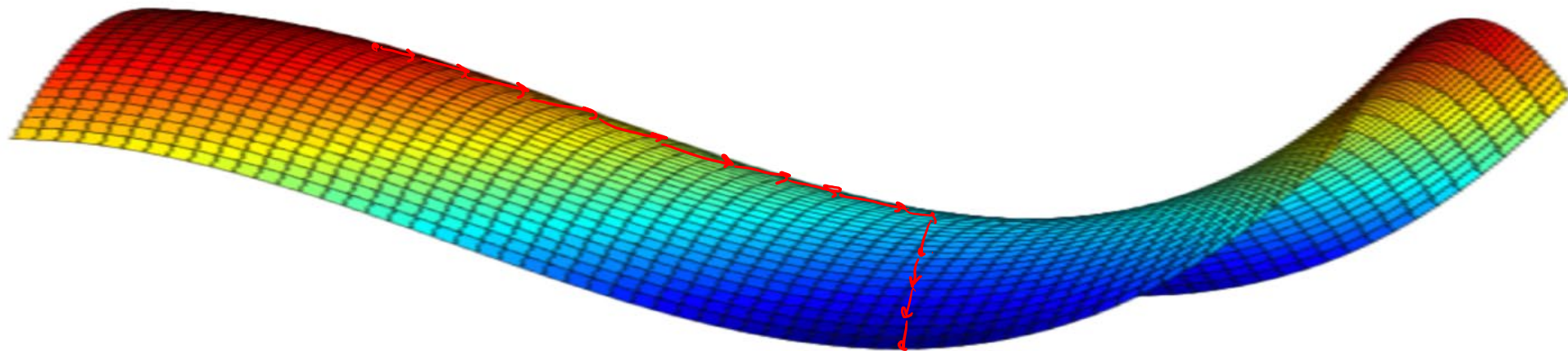


"Visualizing the loss landscape of neural nets". Dec 2017.

