



INSTITUTO
METRÓPOLE
DIGITAL



IMD1122

Special Topics in AI Deep Learning

Training Neural Networks

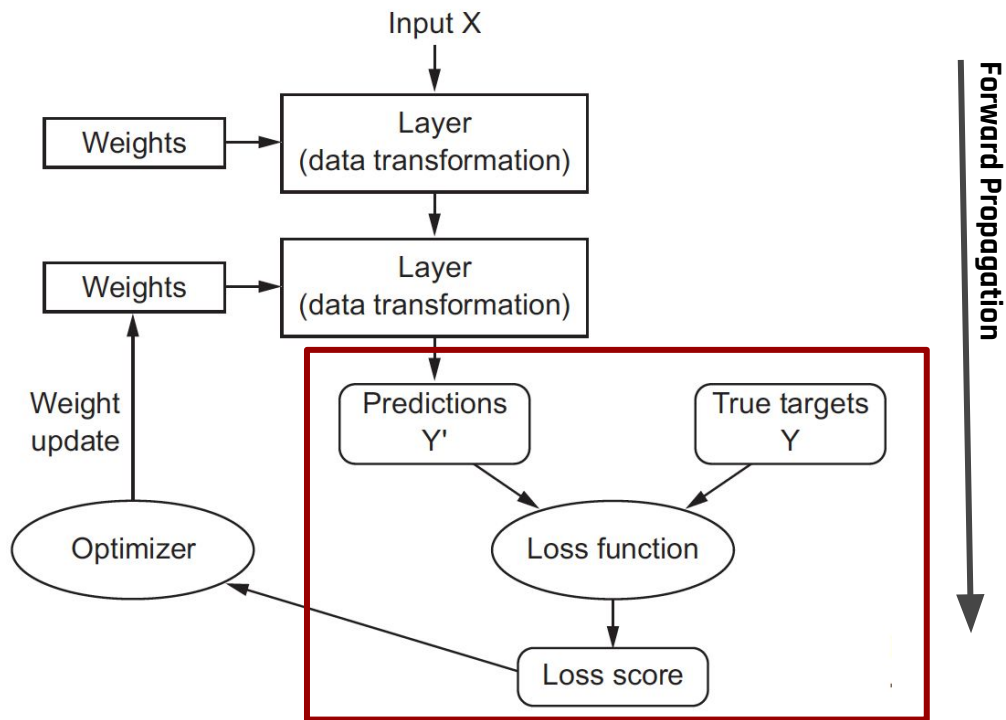
Lesson #03

Training Neural Networks

Part #01



Understanding how DL works



Loss Optimization

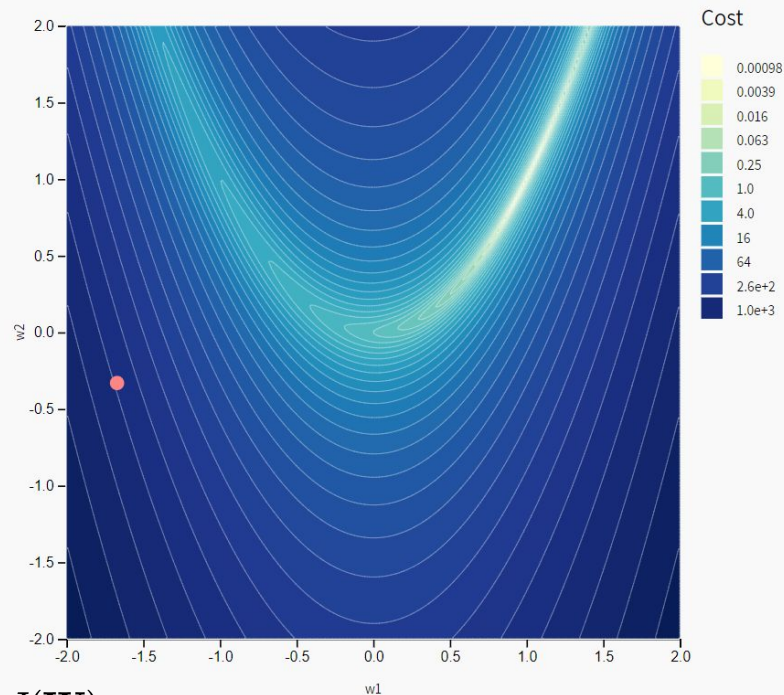
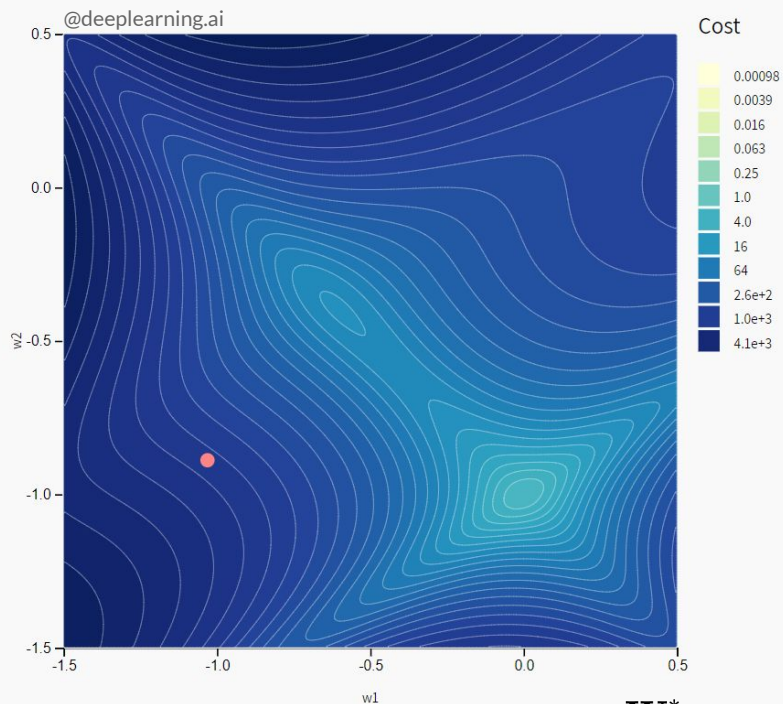
We want to find the network weights that **achieve the lowest loss**

