

Deep-learning Neural Network

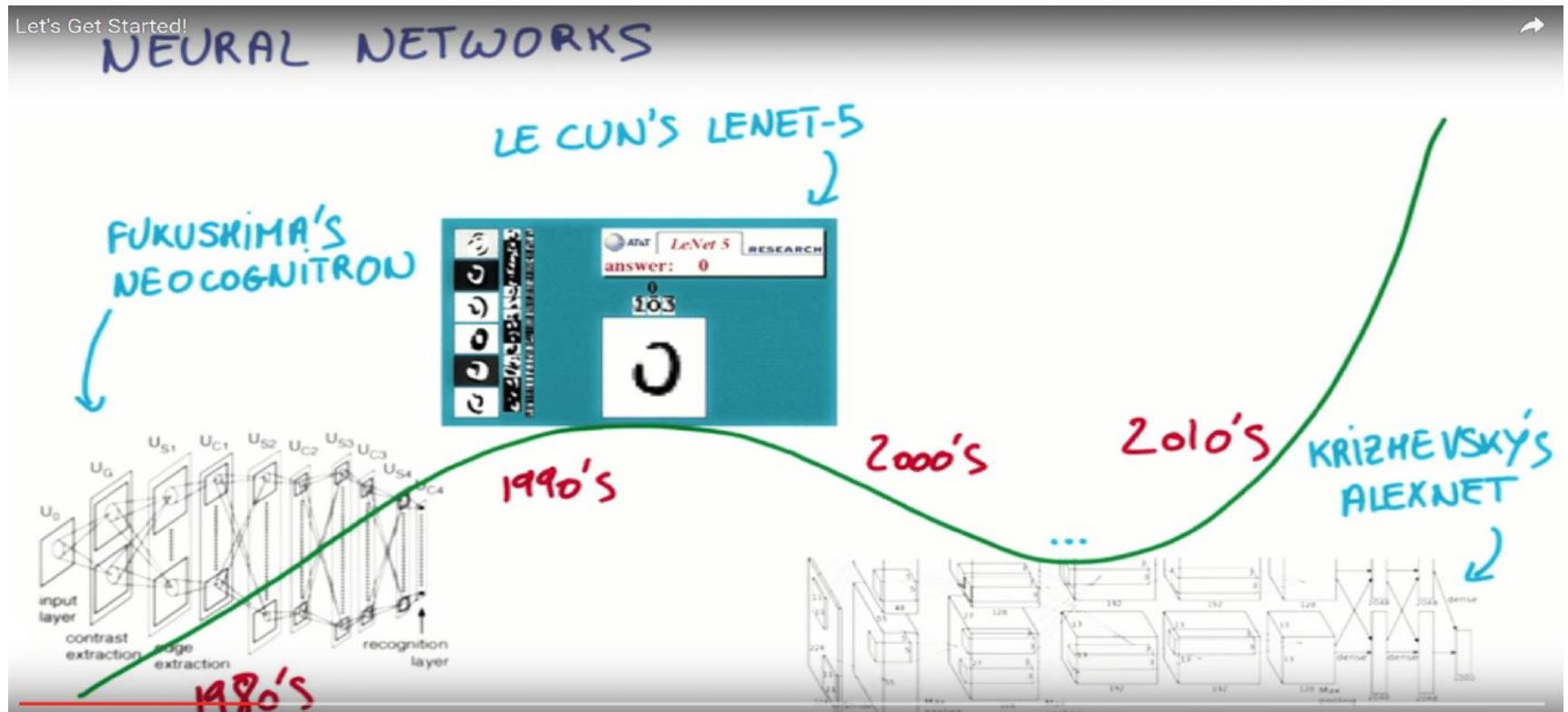
TensorFlow™

MNIST example

**Scientific application:
Higgs CP measurement at LHC**

Neural network

**Since 2010 new era in Machine Learning:
rapidly increasing areas of applications**



Neural network

Since 2010 new era: rapidly increasing areas of applications

NEURAL NETWORKS → DEEP LEARNING

- 2009: SPEECH RECOGNITION
- 2012: COMPUTER VISION
- 2014: MACHINE TRANSLATION

① DATA!



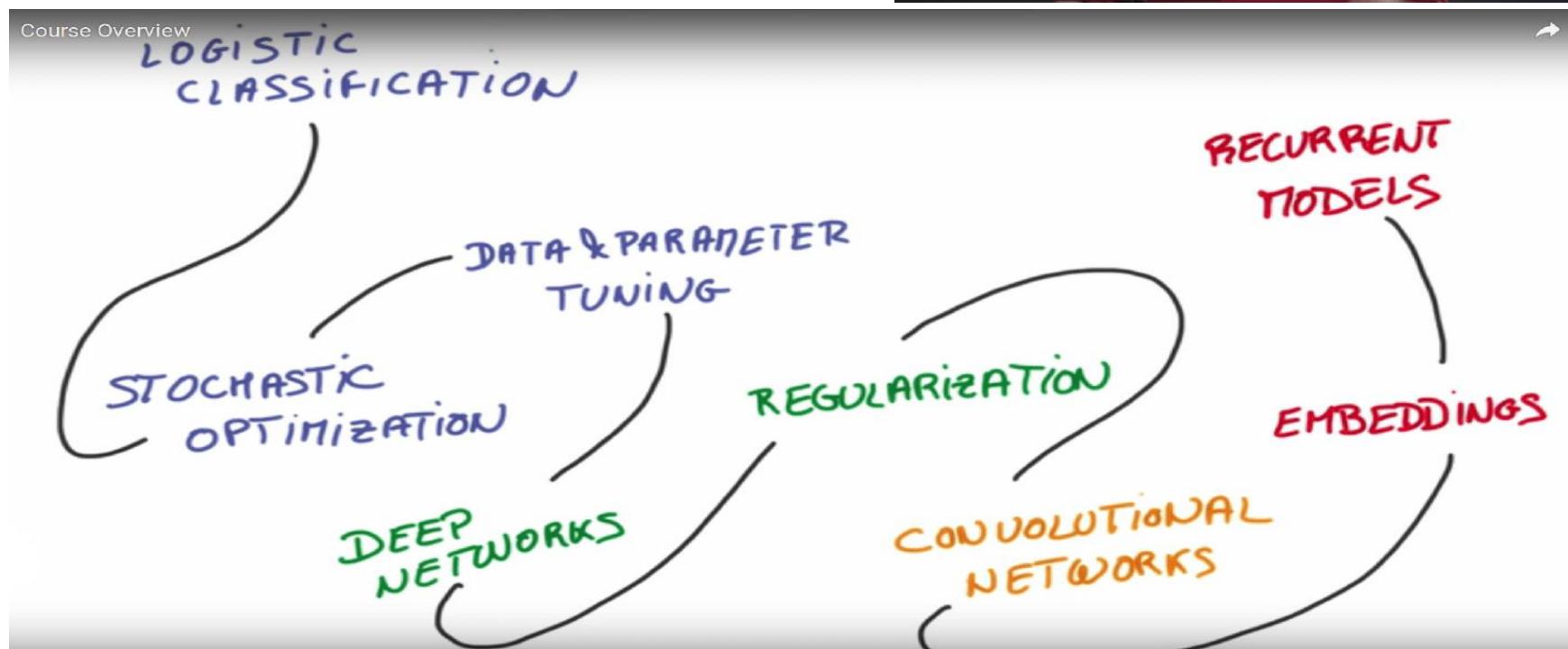
② GPUs!



