

# Q&A Lessons 3,4

01 - 07 Feb 2021

Poll results

# Table of contents

- What videos were engaged?
- What are your main difficulties with this methodology and material?
- How are you consuming the course material?
- Survey

## What videos were engaged?

(1/2)

0 2 5

[Lesson 01] Google Colab Introduction ✓



[Lesson 01] Google Colab Cont. ✓



[Lesson 02] The perceptron ✓



[Lesson 02] Building neural network ✓



[Lesson 02] Applying neural network ✓



[Lesson 03] Training neural network ✓



## What videos were engaged? (2/2)

025

[Lesson 04] Training neural network ✓



What are your main difficulties with this methodology and material?

020

Conciliar com estágio  
Material em língua inglesa  
Falta de prática com a linguagem engajamento  
falta de disciplina  
Densidade do conteúdo  
gestão de tempo **tempo**  
nenhuma Notebook

# Tempo para estudar todos os notebooks

Carga de conhecimento anterior  
tá acelerado, entrei depois  
Sinto falta de hands on  
os primeiros notebooks são enormes  
Tempo para fazer os notebooks  
Dificuldade na língua inglesa  
Quantidade de tópicos no notebook

## How are you consuming the course material?

0 1 4

- Video > notebook
- video > notebook
- Assistir os videos, acompanhando os notebooks e as duvidas vou buscar na bibliografia.
- video e depois o notebook especifico
- Assisto o vídeo e depois o notebook
- Video depois notebooka
- Vídeo e depois notebook
- Assisto o vídeo e depois faço o notebook
- Vídeo depois o notebook
- Vídeos primeiro e notebook depois
- assistindo videos e indo pro notebook; tb lendo coisa por fora
- Assisto os vídeos e depois vou para o notebook
- Vídeo depois notebook
- Assisto o vídeo e logo em seguida tento fazer o notebook

Survey (1/12)

0 2 4

**A demographic dataset with statistics on different cities' population, GDP per capita, economic growth is an example of “unstructured” data because it contains data coming from different sources.**

False ✓



True



## What does a neuron compute?

A neuron computes an activation function followed by a linear function ( $z = Wx + b$ )

 13 %

A neuron computes a linear function ( $z = Wx + b$ ) followed by an activation function



 61 %

A neuron computes a function  $g$  that scales the input  $x$  linearly ( $Wx + b$ )

 17 %

A neuron computes the mean of all features before applying the output to an activation function

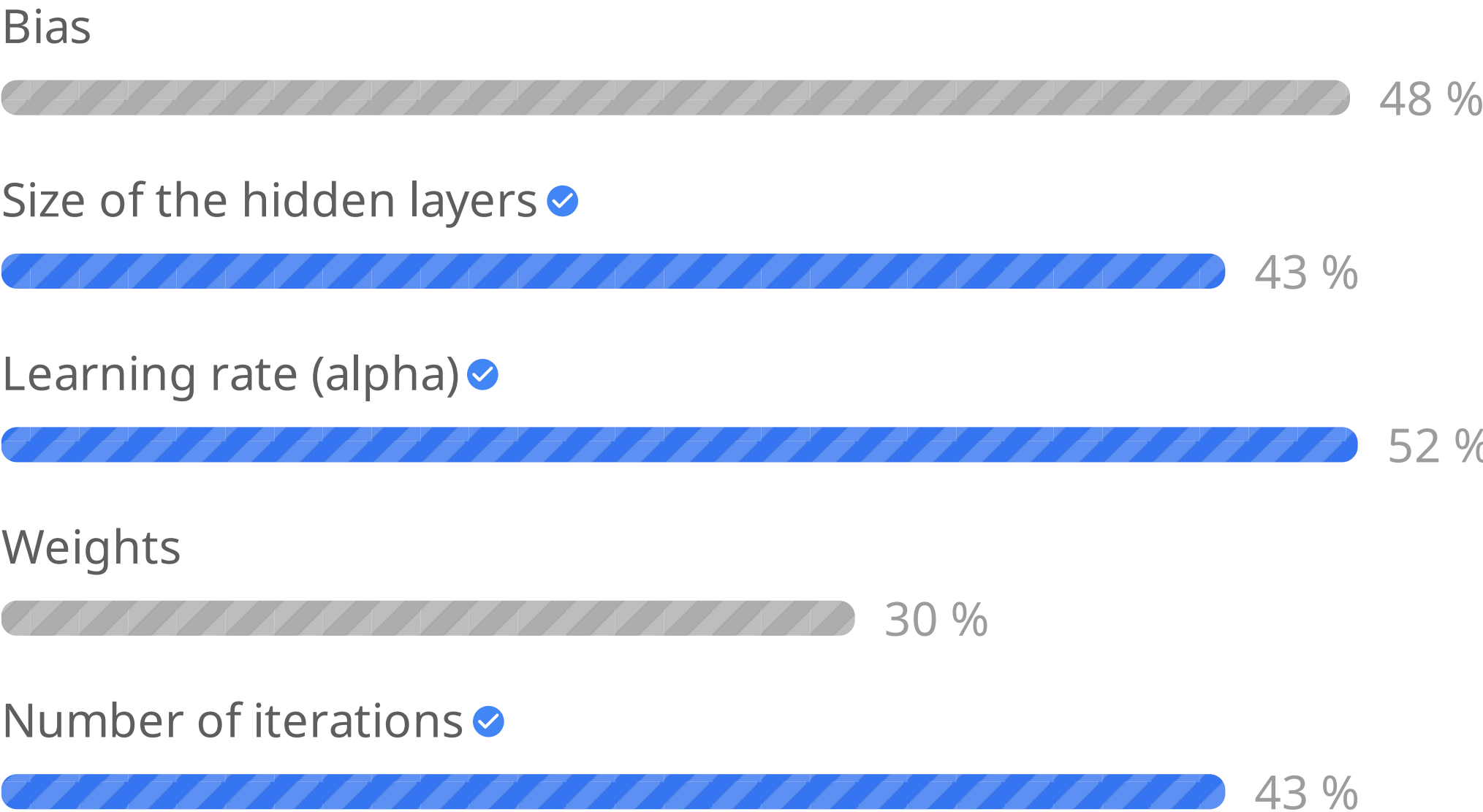
 9 %



Survey (3/12)