$$\mathbf{W}^* = \underset{\mathbf{W}}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m \mathcal{L} \left(f(X^{(i)}; \mathbf{W}), y^{(i)} \right)$$

$$\mathbf{W}^* = \underset{\mathbf{W}}{\operatorname{argmin}} J(\mathbf{W})$$

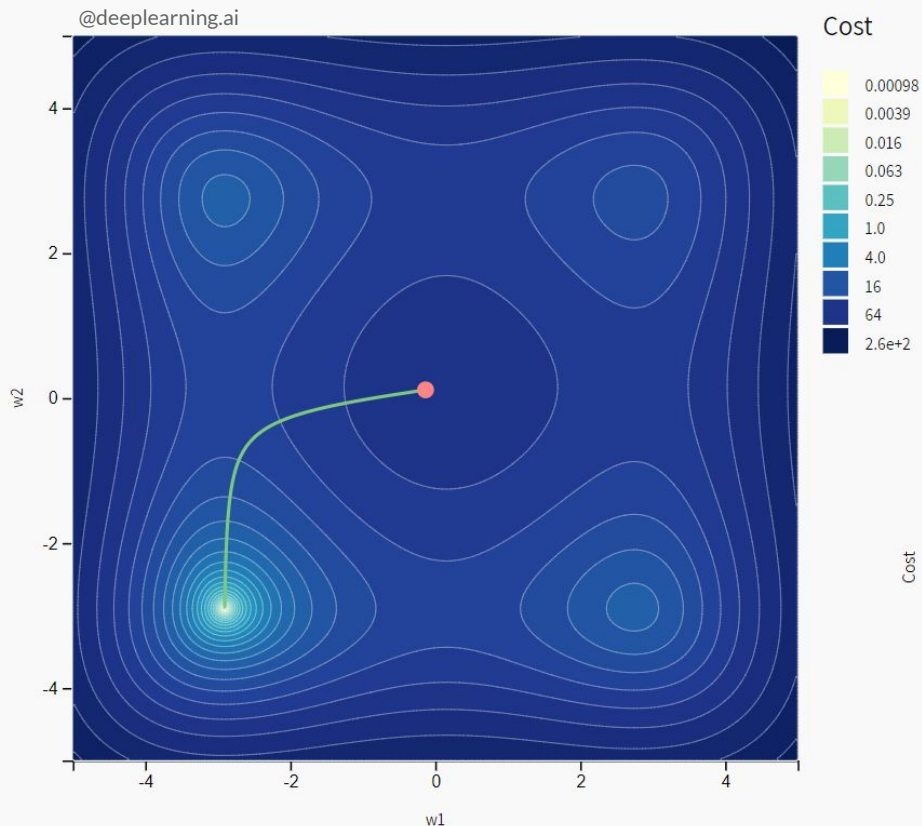
Loss Optimization

Remember: our loss is a function of the network weights!!!

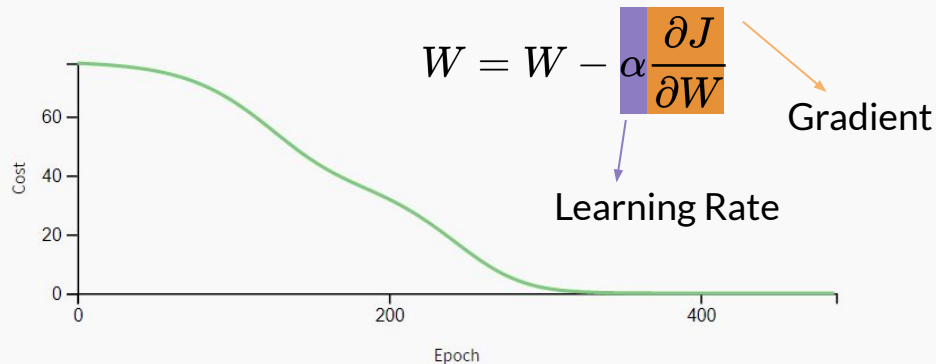


$$\mathbf{W}^* = \underset{\mathbf{W}}{\operatorname{argmin}} J(\mathbf{W})$$

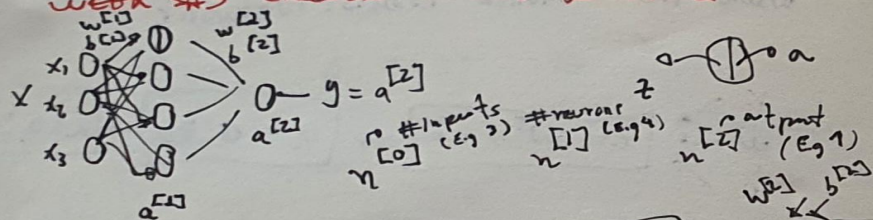
Loss Optimization - Gradient Descent



1. Initialize weights randomly
2. Compute gradient
3. Take small step in opposite direction of gradient
4. Repeat until convergence