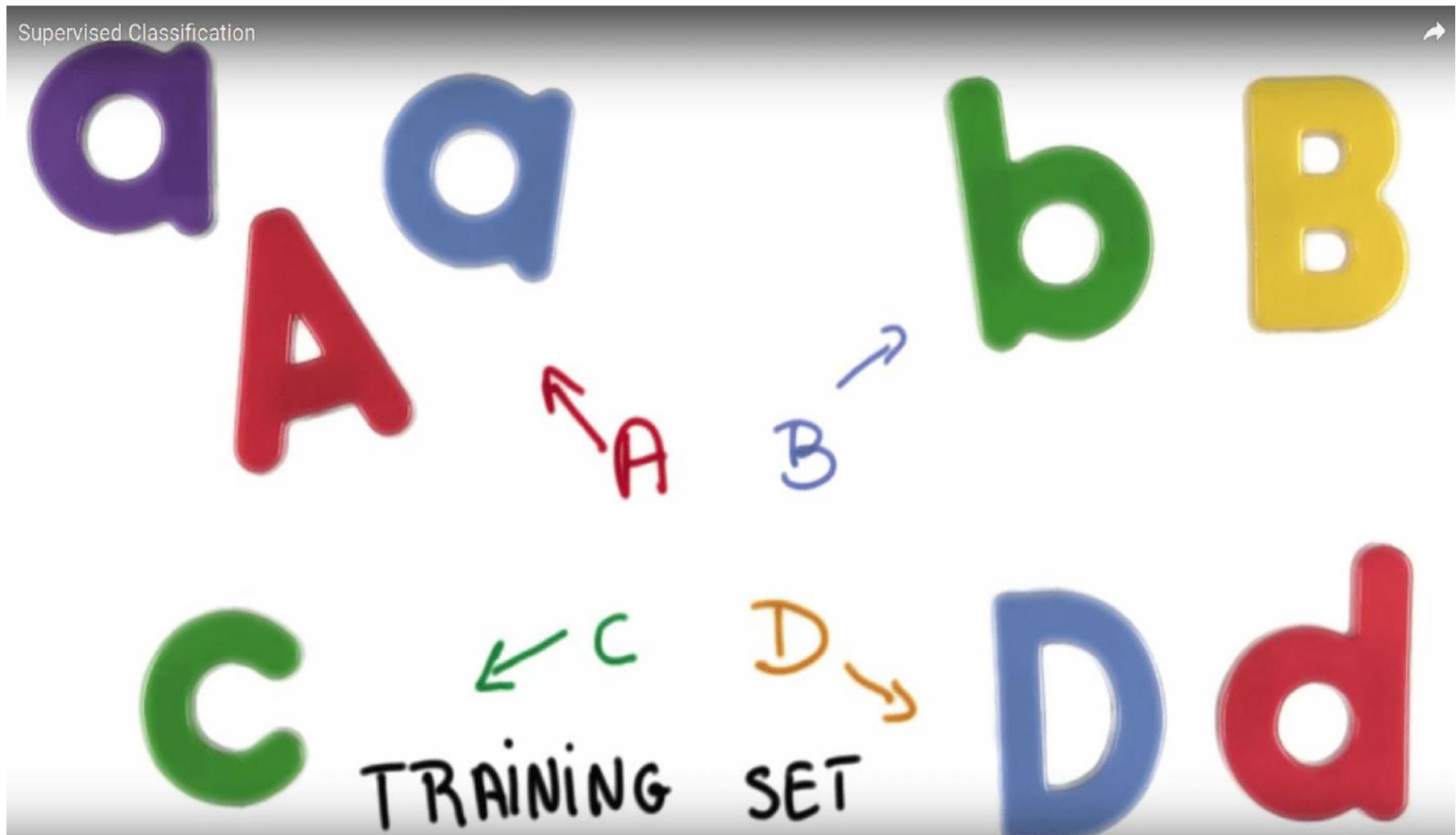
THANK YOU
GAMERS!!!

Deep-Learning tutorial @ udacity

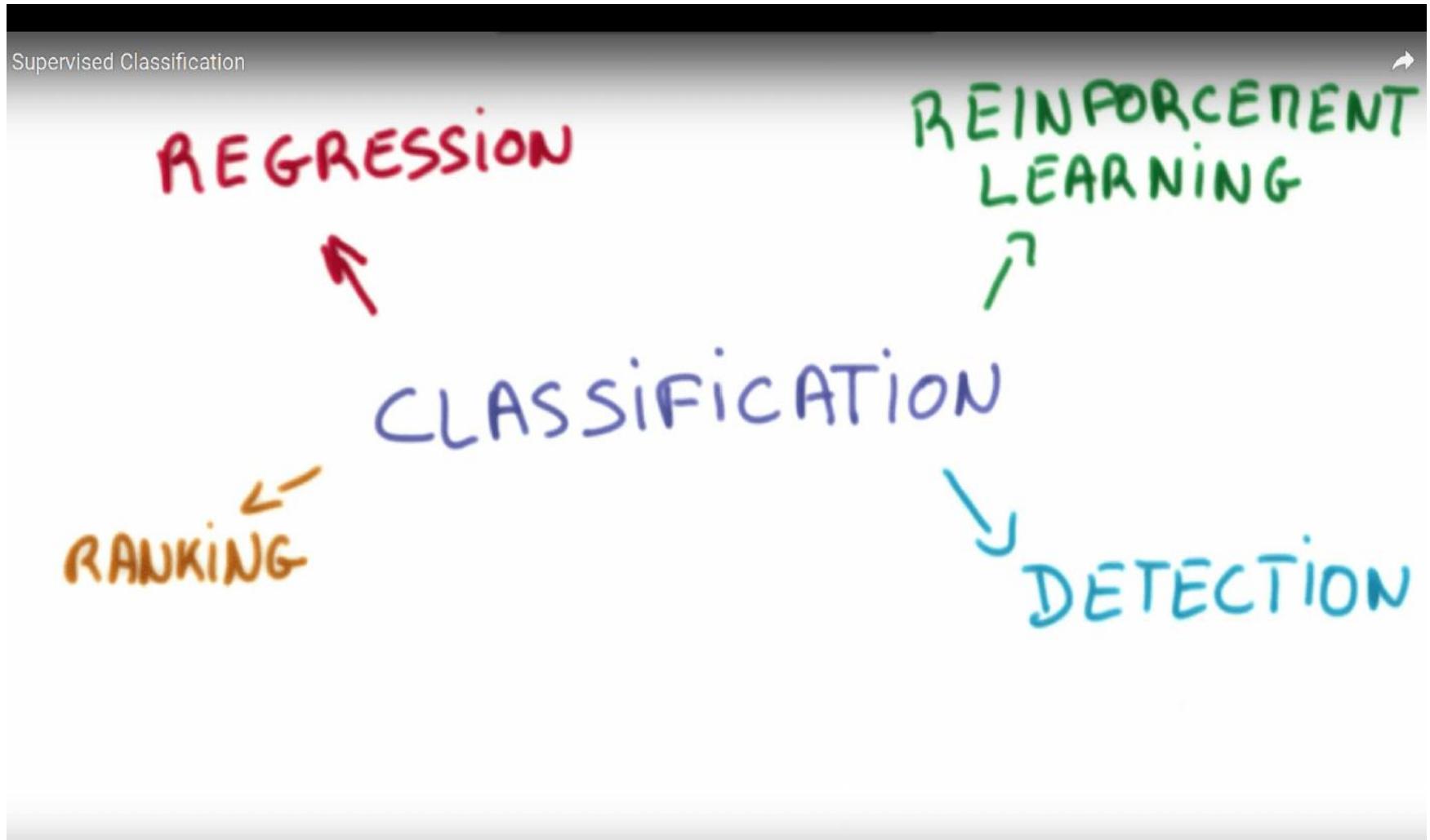
<https://www.udacity.com/course/deep-learning--ud730>



