

Title: Machine Learning Lecture

Speakers: Mohamed Hibat Allah

Collection: Machine Learning 2023/24

Date: April 09, 2024 - 11:30 AM

URL: <https://pirsa.org/24040048>

# Lecture 4

## Last Lecture:

\* Architecture of feed-forward Neural Networks (FFNNs) for supervised learning (SL)

→ Weights  $W_{ij}^{(l)}$ , biases  $b_i^l$ , activation functions

→ Expressivity of FFNNs

→ Basis of NNs training: cost functions

Let's continue our discussion on how NNs learn:

\* Outline for today:

↳ More on NN training ("learning")

→ Learning algorithms

→ Back propagation

↳ Overfitting

## \* Learning algorithms

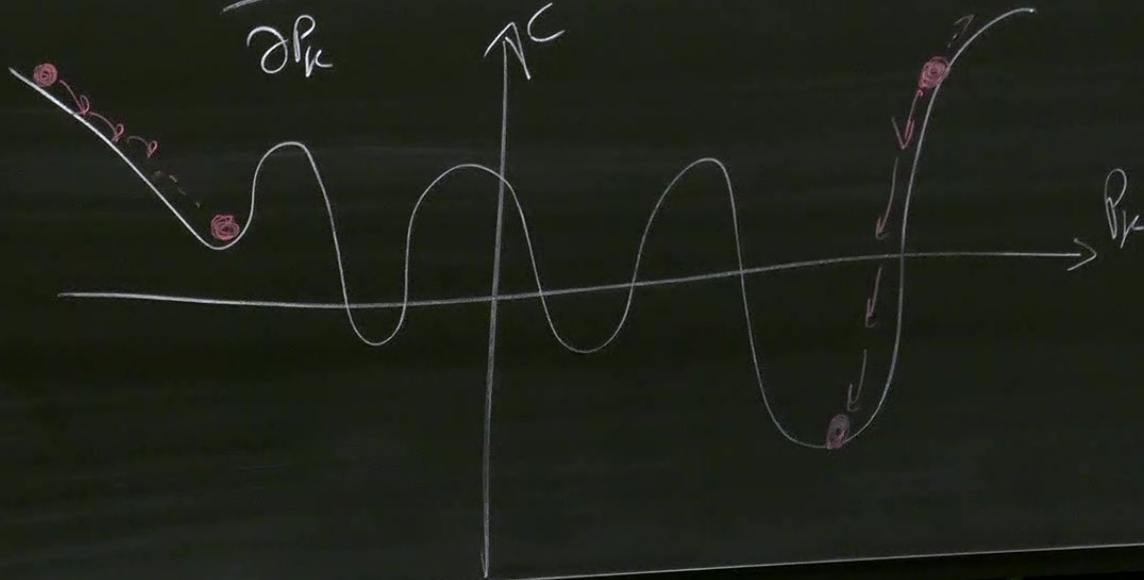
We use a learning algorithm (such as gradient descent) to minimize a cost function with respect to all weights and bias in each layer  $l$  such that:

$$\begin{array}{l} \textcircled{1} \\ \textcircled{2} \end{array} \left\{ \begin{array}{l} \frac{\partial C}{\partial W_{ij}^{(l)}} = 0 \\ \frac{\partial C}{\partial b_j^{(l)}} = 0 \end{array} \right. \rightarrow \text{for } \forall i, j, l$$

Unlike linear regression, # of parameters (Weights & biases) can quickly increase (#P can be  $\gg$   $d_{in}$ )  
 $\hookrightarrow$  It becomes nearly impossible to solve all equations  $\textcircled{1}$  and  $\textcircled{2}$  to get  $W$  &  $b$