Week #3 shallow Neural Networks



$$\begin{aligned} X &= \begin{bmatrix} x_1 & x_2 & x_3 \\ w[0] & w[1] & w[2] \\ b[0] & b[1] & b[2] \end{bmatrix} \\ z &= wX + b \\ a &= \sigma(z) \\ z &= w a + b \end{aligned}$$

$$\mathcal{L}(a, y) = \frac{1}{2} (a - y)^2$$

features
from instances

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

$$w^{[1]} = (n^{[1]}, n^{[0]})$$

$$b^{[1]} = (n^{[1]}, 1)$$

$$z^{[1]} = (n^{[1]}, 1)$$

$$x = (n^{[0]}, 1)$$

$$a^{[1]} = (n^{[1]}, 1)$$

$$\mathcal{L}(a, y) = -y \log(a^{[2]}) - (1-y) \log(1-a^{[2]})$$

$$\frac{d \mathcal{L}}{d a} = \frac{d \mathcal{L}}{d a} = -\frac{y}{a^{[2]}} + \frac{(1-y)}{1-a^{[2]}}$$

$$\frac{d \mathcal{L}}{d z} = \frac{d \mathcal{L}}{d a} \frac{d a}{d z} = \left[-\frac{y}{a^{[2]}} + \frac{(1-y)}{1-a^{[2]}} \right] a^{[2]} (1-a^{[2]})$$

$$\frac{d \mathcal{L}}{d z} = a^{[2]} - y$$

$$\frac{d \mathcal{L}}{d w} = \frac{d \mathcal{L}}{d z} \frac{d z}{d w} =$$

$$\frac{d \mathcal{L}}{d z} \cdot a^{[1]T}$$

$$(n^{[2]}, n^{[1]})$$

$$\frac{d \mathcal{L}}{d b} = \frac{d \mathcal{L}}{d z} \frac{d z}{d b} =$$

$$\frac{d \mathcal{L}}{d z} (n^{[2]}, 1)$$

$$(n^{[1]}, 1)$$

$$\frac{d \mathcal{L}}{d a} = \frac{d \mathcal{L}}{d z} \frac{d z}{d a} = \frac{d \mathcal{L}}{d z} \cdot w^{[2]}$$

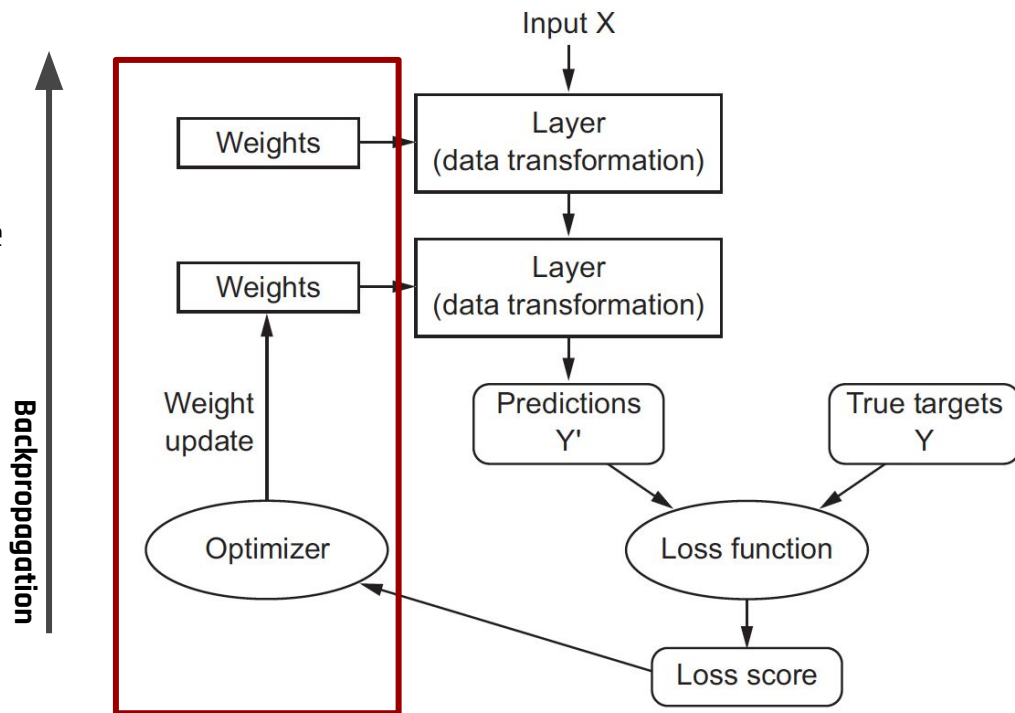
$$\frac{d \mathcal{L}}{d z} = \frac{d \mathcal{L}}{d a} \frac{d a}{d z} = \frac{d \mathcal{L}}{d a} \cdot g'(z^{[2]}) = (n^{[1]}, 1) \cdot (n^{[1]}, 1) = (n^{[1]T} \cdot d \mathcal{L}) \cdot g'(z^{[2]})$$

$$\frac{d \mathcal{L}}{d w} = \frac{d \mathcal{L}}{d z} \frac{d z}{d w} = \frac{d \mathcal{L}}{d z} \cdot x^T$$

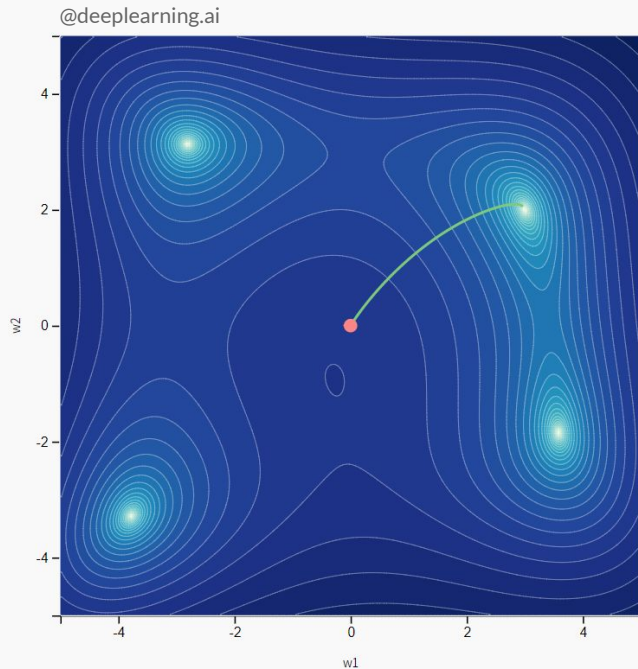
$$\frac{d \mathcal{L}}{d z} \cdot a^{[1]T}$$