Supervised Classifications



Supervised Classifications

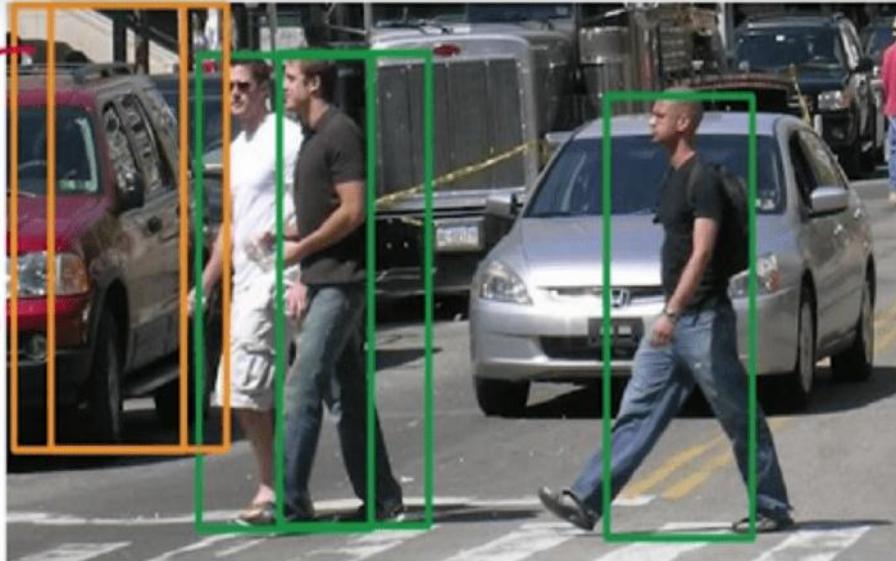


Classifications for Detection

Classification For Detection Solution

DETECTION

PEDESTRIAN
VS.
NO PEDESTRIAN



Use a binary pedestrian / no pedestrian classifier. Slide it over all possible locations in the image.

Classifications for Ranking

The screenshot shows a search results page from a search engine. The query "Udacity deep learning" is entered in the search bar. Below the search bar, there are links for "Web", "News", "Videos", "Images", "Shopping", "More", and "Search tools". A green oval highlights the first result, which is a link to "Udacity deep learning". A large blue arrow labeled "RANKING" points from the top left towards this result. Below the result, it says "About 46,700 results (0.64 seconds)". The second result is highlighted with a green oval and is titled "Machine Learning: Supervised Learning - Udacity". It includes the URL <https://www.udacity.com/wiki/ud675>, the source "Udacity", and the date "Oct 3, 2014 - Ethem Alpaydin, Introduction to Machine Learning, Second Edition. ...".

CLASSIFIER

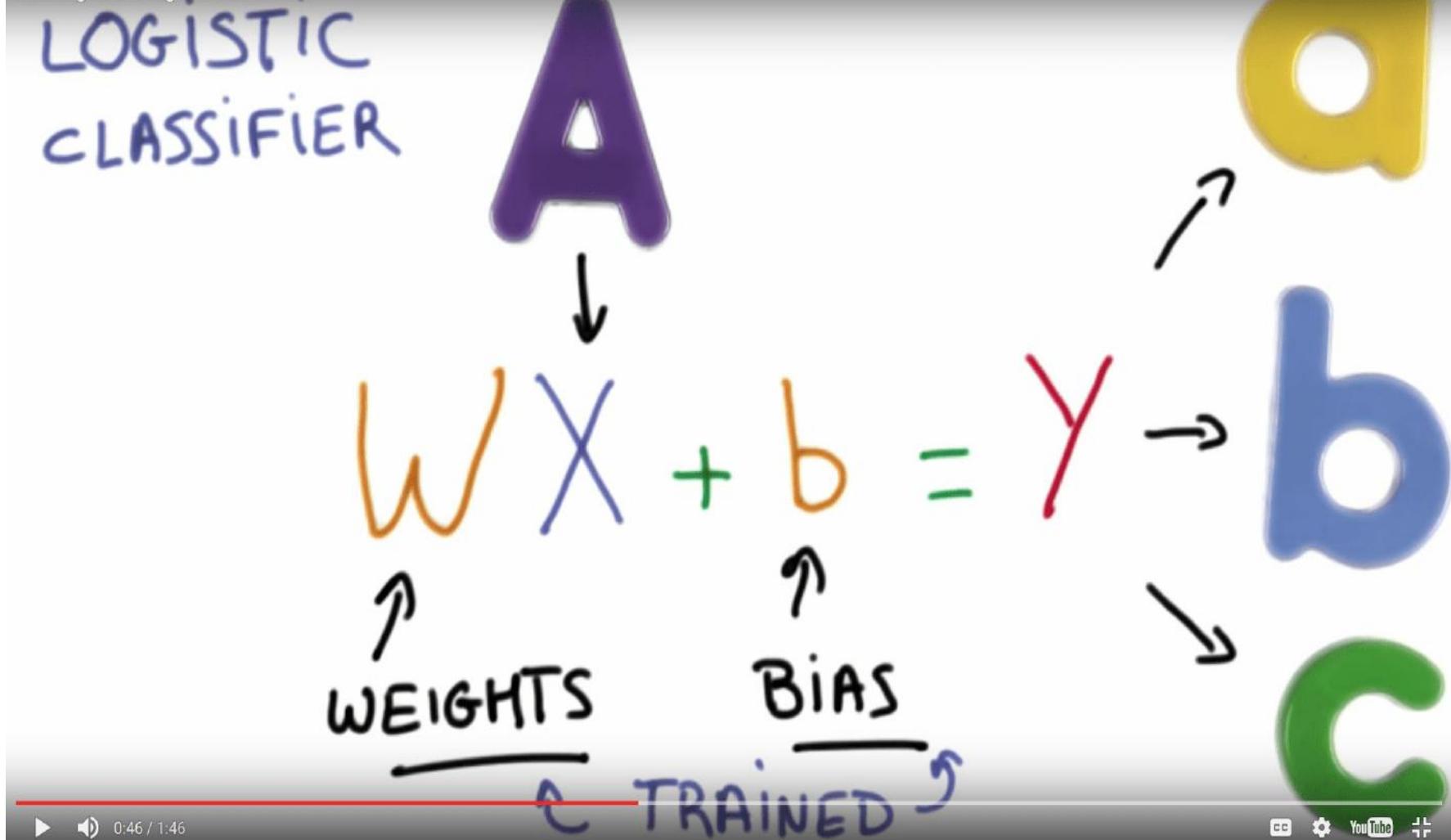
↓

RELEVANT!

Use a classifier that takes the pair
(query, web page) as an input, and
classifies the pair as relevant / not relevant