$$\vec{p} = (p_1, p_2, \dots, p_N)$$

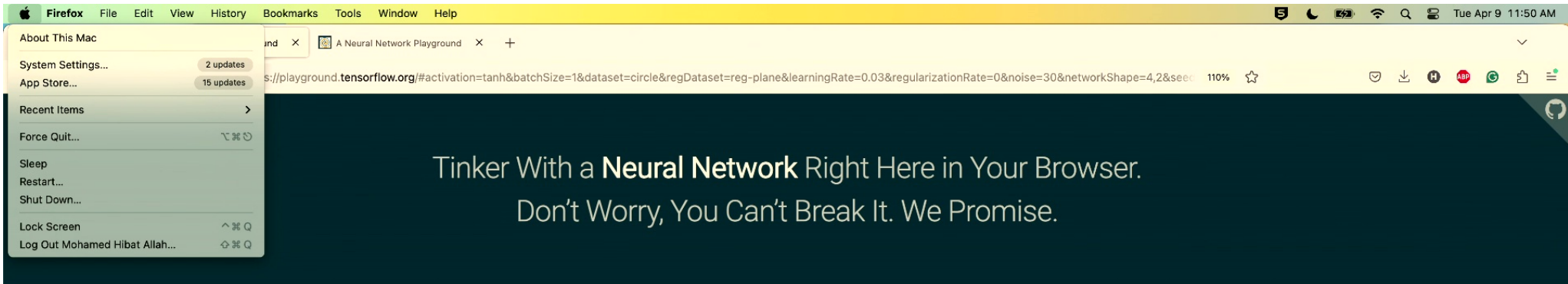
$$\frac{\partial C}{\partial p_k} = 0 \quad \forall (1 \leq k \leq N)$$



Gradient descent (GD) update.

$$p_k \leftarrow p_k - \eta \frac{\partial C}{\partial p_k}$$

Learning rate.



Epoch: 000,000 | Learning rate: 0.03 | Activation: Tanh | Regularization: None | Regularization rate: 0 | Problem type: Classification

**DATA**  
Which dataset do you want to use?  
Ratio of training to test data: 50%  
Noise: 30  
Batch size: 1  
REGENERATE

**FEATURES**  
Which properties do you want to feed in?  
X1, X2, X1<sup>2</sup>, X2<sup>2</sup>, X1X2, sin(X1), sin(X2)

**2 HIDDEN LAYERS**  
4 neurons | 2 neurons  
This is the output from one neuron. Hover to see it larger.  
The outputs are mixed with varying weights, shown by the thickness of the lines.

**OUTPUT**  
Test loss 0.523  
Training loss 0.524  
Colors shows data, neuron and weight values.

Firefox File Edit View History Bookmarks Tools Window Help

https://playground.tensorflow.org/#activation=tanh&batchSize=1&dataset=circle&regDataset=reg-plane&learningRate=0.03&regularizationRate=0&noise=30&networkShape=4,2&seed=133%

Ratio of training to test data: 50%

Noise: 30

Batch size: 1

X1

X2

X12

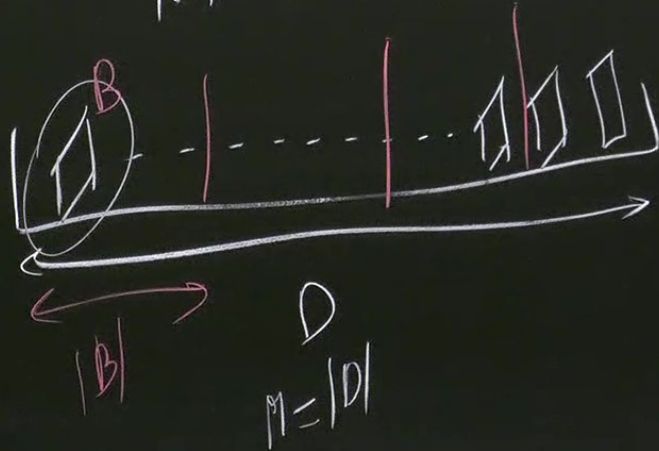
X22

X1X2

This is the output from one neuron. Hover to see it

# Stochastic Gradient Descent (SGD) \* Mini-batch Gradient Descent

$$C = \frac{1}{|D|} \sum_{x \in D} C(x)$$



\* SGD  $|B|=1$

$$P_k \leftarrow P_k - \eta \partial_x C(x^{(i)})$$

$x^{(i)} \in D$   
 Random order  
 $i=1, \dots, M$

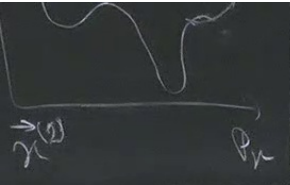
\* Mini-batch

$$C^B = \frac{1}{|B|} \sum_{x \in B} C(x) \rightarrow P_k \leftarrow P_k - \eta \partial_x C^B$$