Understanding how DL works

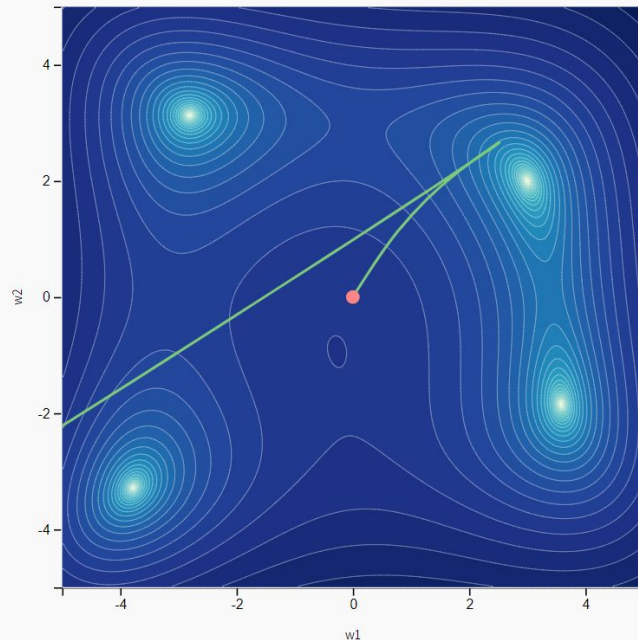
Finding the right values of weights which minimize the error



Loss Functions Can Be Difficult to Optimize



Small learning rate ($lr=0.001$)
converges slowly



Large learning rate ($lr=0.1$) overshoot,
become unstable and diverge

How to deal with this?

Idea 1:

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

Idea 2:

Do something smarter!!

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

Optimization Algorithms

Algorithm

- SGD
- Adam
- Adadelata
- Adagrad
- RMSProp

$$W = W - \alpha \boxed{?}$$

$\beta, \beta_1, \beta_2, learning_decay$

TF Implementation

```
tf.keras.optimizers.SGD
```



```
tf.keras.optimizers.Adam
```



```
tf.keras.optimizers.Adadelata
```



```
tf.keras.optimizers.Adagrad
```



```
tf.keras.optimizers.RMSprop
```



Putting it all together

Mini-Batches

$$1 \leq b \leq m$$

Model

Loss

Evaluation Metrics

Optimizer

```
# Create a source dataset from your training data
dataset = tf.data.Dataset.from_tensor_slices((train_set_x, train_set_y))
dataset = dataset.shuffle(buffer_size=64).batch(32)

# Instantiate a simple classification model
model = tf.keras.Sequential([
    tf.keras.layers.Dense(8, activation=tf.nn.relu, dtype='float64'),
    tf.keras.layers.Dense(8, activation=tf.nn.relu, dtype='float64'),
    tf.keras.layers.Dense(1, activation=tf.nn.sigmoid, dtype='float64')
])

# Instantiate a logistic loss function that expects integer targets.
loss = tf.keras.losses.BinaryCrossentropy()


# Instantiate an accuracy metric.
accuracy = tf.keras.metrics.BinaryAccuracy()

# Instantiate an optimizer.
optimizer = tf.keras.optimizers.SGD(learning_rate=0.001)
```



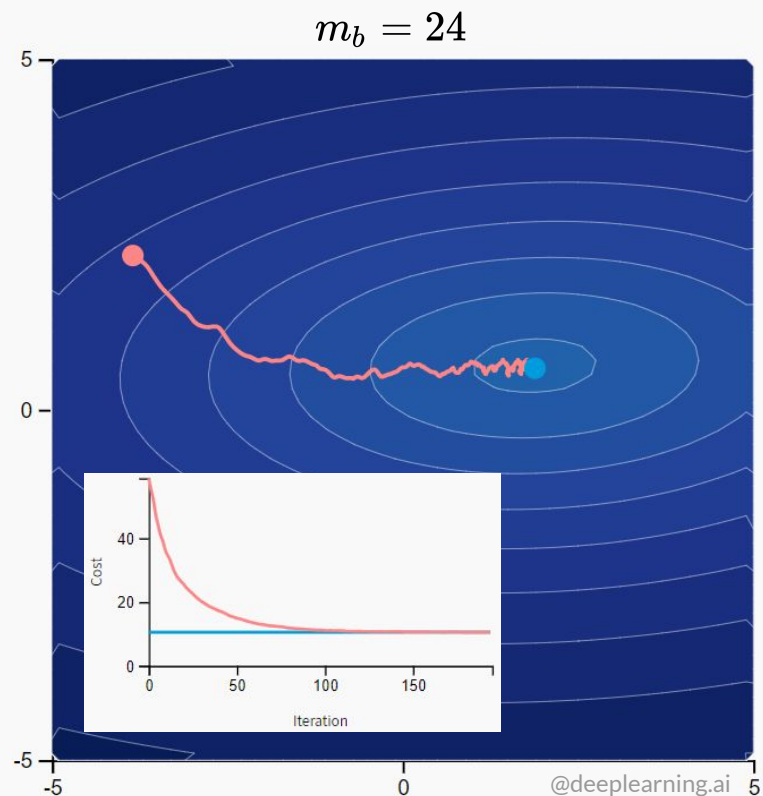
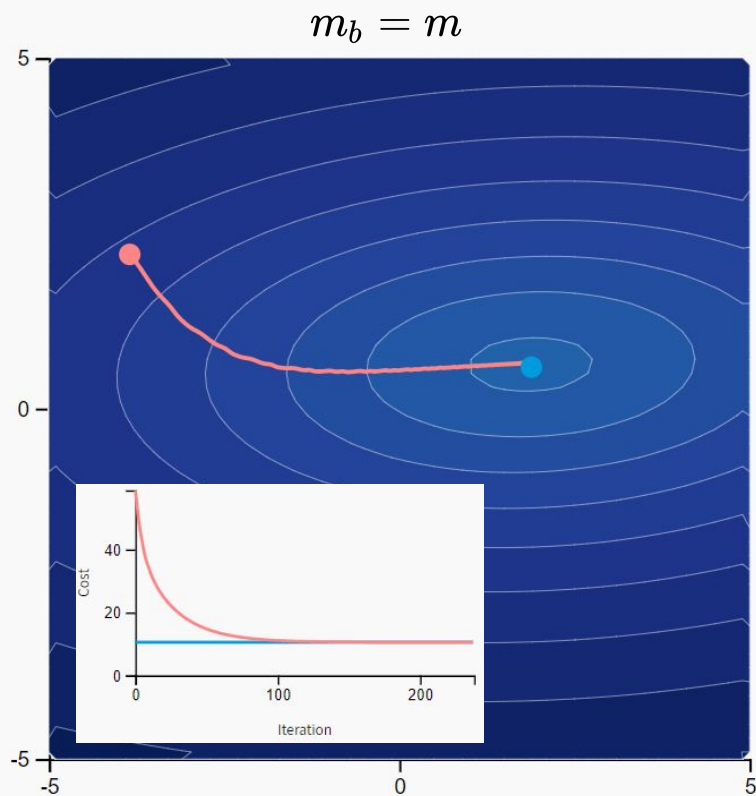
Putting it all together

```
for i in range(500):  
    # Iterate over the batches of the dataset.  
    for step, (x, y) in enumerate(dataset):  
        # Open a GradientTape.  
        with tf.GradientTape() as tape:  
  
            # Forward pass.  
            logits = model(x)  
  
            # Loss value for this batch.  
            loss_value = loss(y, logits)  
  
            # Get gradients of loss wrt the weights.  
            gradients = tape.gradient(loss_value, model.trainable_weights)  
  
            # Update the weights of our linear layer.  
            optimizer.apply_gradients(zip(gradients, model.trainable_weights))  
  
            # Update the running accuracy.  
            accuracy.update_state(y, logits)
```