Logistic classifier: Linear model

Training Your Logistic Classifier



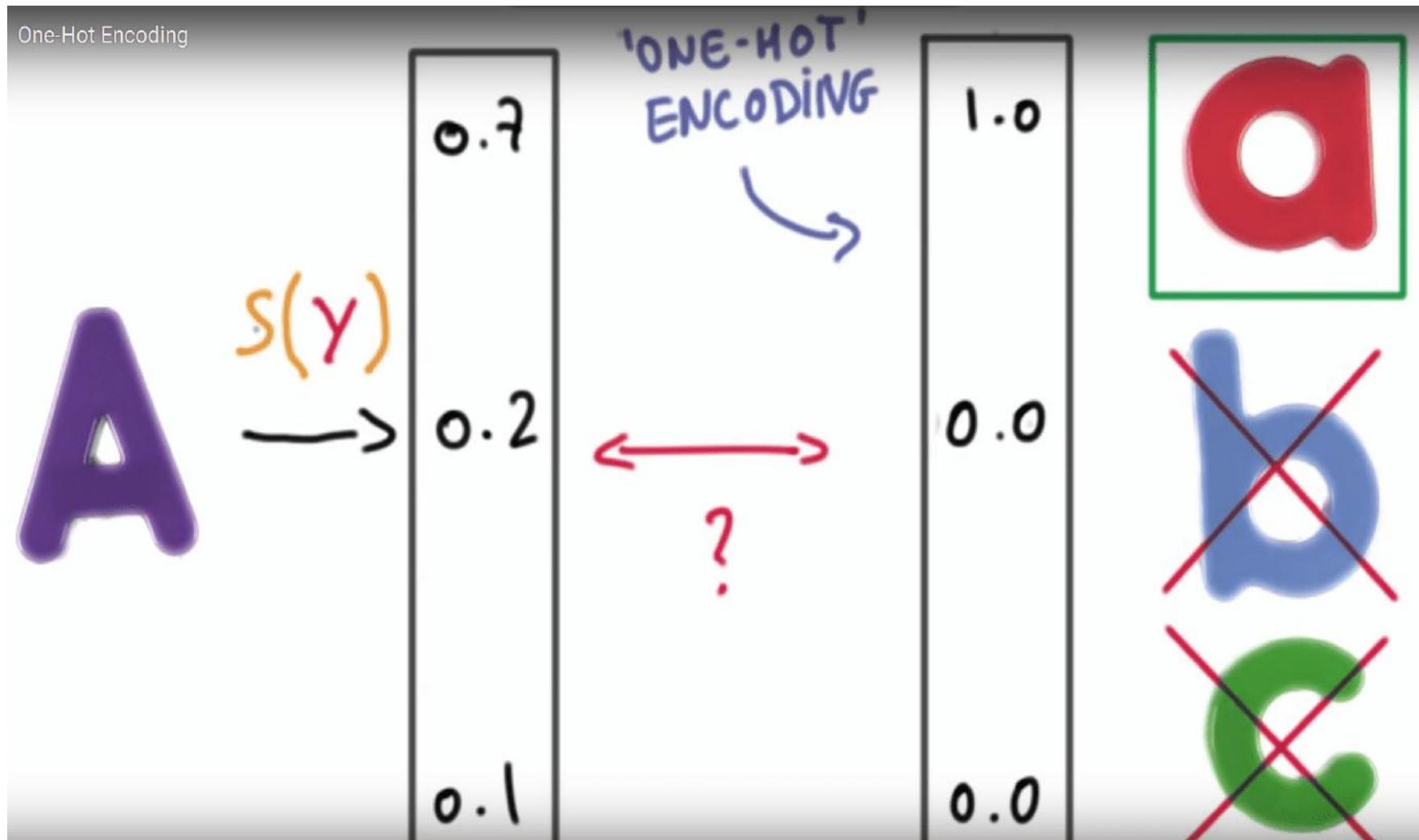
Softmax

Training Your Logistic Classifier

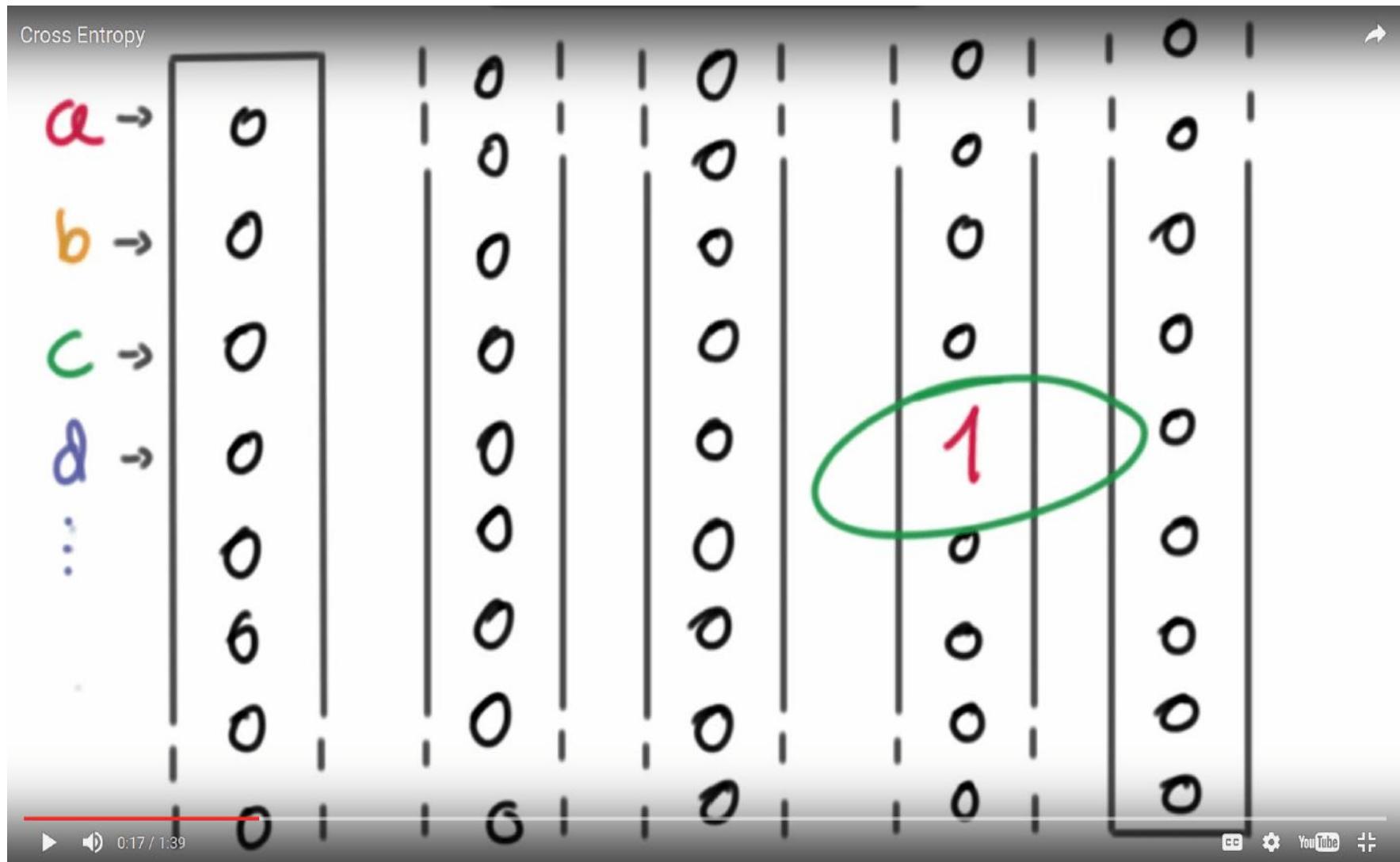
SOFTMAX



