

Copyright Notice

These slides are distributed under the Creative Commons License.

DeepLearning.AI makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite DeepLearning.AI as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>

Optimization Algorithms

Mini-batch
gradient descent

deeplearning.ai



Batch vs. mini-batch gradient descent

卷之三

Vectorization allows you to efficiently compute on m examples.

$$X = \underbrace{[\dots [}_{(m \times n)} \underbrace{\dots [}_{(k \times m)} \underbrace{\dots [}_{(l \times k)} \underbrace{\dots [}_{(n \times l)} \underbrace{\dots [}_{(1 \times n)} \underbrace{\dots [}_{(m \times 1)} X$$

$$X = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

$$m = \sum_{i=1}^n m_i$$

e.g.: X_{t+3}, Y_{t+3} .
batch size = 5, 500, 500?
mini-batches of 1,000 each

Mini-batch t : $X^{\{t\}}, y^{\{t\}}$

Andrew Ng

Mini-batch gradient descent

repeat for $t = 1, \dots, 5000$

Forward prop on $X^{\{t\}}$.

$$\begin{aligned} z^{(t)} &= \theta^{(t)} X^{(t)} + b^{(t)} \\ A^{(t)} &= g^{(t)}(z^{(t)}) \end{aligned}$$

Compute cost $J^{(t)}$ = $\frac{1}{1500} \sum_{i=1}^{1500} l(y^{(t)}, y_i^{(t)})$
 Backprop to compute gradients w.r.t $\theta^{(t)}$ (using $(X^{(t)}, Y^{(t)})$)

$$\theta^{(t+1)} := \theta^{(t)} - \alpha \nabla_{\theta} J^{(t)}, \quad \theta^{(t+1)}_j = \theta^{(t)}_j - \alpha \frac{\partial J^{(t)}}{\partial \theta_j}$$

"1 epoch" pass through training set.
 3

1 step of gradient desc.
 using $X^{(t+1)}, Y^{(t+1)}$
 (as $f_{\text{model}}(x)$)

Optimization Algorithms

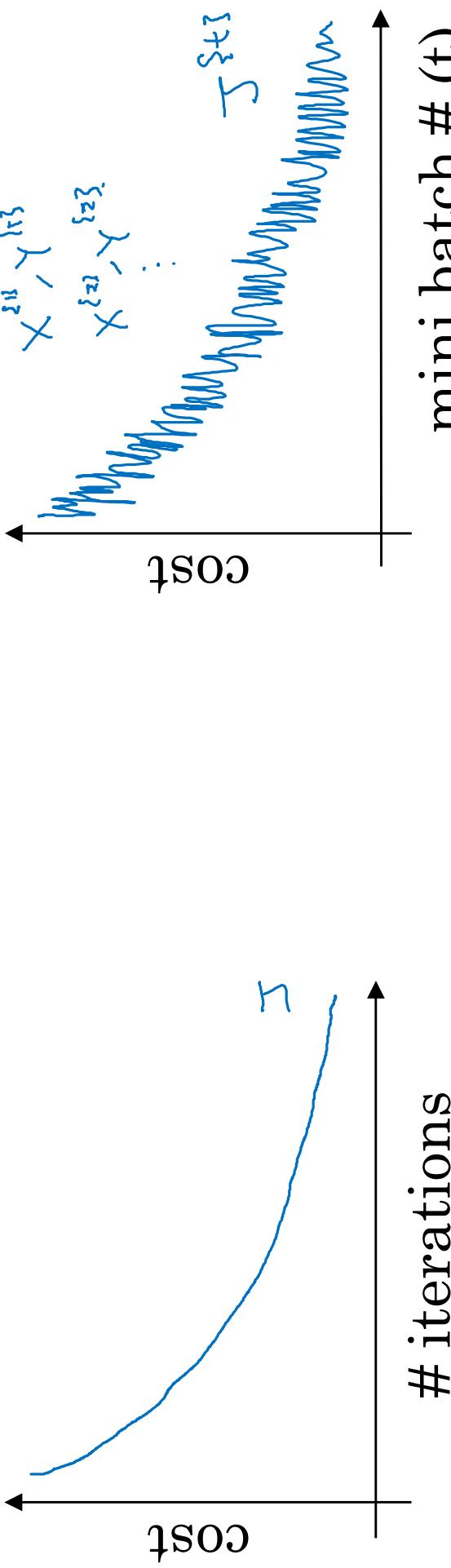
Understanding
mini-batch
gradient descent

deeplearning.ai

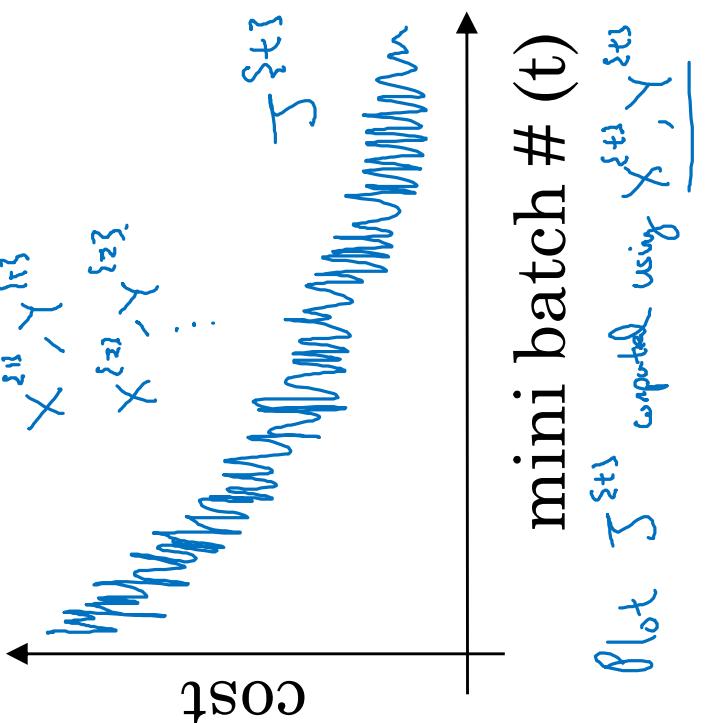


Training with mini batch gradient descent

Batch gradient descent



Mini-batch gradient descent



Andrew Ng

Choosing your mini-batch size

→ If mini-batch size = m : Batch gradient descent. $(X^{(1)}, Y^{(1)}) = (X, Y)$.