$$W = W - \alpha \frac{\partial J}{\partial W}$$



Mini-Batches Challenges



Training Neural Networks

Part #02



TRAIN A DEEP NEURAL NETWORK

=

MINI-BATCH

+

OPTIMIZATION ALGORITHMS



Optimization Algorithms

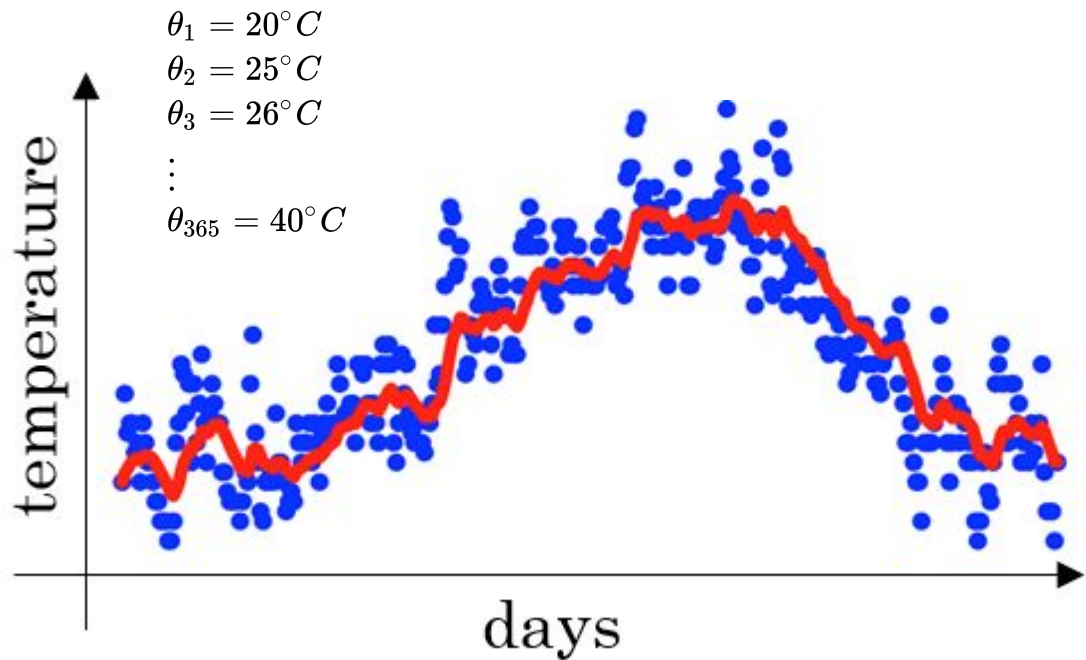
Exponentially Weighted Average

Adam

Momentum

RMSprop

Exponentially Weighted Average



$$V_0 = 0$$

$$V_1 = 0.9V_0 + 0.1\theta_1$$

$$V_2 = 0.9V_1 + 0.1\theta_2$$