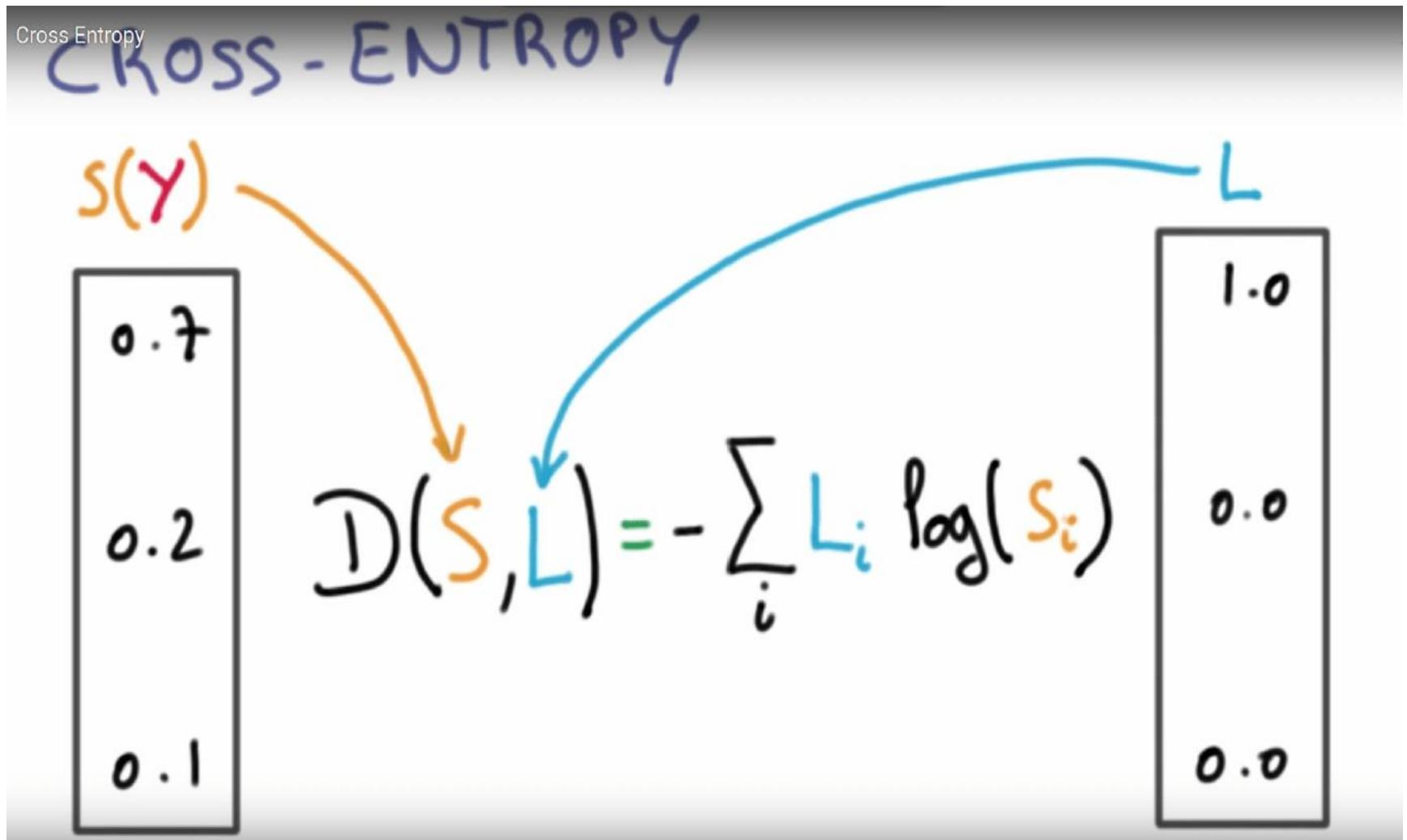
„One hot“ encoding



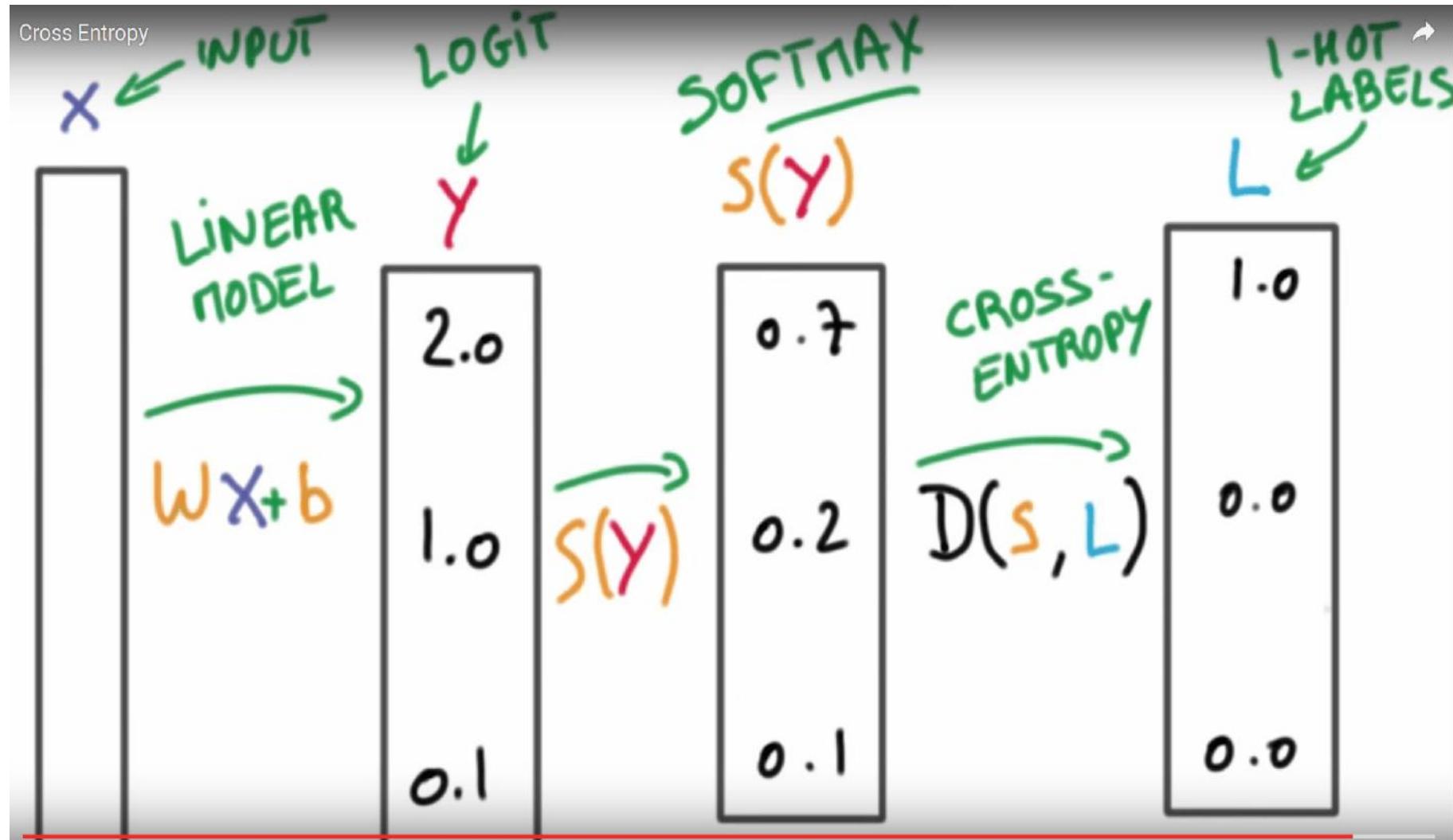
„One hot“ encoding



Optimisation: Cross-Entropy



Multinomial logistic classification



Optimisation of average loss