→ If mini-batch size = 1 : Stochastic gradient descent. Every example is it own
 $(X^{(1)}, Y^{(1)}) = (x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$... $(X^{(1)}, Y^{(1)})$ minis-batch.

In practice: Some in-between 1 and m

Stochastic
gradient
descent

In-between
(mini-batch size
not too big (small))

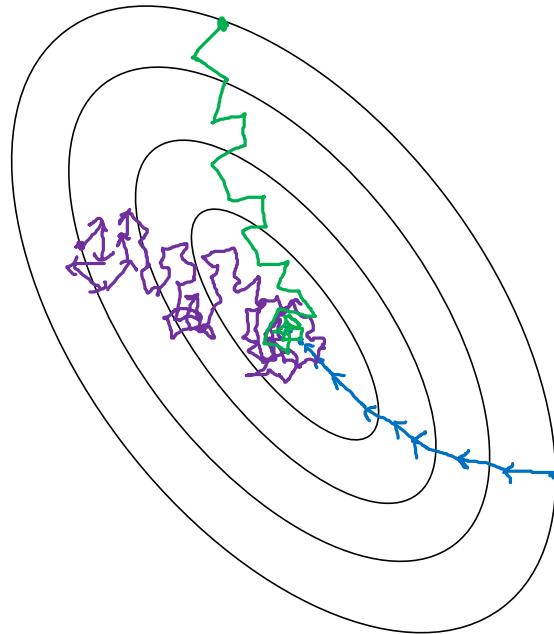
Batch
gradient descent
(mini-batch size = m)

Use sparse
from vectorization

Fastest learning.
• Vectorization.
(n zero)
• Make passes without
processing entire training set.

Two long
per iteration

Andrew Ng



Choosing your mini-batch size

If small tiny set : Use batch gradient descent.
($m \leq 2500$)

Typical mini-batch sizes:

$$\rightarrow \underbrace{64, 128, 256, 512}_{2^6, 2^7, 2^8, 2^9} \quad \frac{1024}{2^{10}}$$

Make sure mini-batch fit in CPU/GPU memory.

$$X^{\{1\}}, Y^{\{1\}}$$

Optimization Algorithms

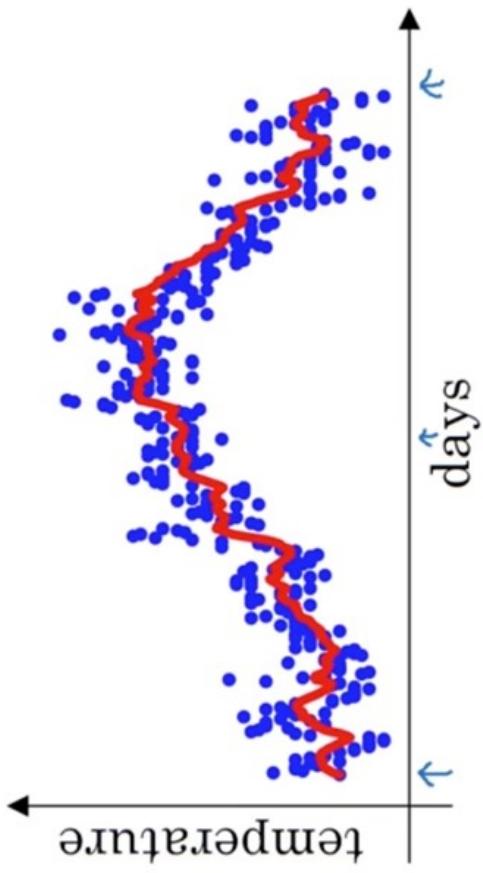
Exponentially
weighted averages

deeplearning.ai



Temperature in London

$$\begin{aligned}\theta_1 &= 40^{\circ}\text{F} \quad 4^{\circ}\text{C} \quad \leftarrow \\ \theta_2 &= 49^{\circ}\text{F} \quad 9^{\circ}\text{C} \\ \theta_3 &= 45^{\circ}\text{F} \quad \vdots \\ &\vdots \\ \theta_{180} &= 60^{\circ}\text{F} \quad 15^{\circ}\text{C} \\ \theta_{181} &= 56^{\circ}\text{F} \quad \vdots\end{aligned}$$



$$V_0 = 0$$

$$\begin{aligned}V_1 &= 0.9 V_0 + 0.1 \Theta_1 \\ V_2 &= 0.9 V_1 + 0.1 \Theta_2 \\ V_3 &= 0.9 V_2 + 0.1 \Theta_3 \\ &\vdots\end{aligned}$$