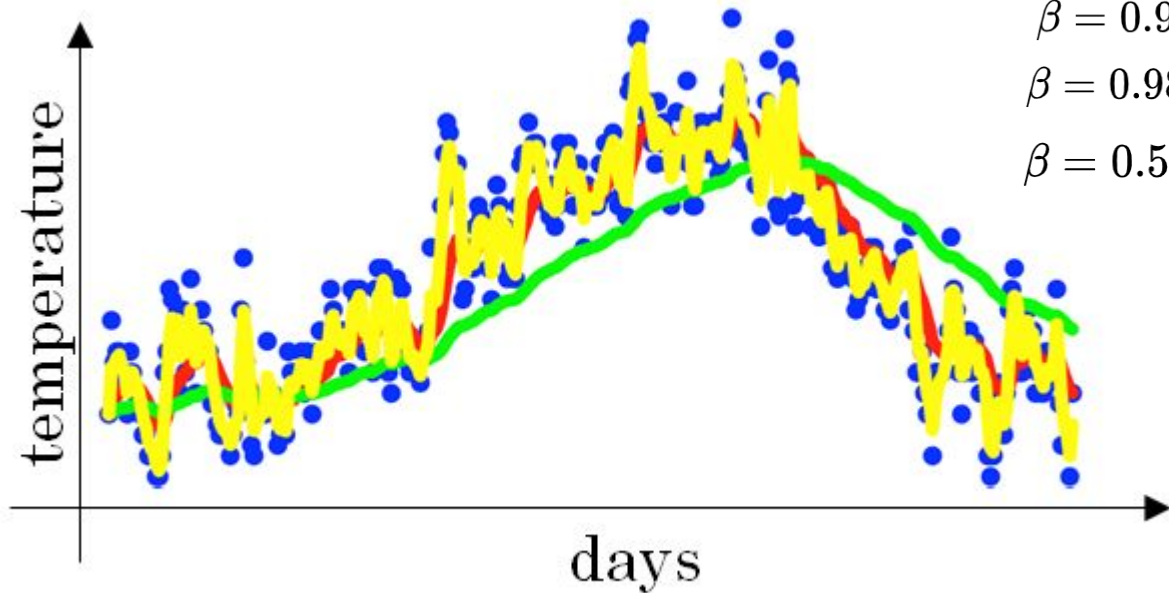
$$V_3 = 0.9V_2 + 0.1\theta_3$$

$$\vdots$$

$$V_{365} = 0.9V_{364} + 0.1\theta_{365}$$

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

Exponentially Weighted Average



$\beta = 0.9 \rightarrow V_t \approx 10 \text{ days of temperature}$

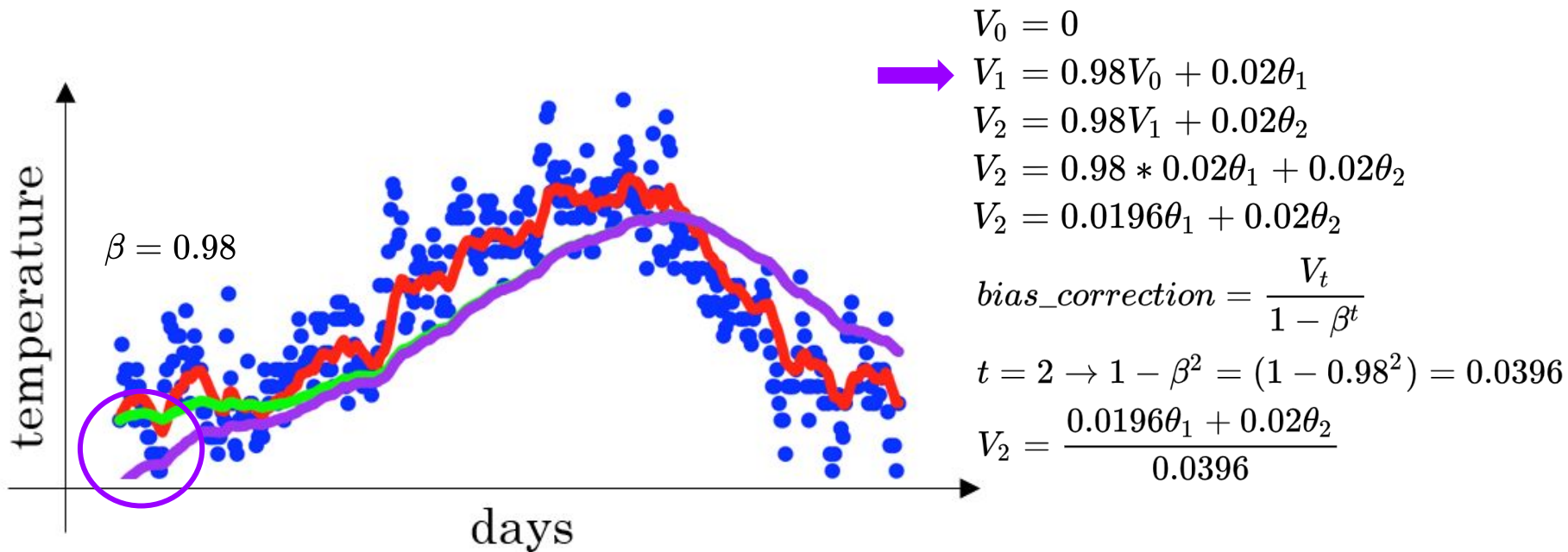
$\beta = 0.98 \rightarrow V_t \approx 50 \text{ days of temperature}$