Minimizing Cross Entropy

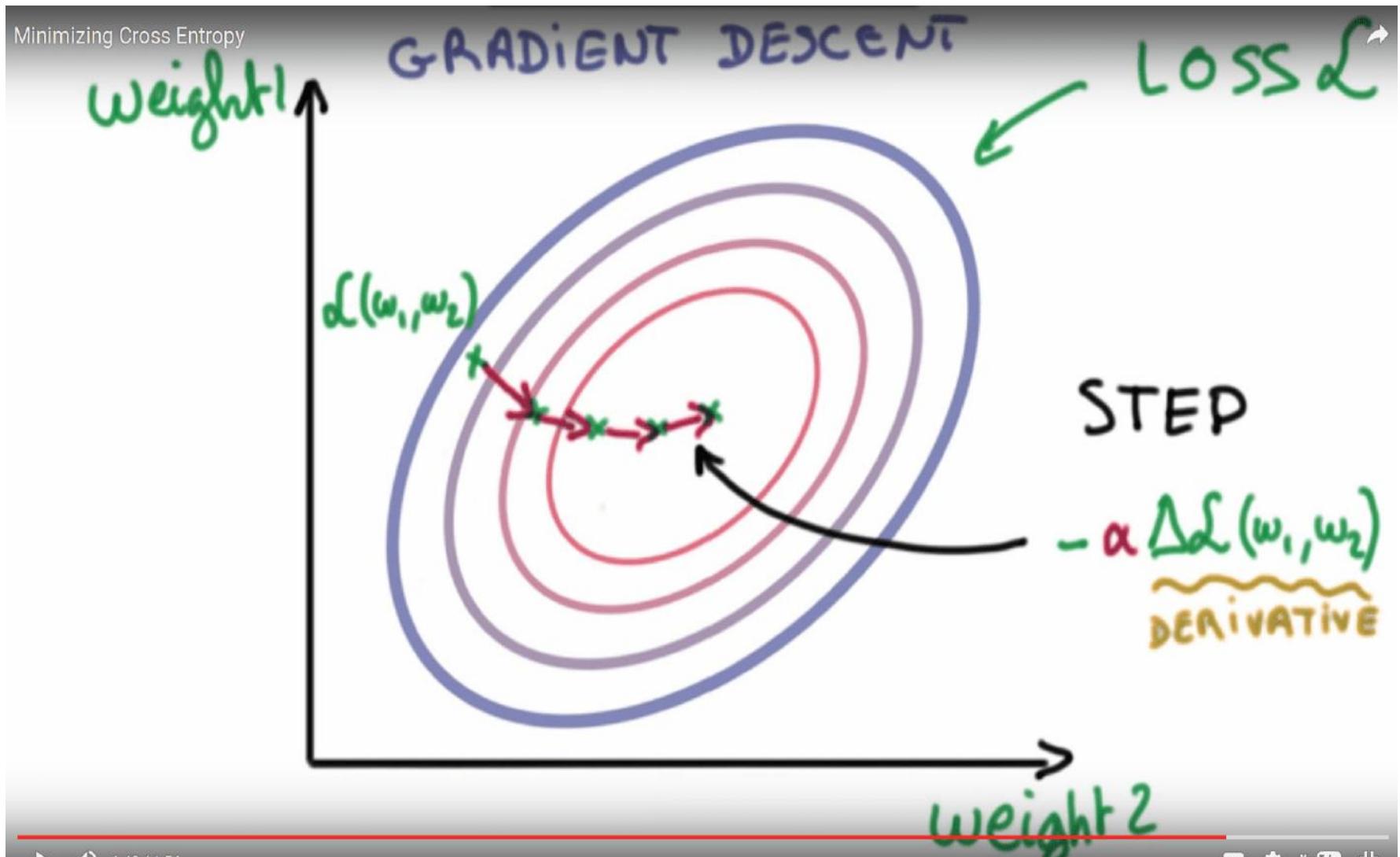
LOSS = AVERAGE CROSS-ENTROPY

$$\mathcal{L} = \frac{1}{N} \sum_i D(s(wx_i + b), L_i)$$

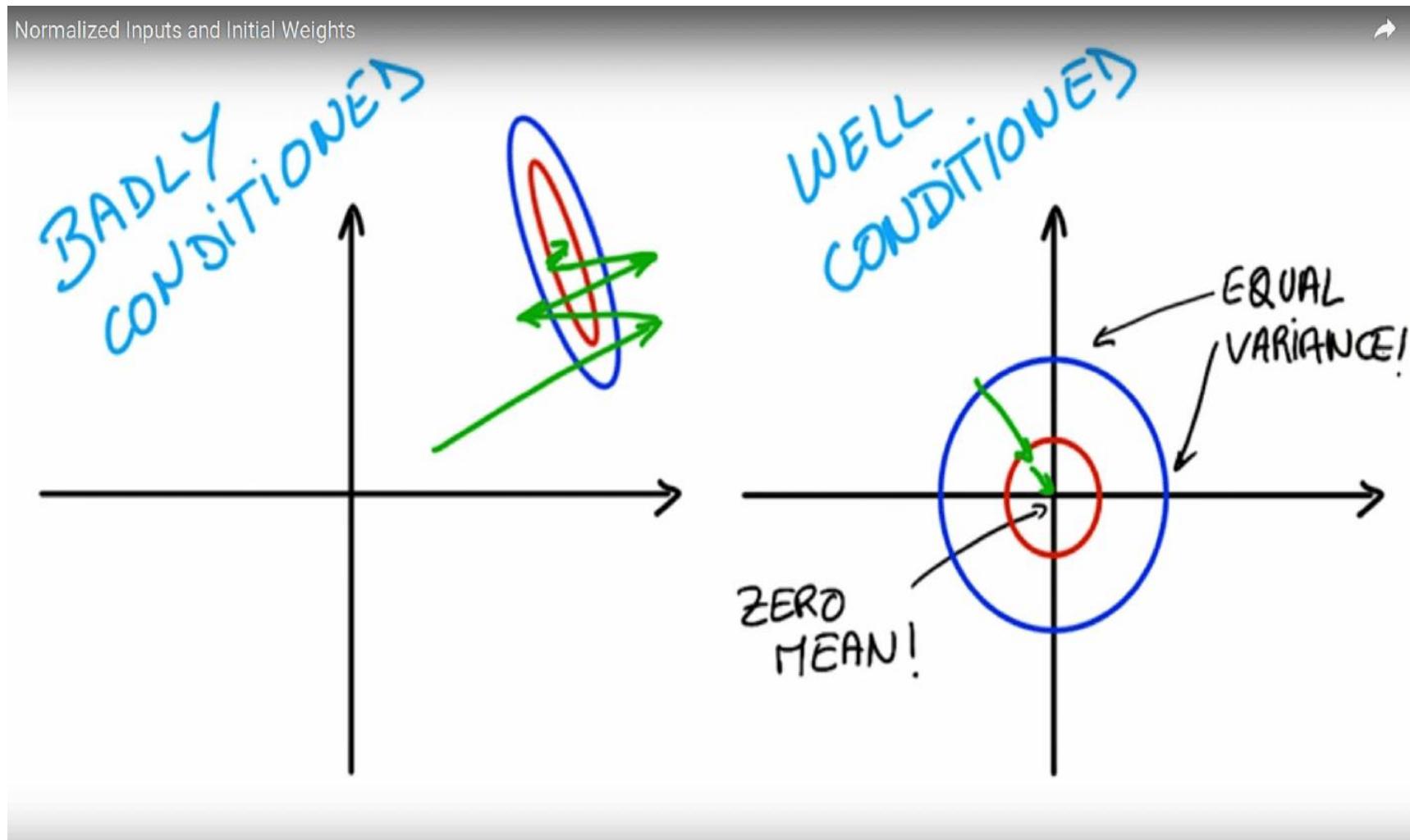
BIG MATRIX!!

BIG SUM!!