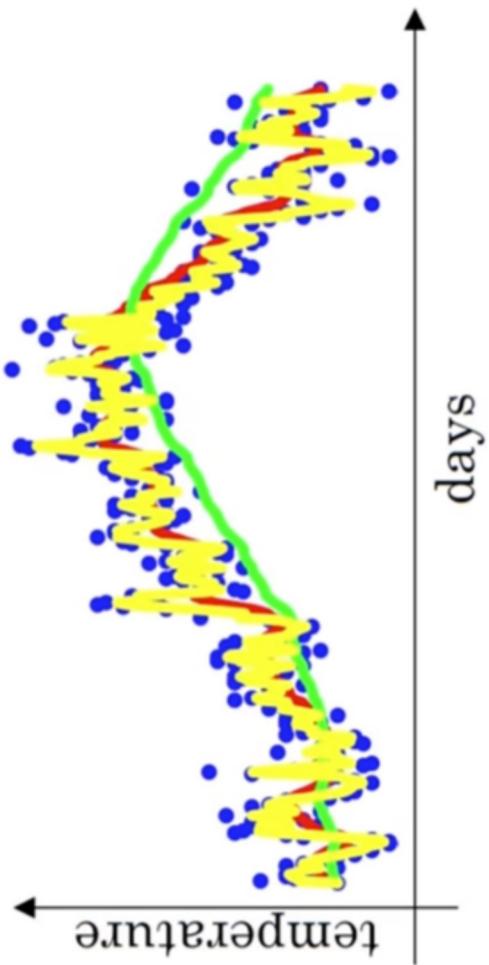
$$V_t = 0.9 V_{t-1} + 0.1 \Theta_t$$

Exponentially weighted averages

$$\bar{V}_t = \frac{\beta}{1-\beta} \frac{V_{t-1}}{1-\beta} + \frac{(1-\beta)}{1-\beta} \theta_t$$

$\beta = 0.9$: ≈ 10 damp' temperature
 $\beta = 0.98$: ≈ 50 damp
 $\beta = 0.5$: ≈ 2 damp

V_t is approximately
over older
 $\Rightarrow V_t \frac{1}{1-\beta}$ damp' temperature.