023

Among the following, which ones are "hyperparameters"?  
(1/2)



Survey (3/12)

0 2 3

**Among the following, which ones are "hyperparameters"?**  
(2/2)

Number of layers  $L$  in the neural network ✓



Activation values of  $g(z)$



Survey (4/12)

024

## Which of the following statements is true?

The deeper layers of a neural network are typically computing more complex features of the input than the earlier layers.



The earlier layers of a neural network are typically computing more complex features of the input than the deeper layers.



Survey (5/12)

0 2 4

**The tanh activation usually works better than sigmoid activation function for hidden units because the mean of its output is closer to zero, and so it centers the data better for the next layer. True/False?**

False



True ✓



Survey (6/12)

0 2 3

**You are building a binary classifier for recognizing "Tapioca" ( $y=1$ ) vs. "Ginga com Tapioca" ( $y=0$ ). Which one of these activation functions would you recommend using for the output layer?**

Relu



Leaky Relu



Sigmoid ✓



Tanh



Survey (7/12)

0 2 2

## Which of these statements about mini-batch gradient descent do you agree with?

You should implement mini-batch gradient descent without an explicit for-loop over different mini-batches, so that the algorithm processes all mini-batches at the same time (vectorization)

 23 %

Training one epoch (one pass through the training set) using mini-batch gradient descent is faster than training one epoch using batch gradient descent.

 27 %

One iteration of mini-batch gradient descent (computing on a single mini-batch) is faster than one iteration of batch gradient descent.



 50 %

Survey (8/12)

020

## Why is the best mini-batch size usually not 1 and not $m$ , but instead something in-between?

If the mini-batch size is  $m$ , you end up with batch gradient descent, which has to process the whole training set before making progress.



 65 %

If the mini-batch size is  $m$ , you end up with stochastic gradient descent, which is usually slower than mini-batch gradient descent.

 20 %

If the mini-batch size is 1, you lose the benefits of vectorization across examples in the mini-batch.



 40 %

If the mini-batch size is 1, you end up having to process the entire training set before making any progress.

 20 %

Survey (9/12)

0 1 5

**Suppose your learning algorithm's cost  $J$ , plotted as a function of the number of iterations, looks like this:**

(1/2)

Whether you're using batch gradient descent or mini-batch gradient descent, something is wrong.

 7 %

If you're using mini-batch gradient descent, this looks acceptable. But if you're using batch gradient descent, something is wrong.



 80 %

Whether you're using batch gradient descent or mini-batch gradient descent, this looks acceptable.

 7 %



Survey (9/12)

0 1 5

**Suppose your learning algorithm's cost  $J$ , plotted as a function of the number of iterations, looks like this:**

(2/2)

If you're using mini-batch gradient descent, something is wrong.  
But if you're using batch gradient descent, this looks acceptable.

 7 %

Survey (10/12)

0 1 6

**Consider this figure: these plots were generated with gradient descent; with gradient descent with momentum ( $\beta = 0.5$ ) and gradient descent with momentum ( $\beta = 0.9$ ). Which curve corresponds to which algorithm?**

(1/2)

(1) is gradient descent with momentum (small  $\beta$ ), (2) is gradient descent with momentum (small  $\beta$ ), (3) is gradient descent.

 19 %

(1) is gradient descent. (2) is gradient descent with momentum (large  $\beta$ ) . (3) is gradient descent with momentum (small  $\beta$ )

 6 %

(1) is gradient descent with momentum (small  $\beta$ ). (2) is gradient descent. (3) is gradient descent with momentum (large  $\beta$ )

 6 %

Survey (10/12)

0 1 6

**Consider this figure: these plots were generated with gradient descent; with gradient descent with momentum ( $\beta = 0.5$ ) and gradient descent with momentum ( $\beta = 0.9$ ). Which curve corresponds to which algorithm?**

(2/2)

(1) is gradient descent. (2) is gradient descent with momentum (small  $\beta$ ). (3) is gradient descent with momentum (large  $\beta$ )



69 %

Survey (11/12)

016

**Suppose batch gradient descent in a deep network is taking excessively long to find a value of the parameters that achieves a small value for the cost function  $J$ . Which of the following techniques could help:**

(1/2)

Try mini-batch gradient descent. ✓



Try initializing all the weights to zero.



Try better random initialization for the weights. ✓



Try tuning the learning rate  $\alpha$  (alpha). ✓



Survey (11/12)

016

**Suppose batch gradient descent in a deep network is taking excessively long to find a value of the parameters that achieves a small value for the cost function  $J$ . Which of the following techniques could help:**

(2/2)

Try using Adam. ✓



Survey (12/12)

018

## Which of the following statements about Adam is False?

Adam combines the advantages of RMSProp and momentum.

 6 %

We usually use “default” values for the hyperparameters  $\beta_1$ ,  $\beta_2$  and  $\epsilon$  in Adam ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $10^{-8}$  )

 6 %

The learning rate hyperparameter  $\alpha$  in Adam usually needs to be tuned.

 28 %

Adam should be used with batch gradient computations, not with mini-batches.



 61 %