Gradient decent



Normalised input and output



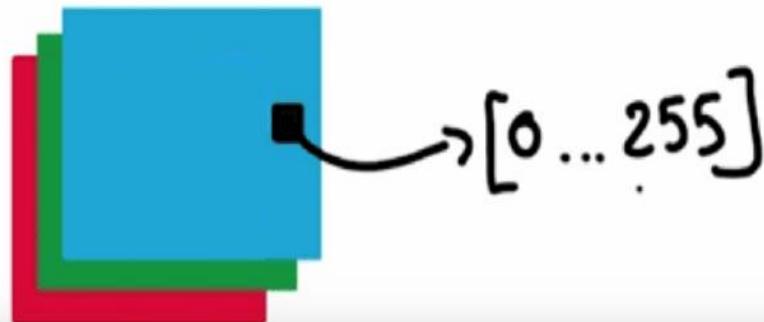
Normalised input and output

IMAGES

$$\frac{R - 128}{128}$$

$$\frac{G - 128}{128}$$

$$\frac{B - 128}{128}$$



Initialisation

Normalized Inputs and Initial Weights

INITIALISATION
OF THE LOGISTIC
CLASSIFIER

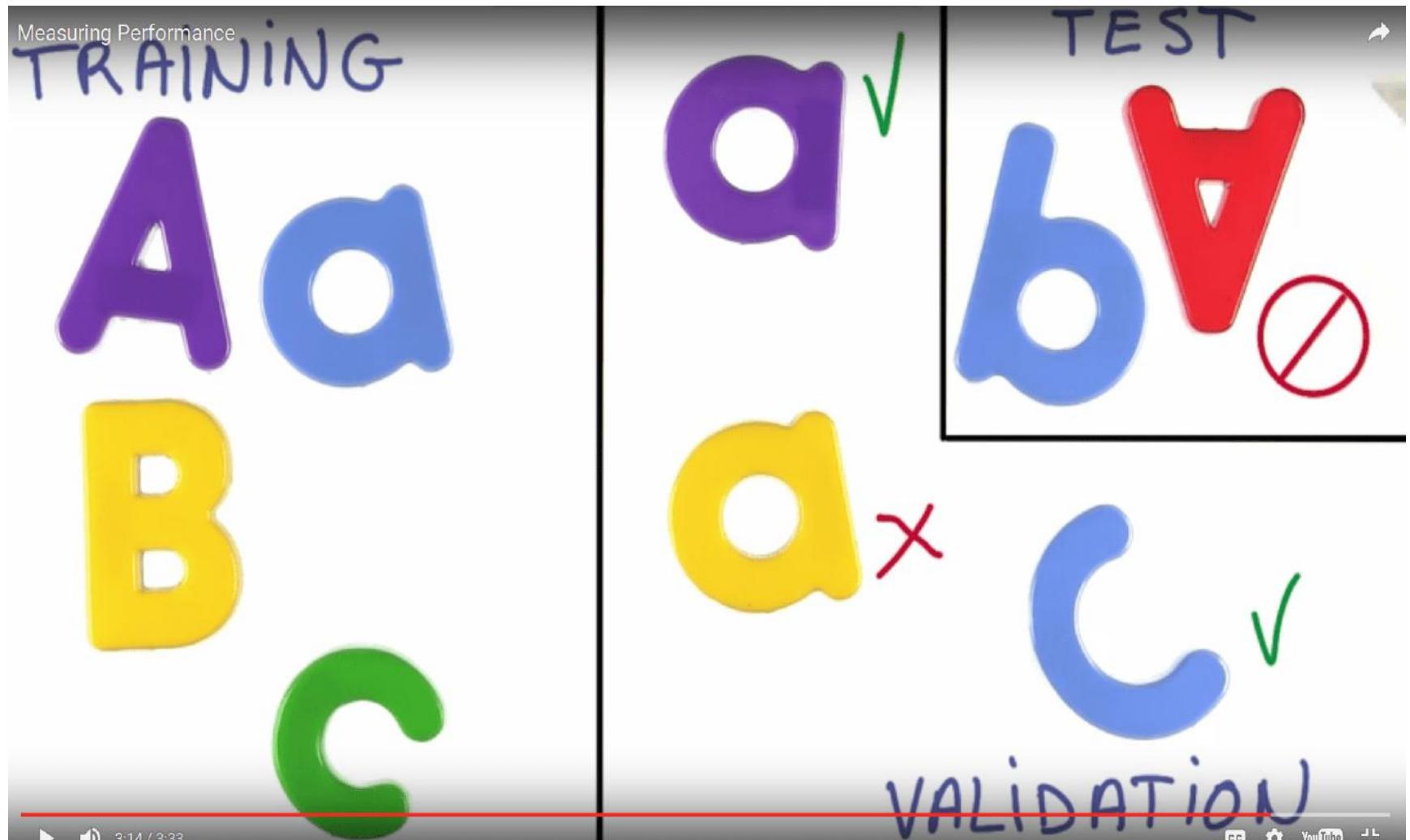