$$\frac{1}{1-0.98} = 50$$

Optimization Algorithms

Understanding
exponentially
weighted averages

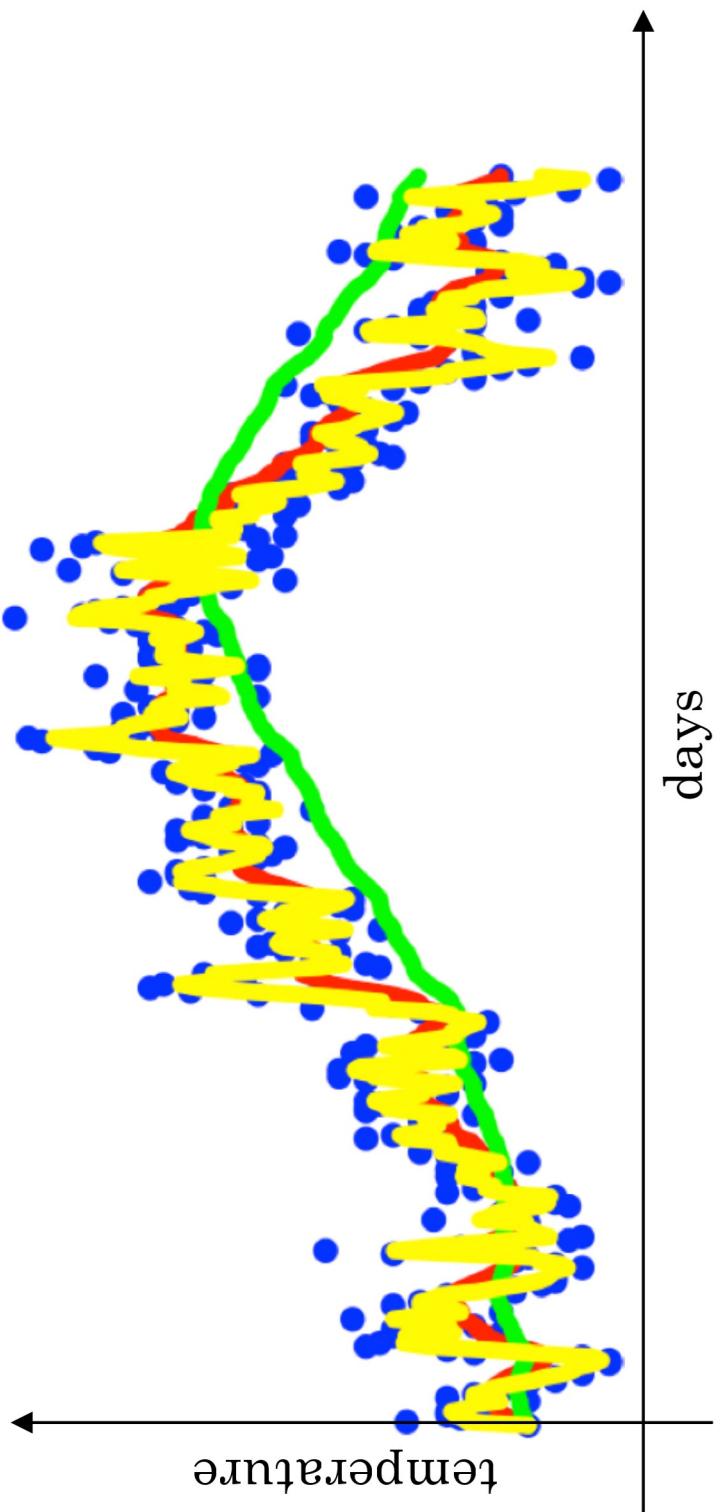
deeplearning.ai



Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

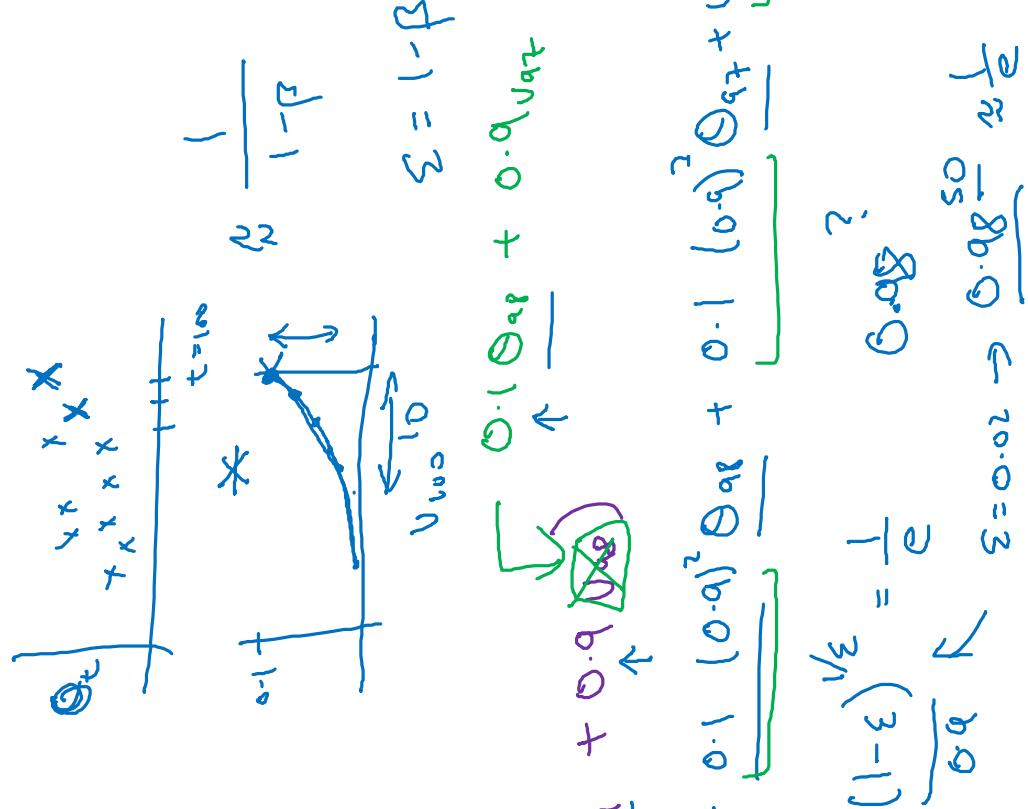
$$\beta = 0.9 \quad 0.98 \quad 0.5$$



Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$\begin{aligned} v_{100} &= 0.9v_{99} + 0.1\theta_{100} \\ v_{99} &= 0.9v_{98} + 0.1\theta_{99} \\ v_{98} &= 0.9v_{97} + 0.1\theta_{98} \end{aligned}$$



Implementing exponentially weighted averages

$$\begin{aligned}v_0 &= 0 \\v_1 &= \beta v_0 + (1 - \beta) \theta_1 \\v_2 &= \beta v_1 + (1 - \beta) \theta_2 \\v_3 &= \beta v_2 + (1 - \beta) \theta_3 \\\dots\end{aligned}$$

$$\begin{aligned}v_{\theta} &:= 0 \\v_{\theta} &:= \beta v + (1 - \beta) \theta_1 \\v_{\theta} &:= \beta v + (1 - \beta) \theta_2 \\&\vdots \\&\xrightarrow{\text{Repeat}} v_{\theta} = 0 \\&\text{Get next } \theta_t \\v_{\theta} &:= \beta v_{\theta} + (1 - \beta) \theta_t\end{aligned}$$

Andrew Ng

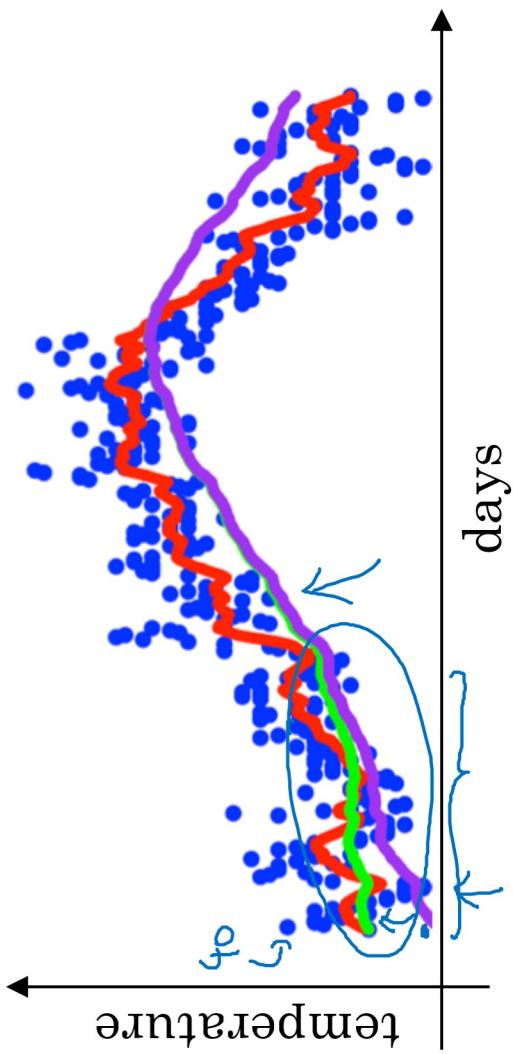
Optimization Algorithms

Bias correction
in exponentially
weighted average

deeplearning.ai



Bias correction



$$\beta = 0.98$$

$$\rightarrow v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$v_0 = 0$$

$$\frac{v_1}{v_1} = \underbrace{0.98 \times 0.98}_{0.96} + \underbrace{0.02 \theta_1}_{0.02}$$