$\beta = 0.5 \rightarrow V_t \approx 2 \text{ days of temperature}$

$$V_t = \beta V_{t-1} + (1 - \beta)\theta_1$$

$$V_t \approx \frac{1}{1 - \beta} \text{ days of temperature}$$

Exponentially Weighted Average Bias Correction

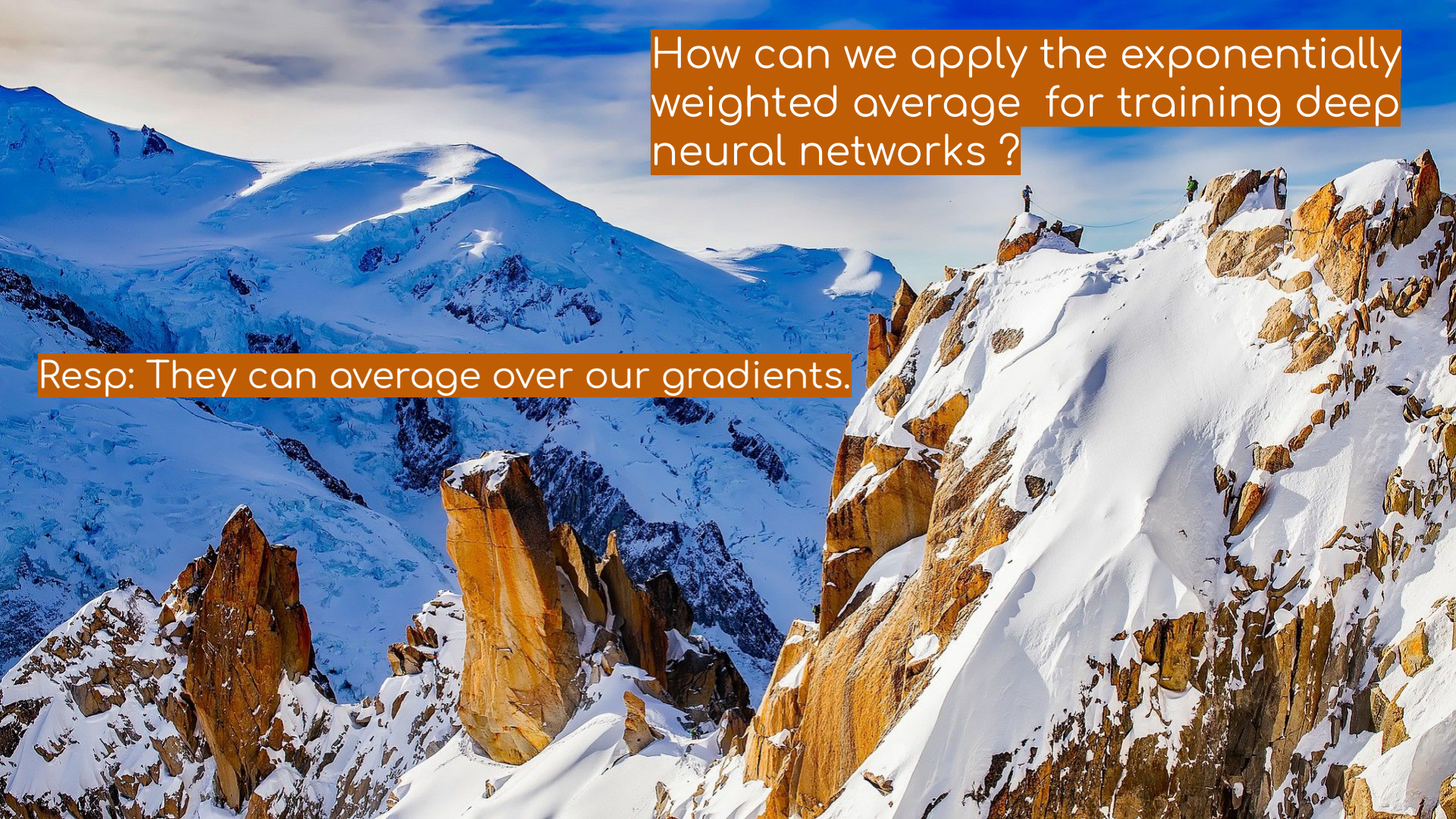


PETROBRAS PN N2 - D - BMFBOVESPA O29.51 H29.68 L29.03 C29.14 -0.78 (-2.61%)

EMA 50 close 29.37

SELL 29.08 0.09 BUY 29.17  + - losed Delayed

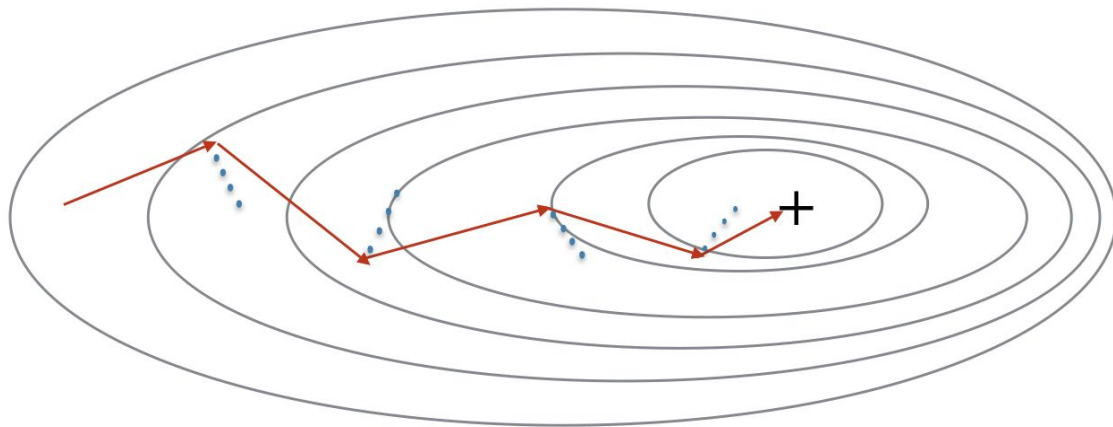


A high-altitude mountain landscape. In the foreground, a steep, rocky ridge is covered in patches of snow. The rocks are a warm, golden-brown color. In the background, more snow-covered mountain peaks are visible under a clear blue sky with some light clouds. The overall scene is bright and crisp.

How can we apply the exponentially weighted average for training deep neural networks ?

Resp: They can average over our gradients.

Gradient Descent with Momentum



- Momentum takes into account the past gradients to smooth out the update.
- Formally, this will be the exponentially weighted average of the gradient on previous steps.