Local Optima in Neural Networks



Problem of plateaus



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow

Hyperparameter Tuning: Tuning Process

Hyperparameters

- Hyperparameter Tuning

α

$\beta \sim 0.9$

$\beta_1, \beta_2, \epsilon$
 $0.9 \quad 0.999 \quad 10^{-8}$

layers

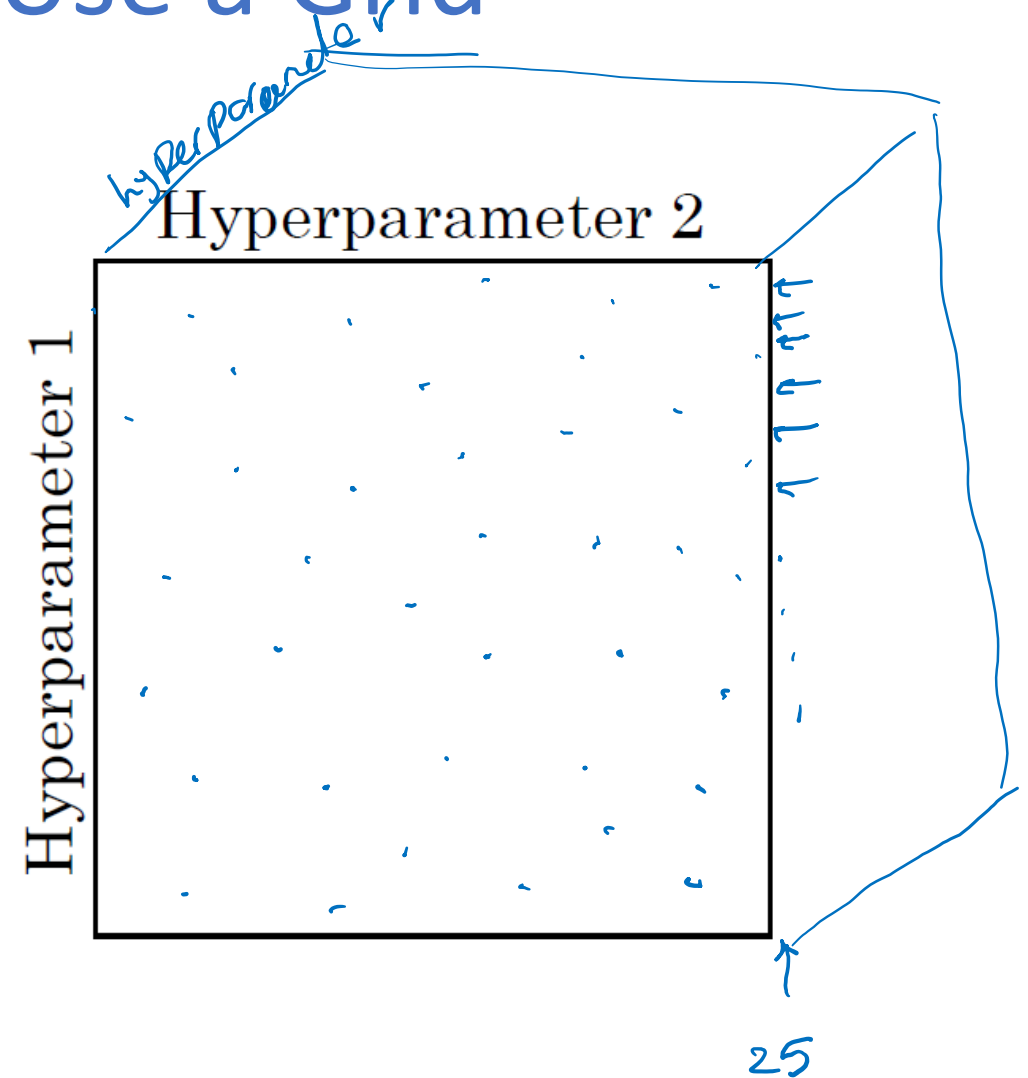
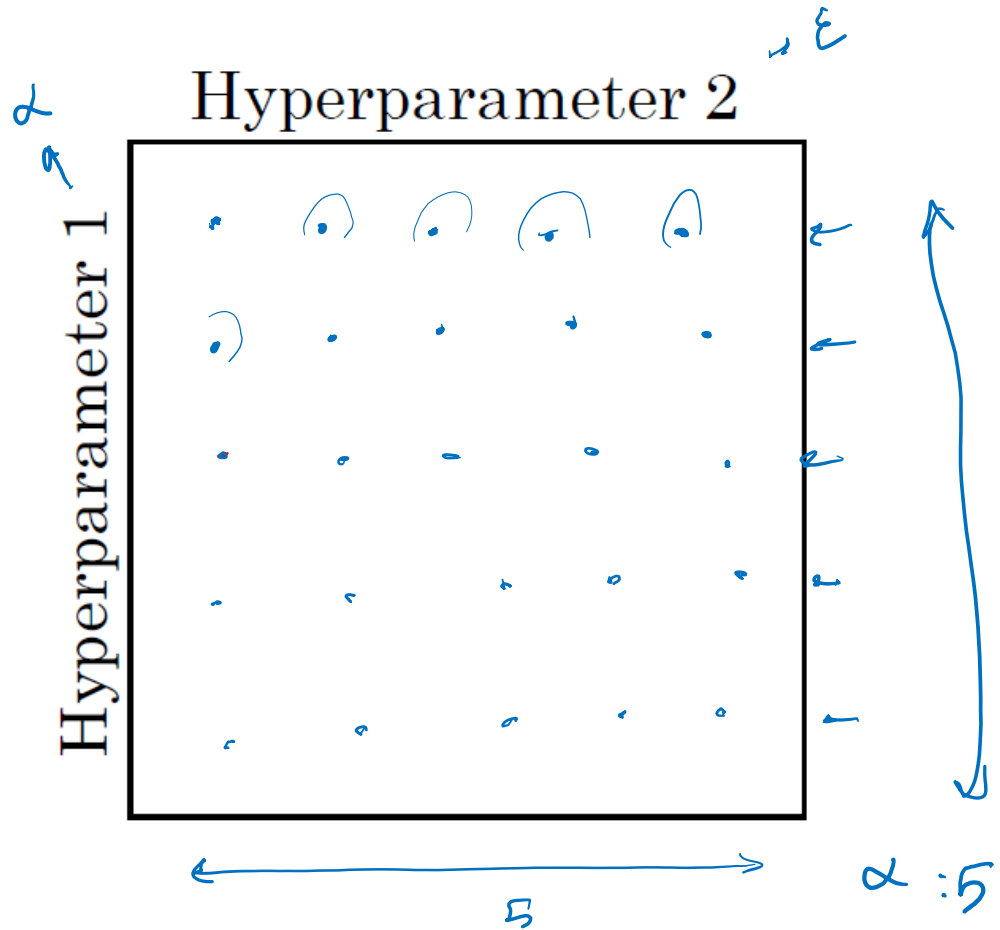
hidden units