$$v_2 = 0.98 \times 0.96 \times \theta_1 + 0.02 \theta_2$$

$$= 0.98 \times 0.96 \times \theta_1 + 0.02 \theta_2$$

$$= 0.96 \theta_1 + 0.02 \theta_2$$

$$\frac{v_t}{1 - \beta^t} = 1 - \beta^t$$

$$t=2: \quad \frac{v_2}{0.96} = 1 - (0.98)^2 = 0.0396$$

$$\frac{v_t}{1 - \beta^t} = 1 - (0.98)^2 = 0.0396$$

$$\frac{0.0396 \theta_1 + 0.02 \theta_2}{0.0396} =$$

$$0.0396$$

Andrew Ng

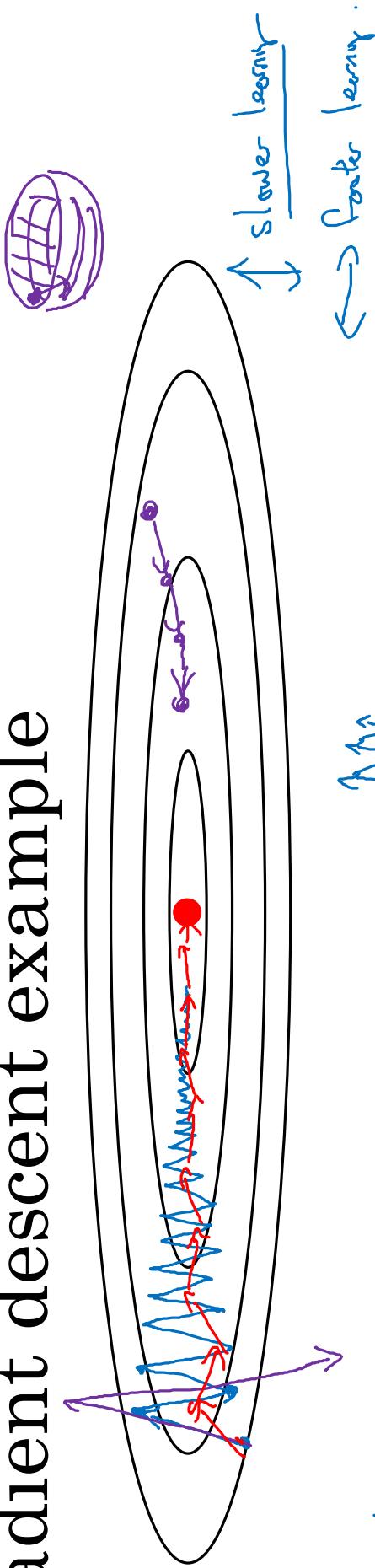
Optimization Algorithms

Gradient descent with momentum

deeplearning.ai



Gradient descent example



Momentum:

On iteration t :

Compute $\Delta w_t, \Delta b_t$ on current mini-batch.

$$v_{wt} = \beta v_{wt} + (1-\beta) \frac{\Delta w_t}{\|\Delta w_t\|}$$

$$\Delta w_t = v_{wt} + (1-\beta) \frac{\Delta w_t}{\|\Delta w_t\|}$$

Friction β \uparrow velocity

$$w_t = w_{t-1} - \alpha \Delta w_t, \quad b_t = b_{t-1} - \alpha \Delta b_t$$

$$v_t = \beta v_t + (1-\beta) \nabla J(w)$$

Implementation details

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

On iteration t :

Compute dW, db on the current mini-batch

$$\begin{aligned} \rightarrow v_{dw} &= \beta v_{dw} + \frac{(1 - \beta) dW}{\underbrace{}}_{\text{v}_{dw} = \beta v_{dw} + dW} \\ \rightarrow v_{db} &= \beta v_{db} + \frac{(1 - \beta) db}{\underbrace{}}_{\text{v}_{db} = \beta v_{db} + db} \\ W &= W - \alpha v_{dw}, \quad b = b - \alpha v_{db} \end{aligned}$$