$$|B| = O(100)$$

$$|B|=2 \rightarrow |D|$$





① Random order  
↓  
i=1, ..., M

→  $P_k \leftarrow P_k - \eta \frac{\partial L}{\partial P_k}$

# Automatic Differentiation (AD)

$$f(x) = x^2 \rightarrow f'(x) = 2x$$



\* Adding momentum, velocity Acc  $\beta=0 \rightarrow GD$

$$\vec{v} \rightarrow \beta \vec{v} + (1 - \beta) \frac{\partial C}{\partial \mathbf{p}}$$

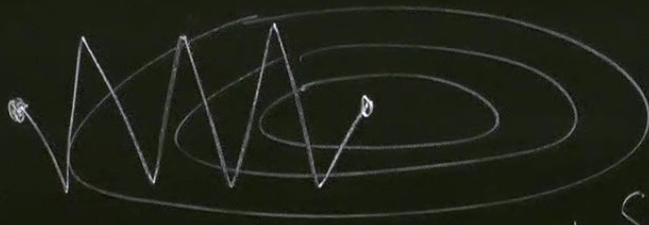
momentum decay

$$0 < \beta < 1$$

$$\vec{p} \leftarrow \vec{p} - \eta \vec{v}$$

Learning rate

No momentum



Oscillations

=

Slow learning

Momentum



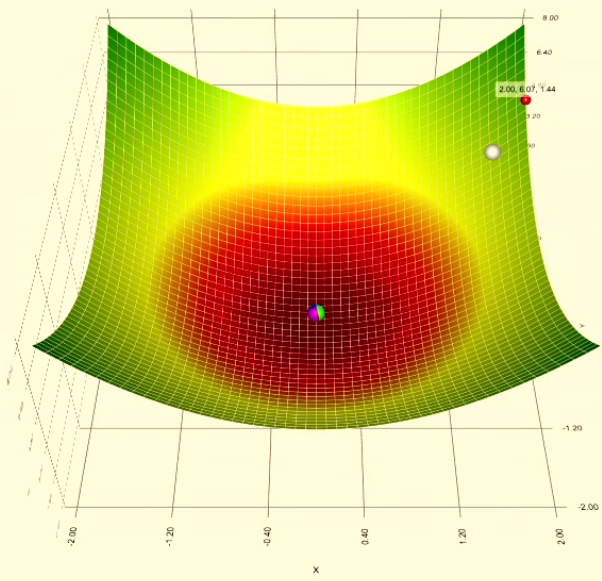
\* Second Order optimizers.

\* Adaptive learning rate.

\* RMS Prop

\* Adam.

→ Pytorch (torch.optim)



Global Minimum ⌵

**Overview** Step-by-Step

- Gradient Arrows
- Adjusted Gradient Arrows
- Momentum Arrows
- Sum of Gradient Squared
- Path