learning rate decay

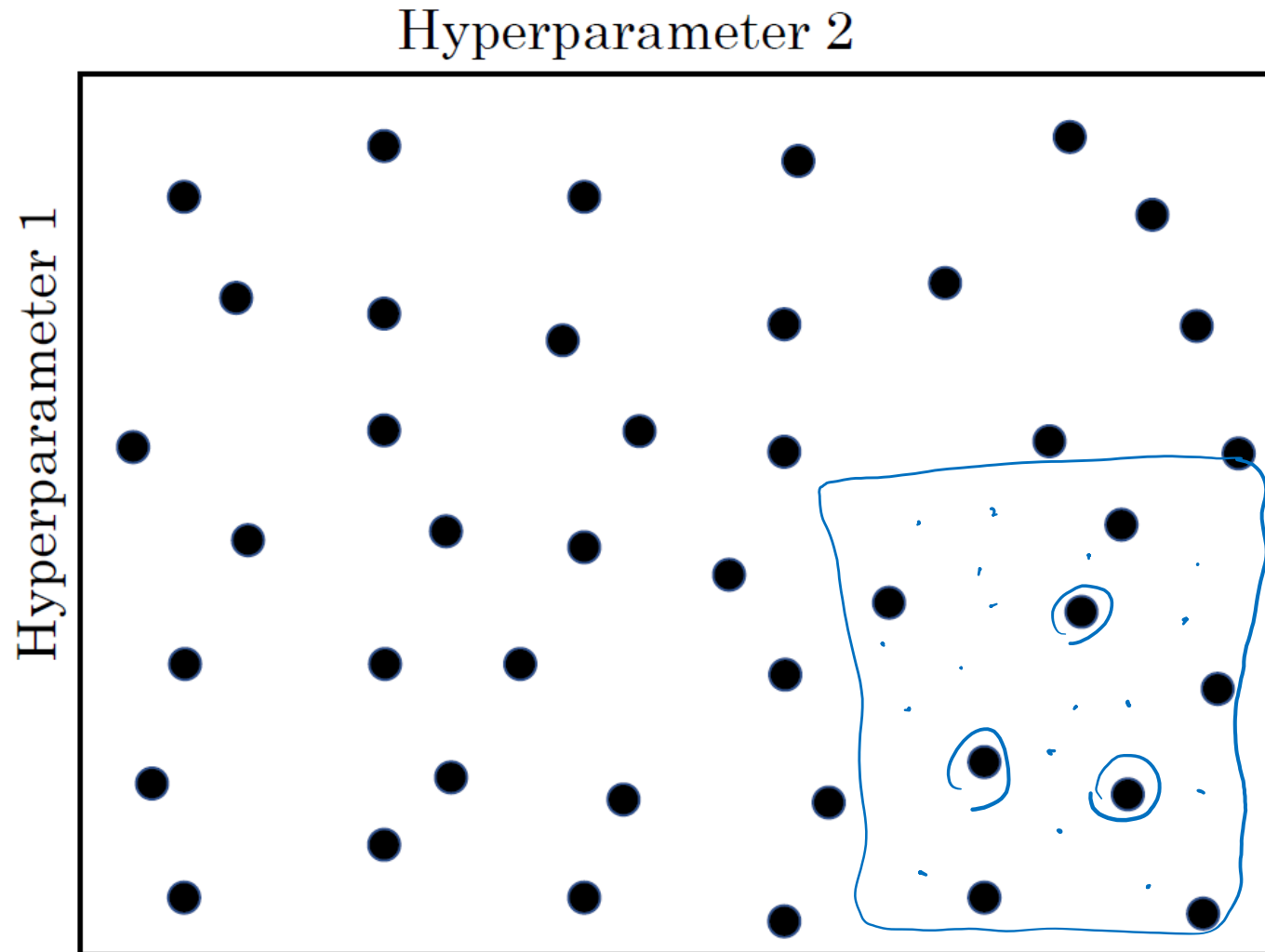
mini-batch size

/

Try Random Values: Don't Use a Grid



Coarse to fine



Hyperparameter Tuning:

Using an appropriate scale to pick hyperparameters

Picking Hyperparameters at Random

$$n^{[2]} = 50 \dots 100$$

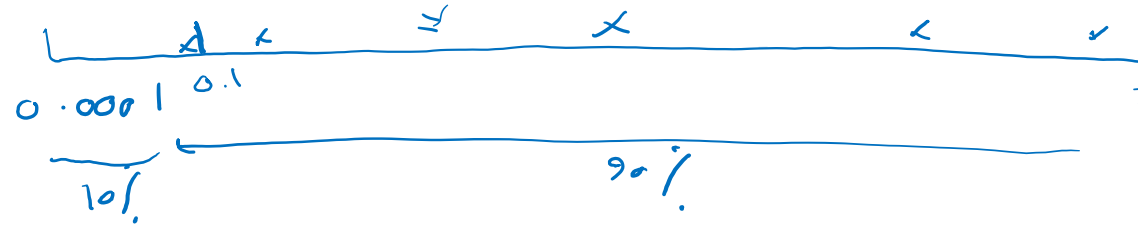


#layer $L: 2 - 4$

2, 3, 4

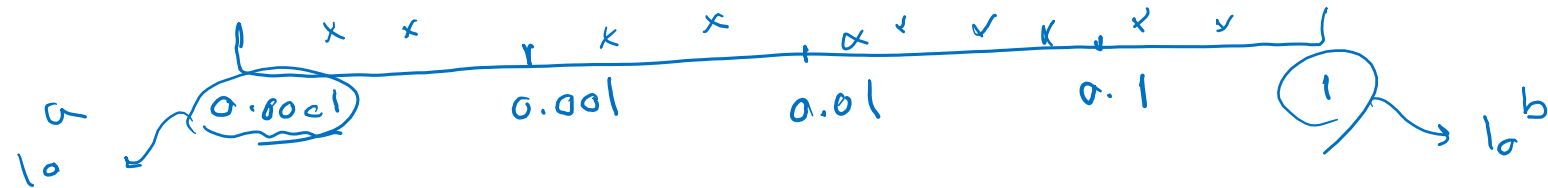
Appropriate Scale for Hyperparameters

$$\alpha = 0.0001 \quad \dots \quad 1$$



$$\alpha = \log_{10} 0.0001$$

$$= -4$$



$$\log_{10} b = 0$$

$$\log_{10} 1 = 0$$

$$r = -4 * np.random.rand() \leftarrow [-4, 0]$$

$$\leftarrow 10^{-4} \dots 10^0$$

$$0.0001 \dots 1$$

$$\alpha = 10^r$$

$$\alpha \in 10^a \dots 10^b$$

$$r \in [a, b]$$

Hyperparameters for Exponentially Weighted Averages

$$\beta = 0.9 \quad \dots \quad 0.999$$

\downarrow \downarrow
 10 1000

$$1-\beta = 0.1 \quad \dots \quad 0.001$$

$$\beta = 0.9000 \rightarrow 0.9005 \quad \sim 10$$

$$\beta = 0.999 \rightarrow 0.9995 \quad \sim 2000$$

~ 1000

$$\left[\begin{array}{cccccc} \times & \times & \times & \times & \times & \times \end{array} \right]$$

$0.9 \quad \quad \quad 0.999$

$$\left[\begin{array}{ccc} & \cdot & \end{array} \right]$$

$0.9 \quad \quad 0.99 \quad \quad 0.999$

$$\left[\begin{array}{ccc} \underbrace{}_{10^{-1}} & \underbrace{}_{10^{-2}} & \underbrace{}_{10^{-3}} \end{array} \right]$$

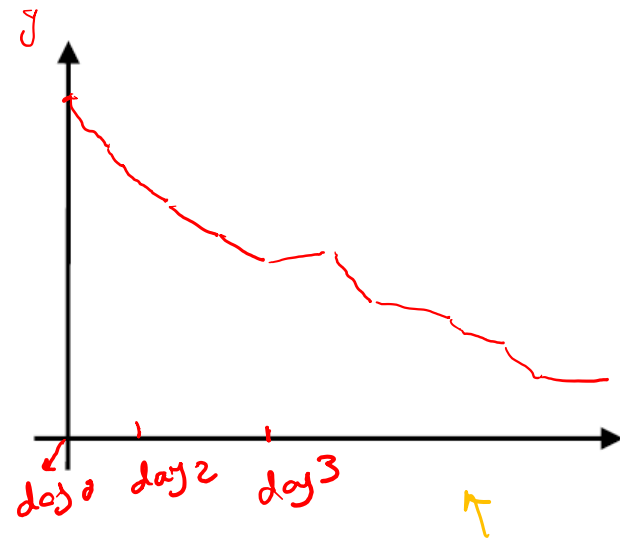
$$r = [-3, -1]$$

$$1-\beta = 10^r$$

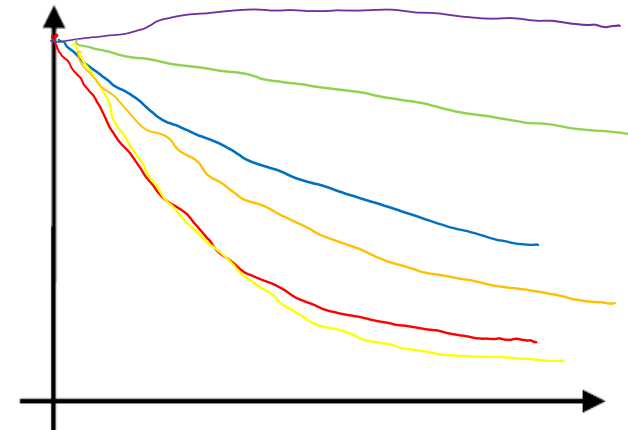
$$\beta = 1 - 10^r$$

Hyperparameters tuning in practice

Babysitting one model



Training many models in parallel



Outline

- Optimization Algorithms
 - Understanding Exponentially Weighted Averages
 - Bias Correction in Exponentially Weighted Averages
 - Gradient Descent with Momentum
 - RMSprop
 - • Adam Optimization Algorithm
 - • Learning Rate Decay
 - • The Problem of Local Optima
- Hyperparameter Tuning
 - Tuning Process
 - Using an Appropriate Scale to pick Hyperparameters
 - Hyperparameters Tuning in Practice