~~v_{dw}~~ ~~v_{db}~~ ~~t~~

Hyperparameters: α, β

$$\beta = 0.9$$

~~our~~ lost ≈ 10 gradients

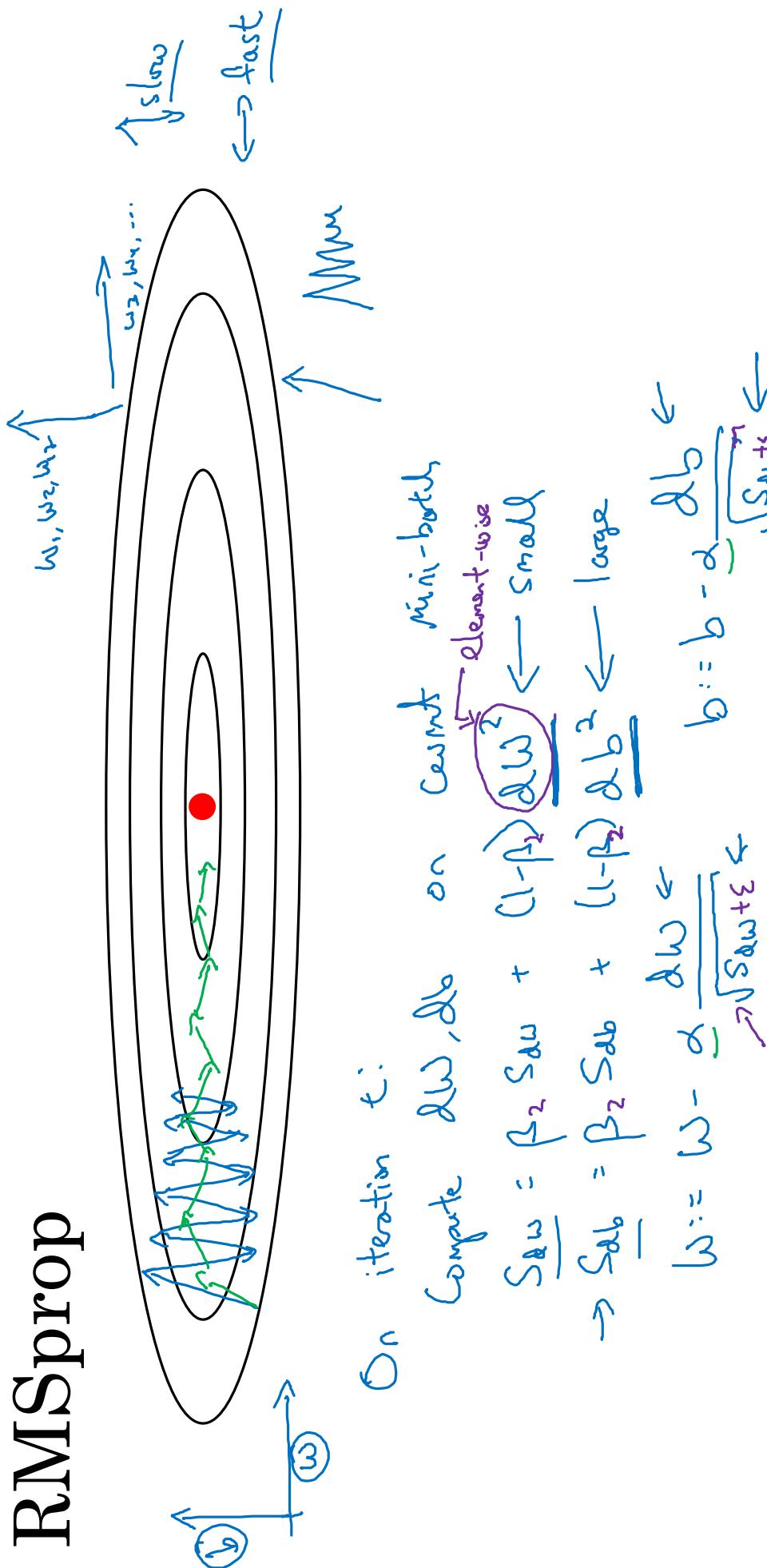
Optimization Algorithms

RMSprop



deeplearning.ai

RMSprop



Andrew Ng

Optimization Algorithms

Adam optimization algorithm

deeplearning.ai



Adam optimization algorithm

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$$

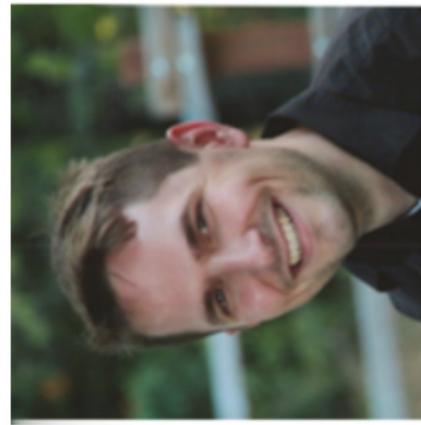
On iteration t :

Compute δ_{dw}, δ_{db} using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) \delta_{dw}, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) \delta_{db} \quad \leftarrow \text{"moment"} \quad (\beta_1)$$
$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) \delta_{dw}^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) \delta_{db}^2 \quad \leftarrow \text{"RMSprop"} \quad (\beta_2)$$
$$\hat{y} = np.array([.9, .02, .01, .4, .9])$$
$$\sqrt{\delta_{dw}} = V_{dw} / (1 - \beta_1^{t+1}), \quad \sqrt{\delta_{db}} = V_{db} / (1 - \beta_1^{t+1})$$
$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2^{t+1}), \quad S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2^{t+1})$$
$$b := b - \alpha \frac{\sqrt{S_{dw}^{\text{corrected}}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

Hyperparameters choice:

- α : needs to be tune
 - β_1 : 0.9 → $(\frac{d\omega}{\omega})$
 - β_2 : 0.999 → $(\frac{d\omega^2}{\omega^2})$
 - Σ : 10^{-8}
- Adam : Adaptive moment estimation