**Gradient Descent**

Learning Rate:  $1e^{-2}$  ⌵

**Momentum**

Learning Rate:  $1e^{-2}$  ⌵

Decay rate: 0.800 ⌵

**Adagrad**

Learning Rate:  $1e^{-3}$  ⌵

**RMSprop**

Learning Rate:  $1e^{-2}$  ⌵

Decay rate: 0.990 ⌵

**Adam**

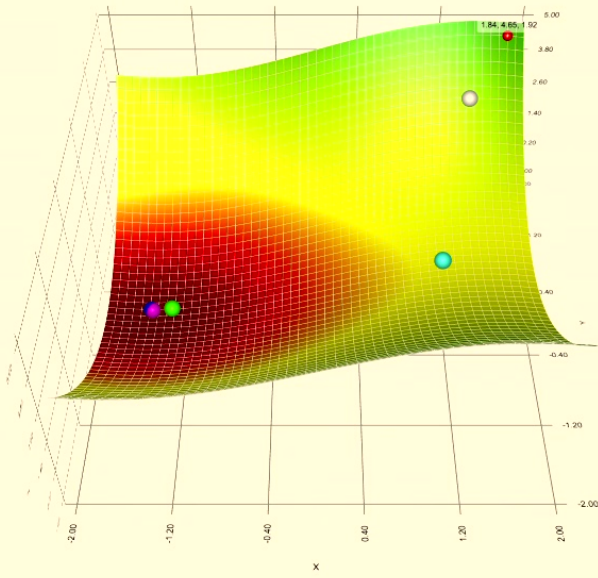
Learning Rate:  $1e^{-2}$  ⌵

Beta1: 0.900 ⌵

Beta2: 0.999 ⌵

Pause Restart Playback speed: 0.2x 🔍 🔍





Saddle Point ⌵

Overview Step-by-Step

- Gradient Arrows
- Adjusted Gradient Arrows
- Momentum Arrows
- Sum of Gradient Squared
- Path

Gradient Descent

Learning Rate:  $1e^{-2}$  ⌵

Momentum

Learning Rate:  $1e^{-2}$  ⌵

Decay rate: 0.800 ⌵

Adagrad

Learning Rate:  $1e^{-3}$  ⌵

RMSprop

Learning Rate:  $1e^{-2}$  ⌵

Decay rate: 0.990 ⌵

Adam

Learning Rate:  $1e^{-2}$  ⌵

Beta1: 0.900 ⌵

Beta2: 0.999 ⌵

Pause Restart Playback speed: 0.2x ⌵ ⌶



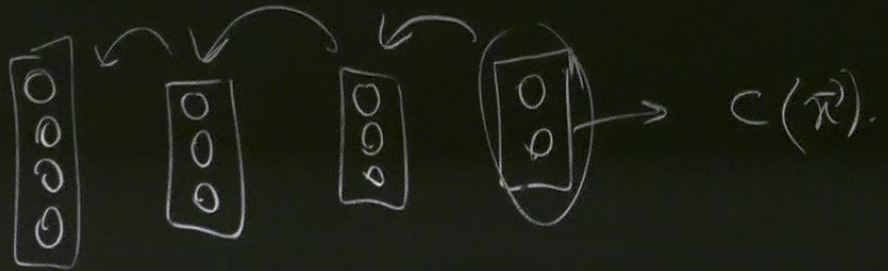
Test

# Back propagation in FFNNs

$$\frac{\partial C}{\partial W_{ij}^{(l)}} , \frac{\partial C}{\partial b_i^{(l)}}$$

$$C = \frac{1}{|D|} \sum_{x \in D} C(\vec{x})$$

One training sample  $\vec{x}$ .



\* Adam.

$$a_j^{(l)} = g_l(z_j^{(l)}) \quad 0 \leq l < L$$

$$z_j^{(l)} = \sum_{i=1}^{n_{l-1}} a_i^{(l-1)} W_{ij} + b_j^{(l)}$$

$$\delta_j^{(l)} = \frac{\partial c}{\partial z_j^{(l)}}$$

$$c = c(\vec{a}^{(L)})$$

$$\delta_j^{(l)}$$

$$= \frac{\partial c}{\partial z_j^{(l)}}$$

$$= \frac{\partial c}{\partial a_j^{(l)}} \cdot \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}}$$

$$= \frac{\partial c}{\partial a_j^{(l)}} \cdot g_l'(z_j^{(l)}) \quad \uparrow$$

$$\delta_i^{(l)} = \sum_{k=1}^{n_{l+1}} \delta_k^{(l+1)} W_{kj}^{(l+1)T} g'_l(z_j^{(l)}) \quad (2)$$

$$\textcircled{3} \quad \frac{\partial C}{\partial w_{ij}^{(l)}} = a_i^{(l-1)} \delta_j^{(l)}$$

$$\frac{\partial C}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad (4)$$



Firefox File Edit View History Bookmarks Tools Window Help

https://playground.tensorflow.org/#activation=tanh&batchSize=10&dataset=gauss&regDataset=reg-plane&learningRate=0.03&regularizationRate=0&noise=35&networkShape=4,2&se 133%

Epoch: 000,000 | Learning rate: 0.03 | Activation: Tanh | Regularization: None | Regularization rate: 0 | Problem type: Classification

**DATA**

Which dataset do you want to use?