Gradient Descent with Momentum

On iteration t :

Compute dW, db on the current mini-batch

$$v_{dW} = \beta v_{dW} + (1 - \beta) dW$$

$$v_{db} = \beta v_{db} + (1 - \beta) db$$

$$W = W - \alpha v_{dW}, \quad b = b - \alpha v_{db}$$

Hyperparameters: α, β $\beta = 0.9$

Gradient Descent with RMSprop

On iteration t :

Compute dW, db on the current mini-batch

$$s_{dW} = \beta s_{dW} + (1 - \beta) dW^2$$

$$s_{db} = \beta s_{db} + (1 - \beta) db^2$$

$$W = W - \alpha \frac{dW}{\sqrt{s_{dW}} + \epsilon} \quad b = b - \alpha \frac{db}{\sqrt{s_{db}} + \epsilon}$$

$$\epsilon = 10^{(-8)}$$

Hyperparameters: α, β

$$\beta = 0.9$$

Gradient Descent with Adam

On iteration t :

Compute dW, db on the current mini-batch

$$v_{dW} = \beta v_{dW} + (1 - \beta) dW$$

$$v_{db} = \beta v_{db} + (1 - \beta) db$$

$$s_{dW} = \beta s_{dW} + (1 - \beta) dW^2$$

$$s_{db} = \beta s_{db} + (1 - \beta) db^2$$

$$v_{dW}^{correct} = v_{dW} / (1 - \beta_1^t)$$

$$v_{db}^{correct} = v_{db} / (1 - \beta_1^t)$$

$$s_{dW}^{correct} = s_{dW} / (1 - \beta_2^t)$$

$$s_{db}^{correct} = s_{db} / (1 - \beta_2^t)$$

$$W = W - \alpha \frac{v_{dW}^{correct}}{\sqrt{s_{dW}^{correct} + \epsilon}}$$

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

$\epsilon = 10^{-8}$

Hyperparameters: α, β_1, β_2

$\beta_1 = 0.9, \beta_2 = 0.999$

Learning Rate Decay

Idea: reduce α for each mini-batch t in order to smooth the gradient descent.

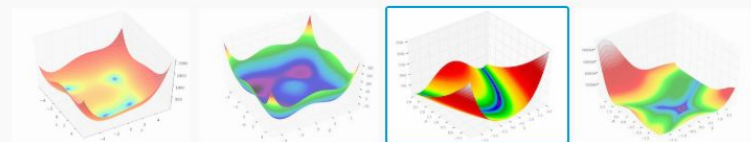
$$\alpha = \frac{1}{1 + \textit{learning_decay} * \textit{epoch_num}} * \alpha_0$$

$$\alpha = 0.95 e^{\textit{epoch_num}} \alpha_0$$

$$\alpha = \frac{k}{\sqrt{\textit{epoch_num}}} * \alpha_0$$

1. Choose a cost landscape

Select an artificial landscape $\mathcal{J}(w_1, w_2)$.



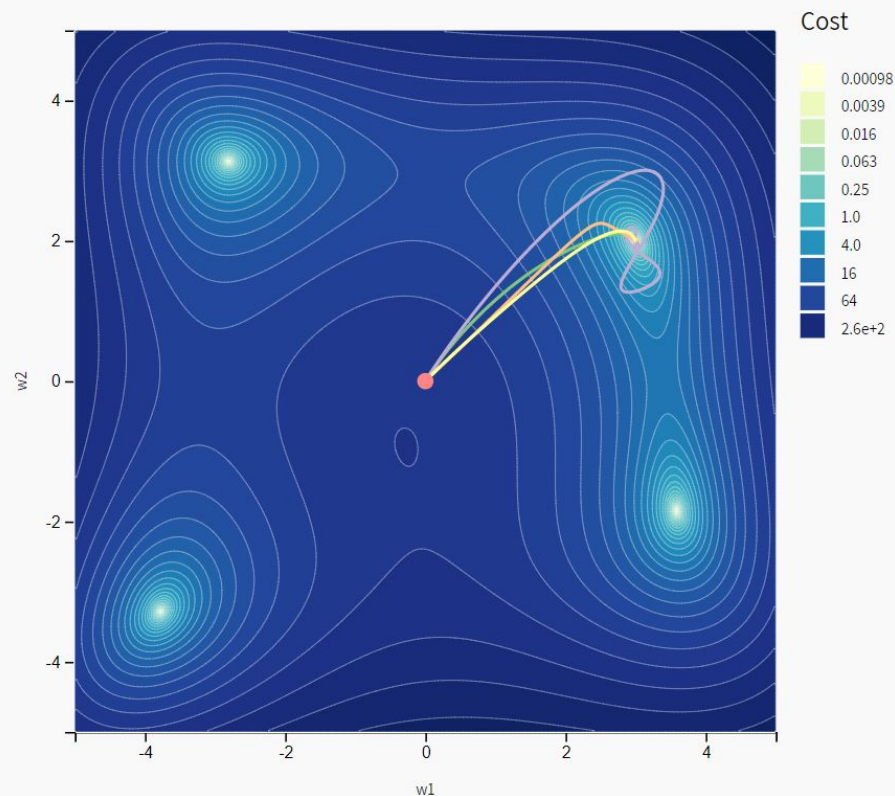
2. Choose initial parameters

On the cost landscape graph, drag the **red dot** to choose initial parameter values and thus the initial value of the cost.

3. Choose an optimizer

Select the optimizer(s) and hyperparameters.

Optimizer	Learning Rate	Learning Rate Decay
<input checked="" type="checkbox"/> Gradient Descent	<input type="text" value="0,001"/>	<input type="text" value="0"/>
<input checked="" type="checkbox"/> Momentum	<input type="text" value="0,001"/>	<input type="text" value="0"/>
<input checked="" type="checkbox"/> RMSprop	<input type="text" value="0,001"/>	<input type="text" value="0"/>
<input checked="" type="checkbox"/> Adam	<input type="text" value="0,001"/>	<input type="text" value="0"/>



<https://www.deeplearning.ai/ai-notes/optimization/>

Neural Networks in Practice #01

Splitting Data & Regularization

Next