Adam Coates

Andrew Ng

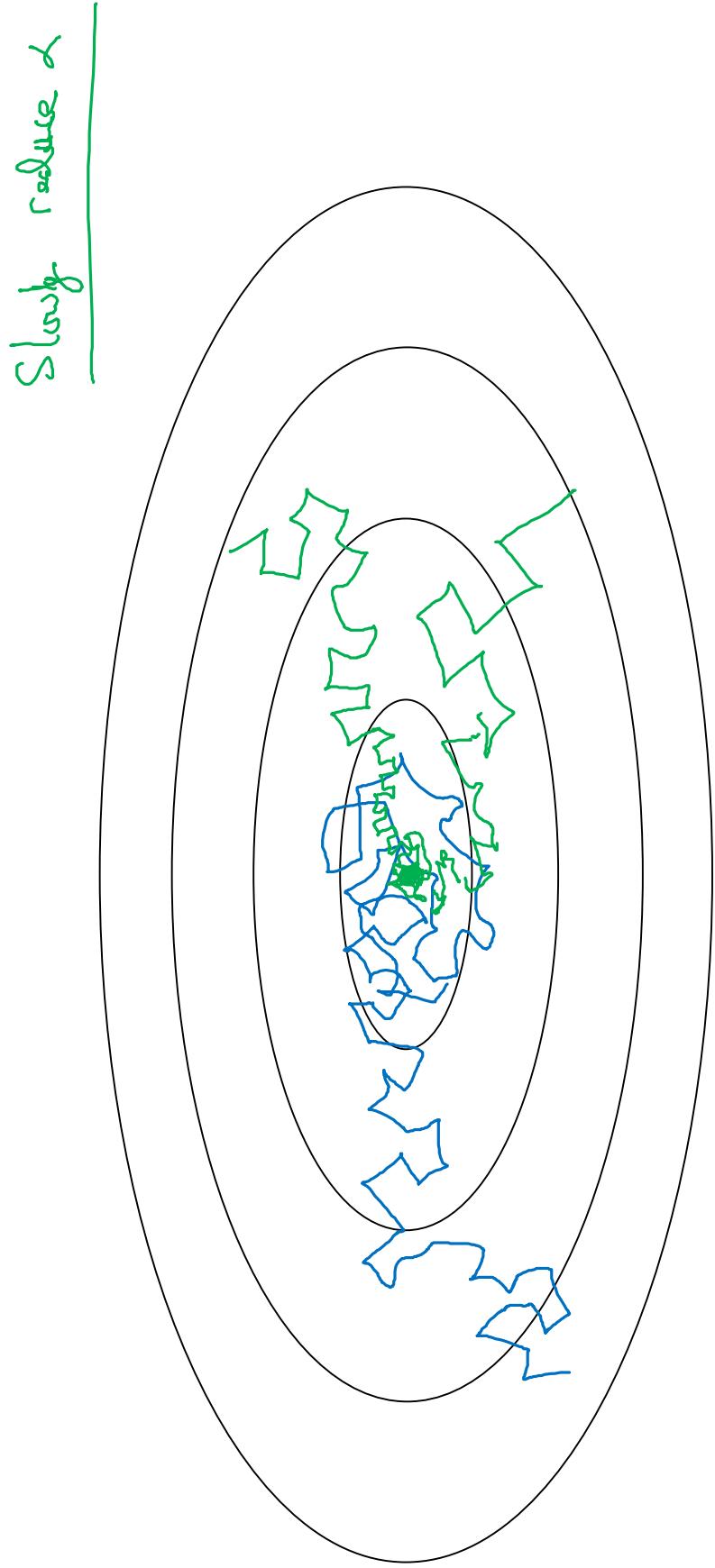
Optimization Algorithms

Learning rate decay



deeplearning.ai

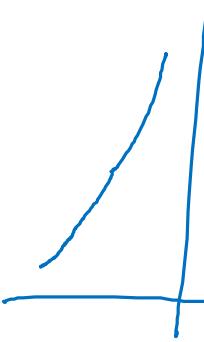
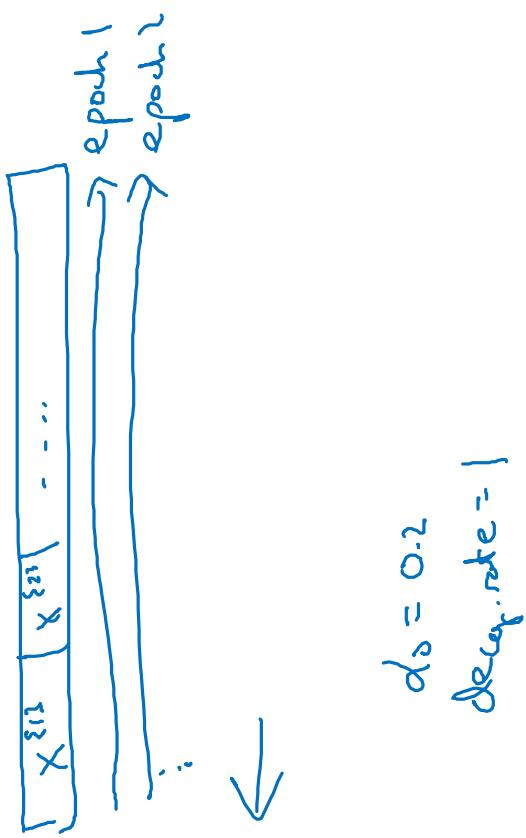
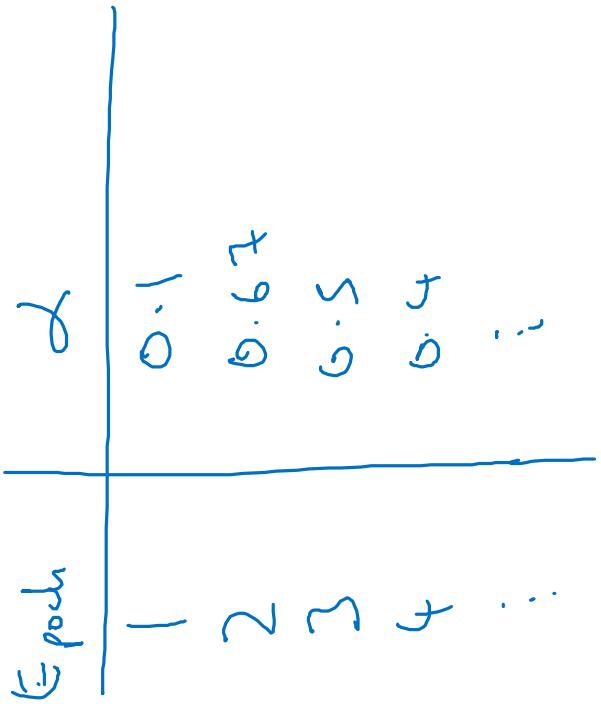
Learning rate decay