Ratio of training to test data: 20%

Noise: 35

Batch size: 10

REGENERATE

**FEATURES**

Which properties do you want to feed in?

X1, X2, X1<sup>2</sup>, X2<sup>2</sup>, X1X2, sin(X1), sin(X2)

+ - 2 HIDDEN LAYERS

+ -

4 neurons

+ -

2 neurons

*This is the output from one neuron. Hover to see it larger.*

*The outputs are mixed with varying weights, shown by the thickness of the lines.*

**OUTPUT**

Test loss 0.495  
Training loss 0.481

Colors shows data, neuron and weight values.

Show test data  Discretize output

Don't worry, you can't break it. We promise.



Epoch  
000,000

Learning rate  
0.03

Activation  
Tanh

Regularization  
None

Regularization rate  
0

Problem type  
Classification

DATA

Which dataset do you want to use?



Ratio of training to test data: 20%



Noise: 35



Batch size: 10



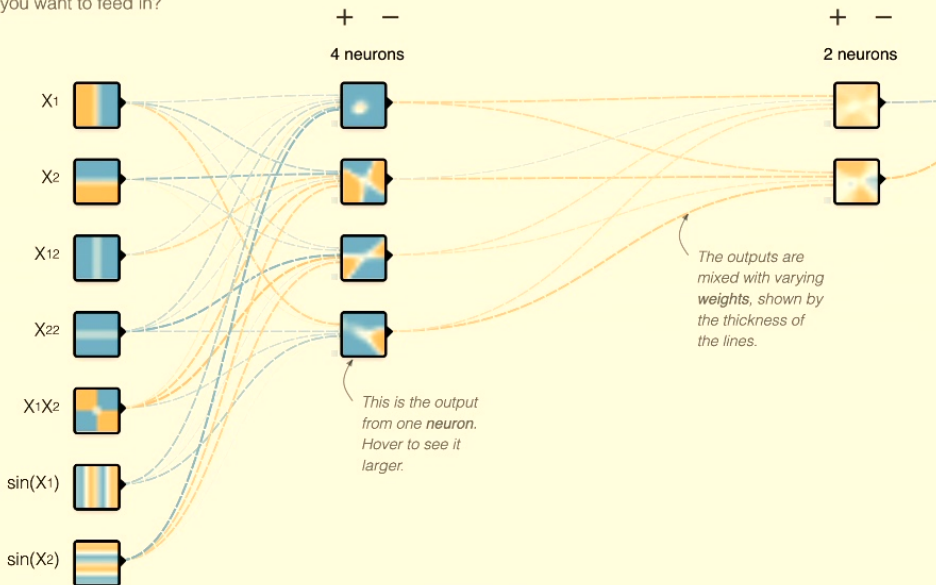
REGENERATE

FEATURES

Which properties do you want to feed in?

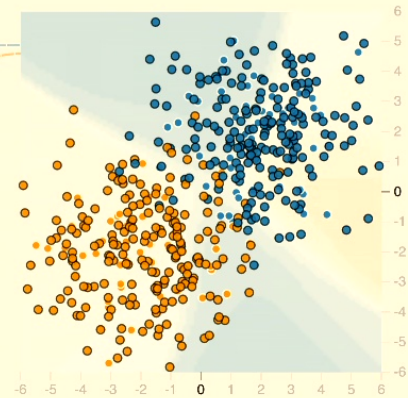
- X1
- X2
- X1X2
- sin(X1)
- sin(X2)

2 HIDDEN LAYERS



OUTPUT

Test loss 0.495  
Training loss 0.481



Colors shows data, neuron and weight values.

Show test data  Discretize output



Firefox File Edit View History Bookmarks Tools Window Help

https://playground.tensorflow.org/#activation=tanh&batchSize=10&dataset=gauss&regDataset=reg-plane&learningRate=0.03&regularizationRate=0&noise=35&networkShape=4,2&... 133%

Learning rate: 0.03    Activation: Tanh    Regularization: None    Regularization rate: 0    Problem type: Classification

+ - 2 HIDDEN LAYERS

4 neurons    2 neurons

Test loss 0.135  
Training loss 0.024

The outputs are mixed with varying weights, shown by the thickness of the lines.

This is the output from one neuron.