Learning rate decay

$\lfloor \text{epoch} = \lfloor \text{pass through data} \rfloor$

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



Other learning rate decay methods

$$\alpha = 0.95^{\text{epoch_num}} \cdot \alpha_0$$

$$\alpha = \frac{k}{\text{epoch_num}} \cdot \alpha_0$$

Linear

$$\text{or } \alpha = \frac{k}{\sqrt{t}} \cdot \alpha_0$$



discrete staircase

Manual decay

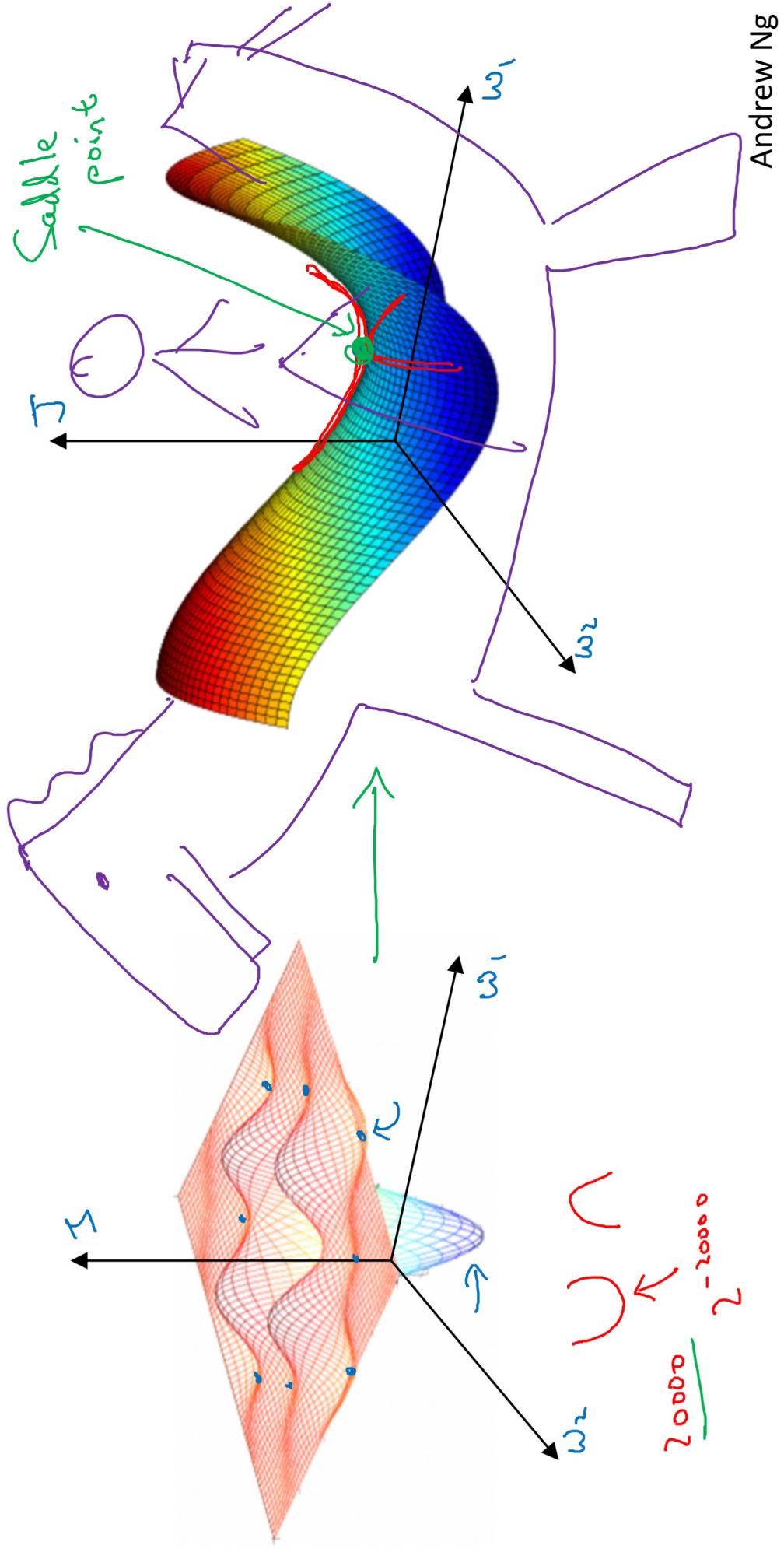
Optimization Algorithms

The problem of
local optima

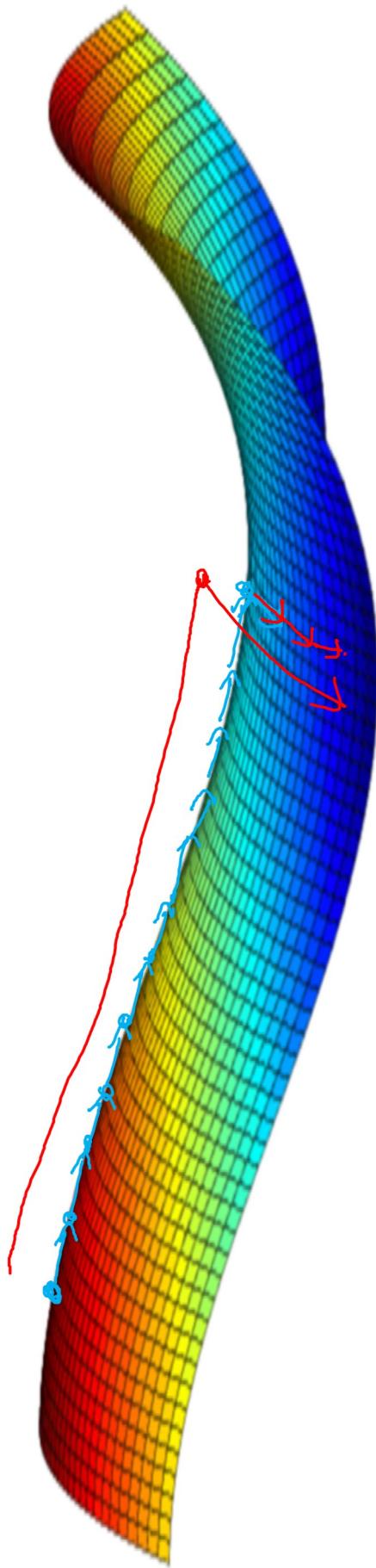


deeplearning.ai

Local optima in neural networks



Problem of plateaus



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