$$\frac{\text{pixels} - 128}{128}$$

$$\mathcal{L} = \frac{1}{N} \sum_i \mathcal{D}(S(\omega x_i + b), L_i)$$

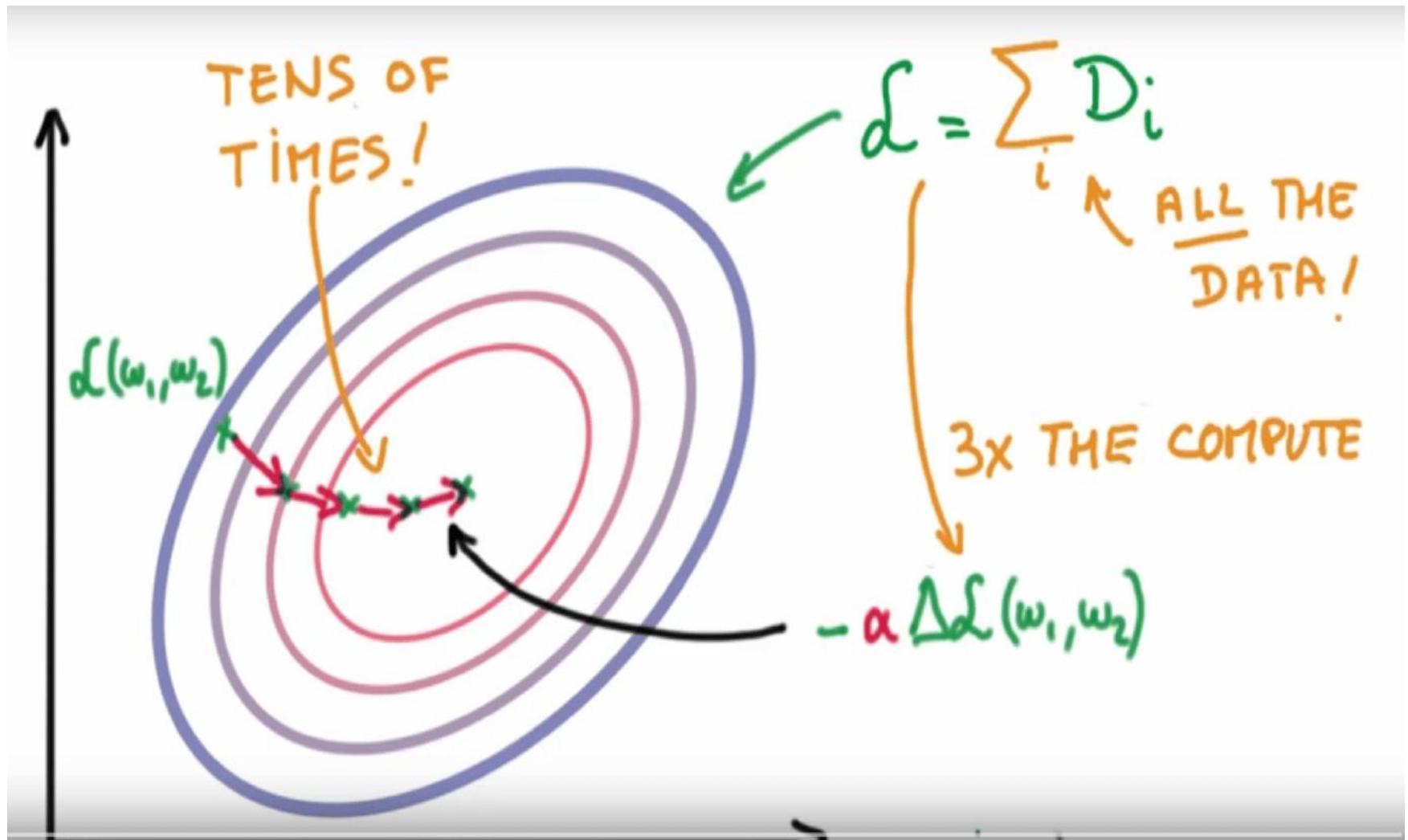
$$\phi$$



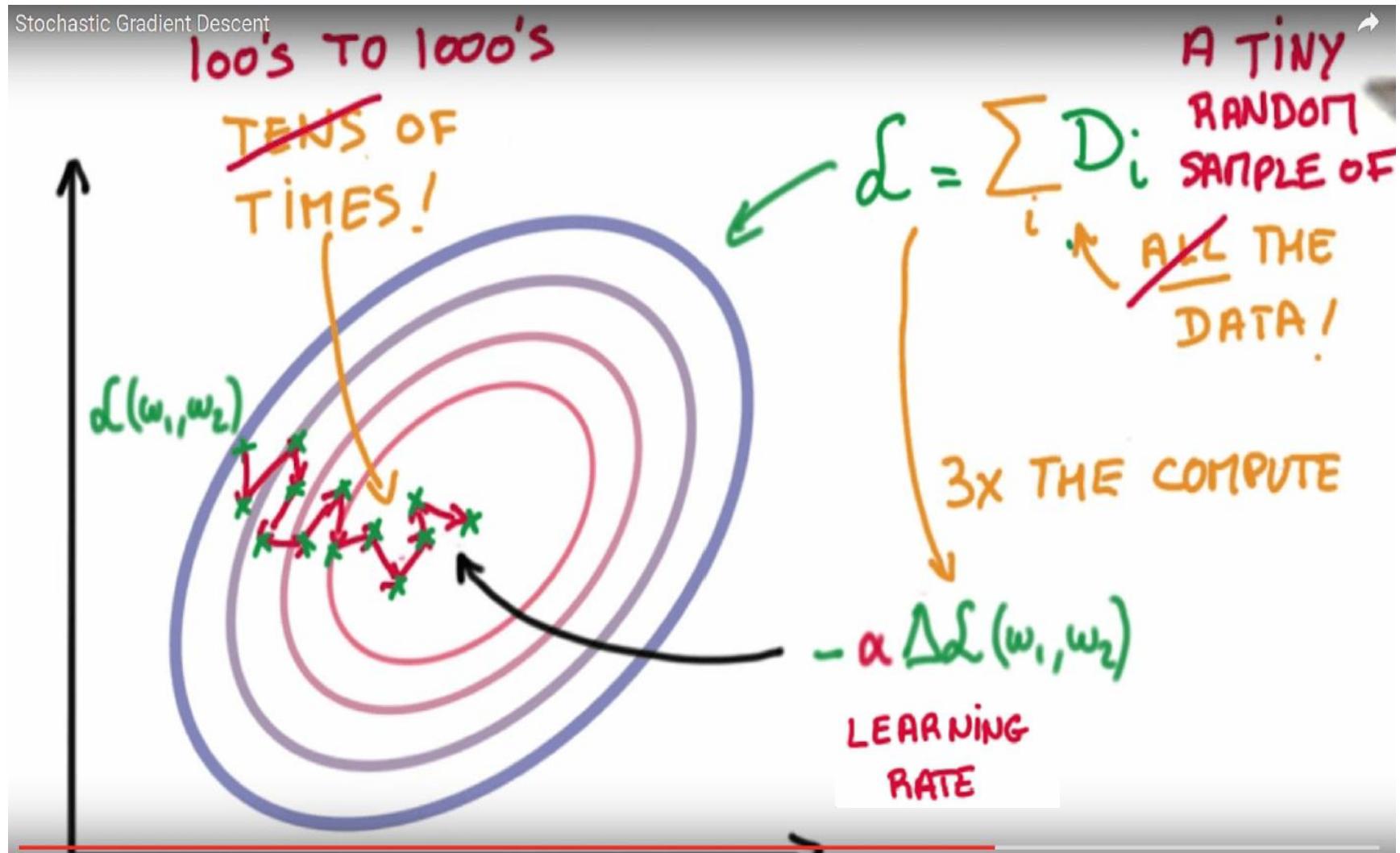
Training, validation, testing



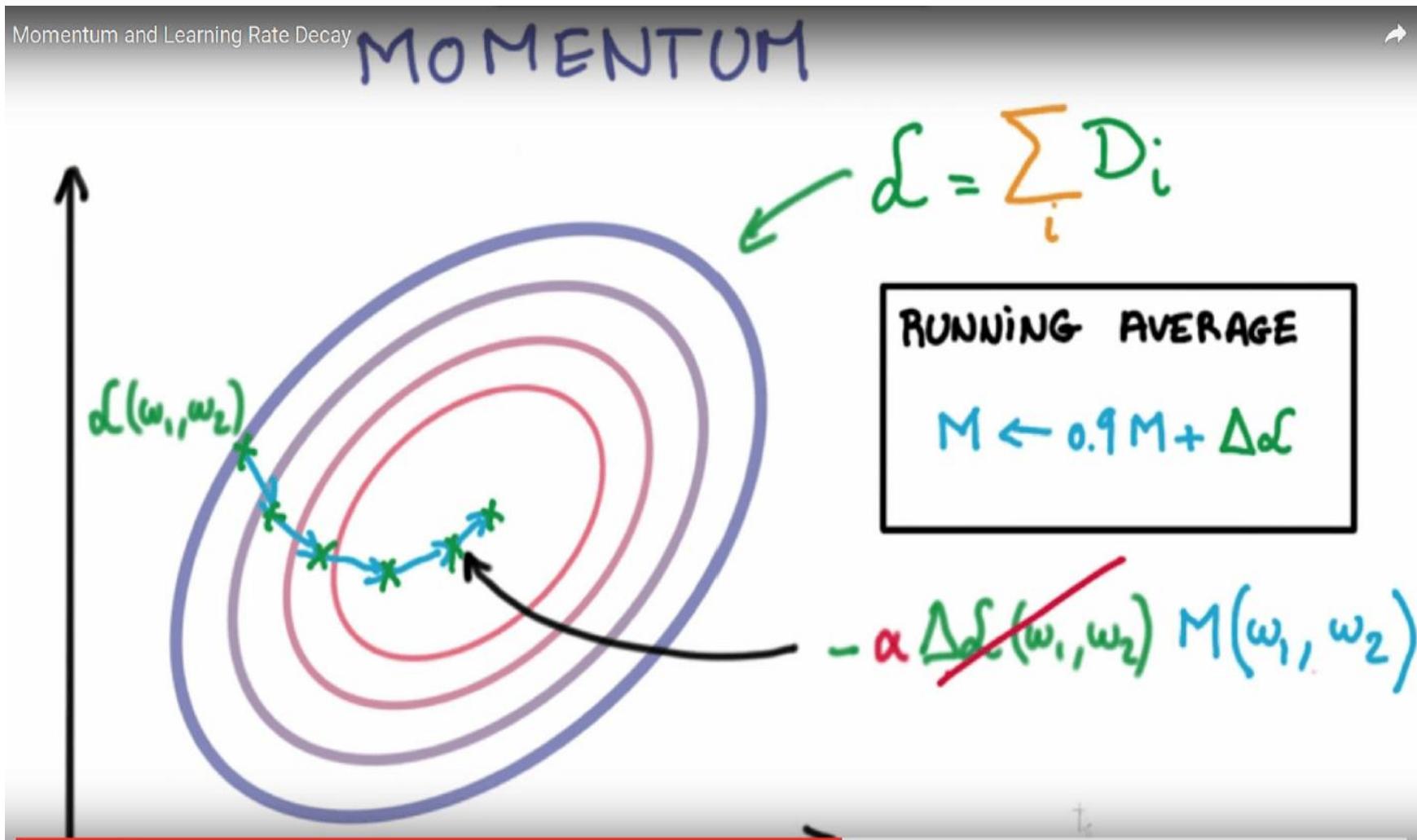
Gradient Descent



Stochastic Gradient Descent



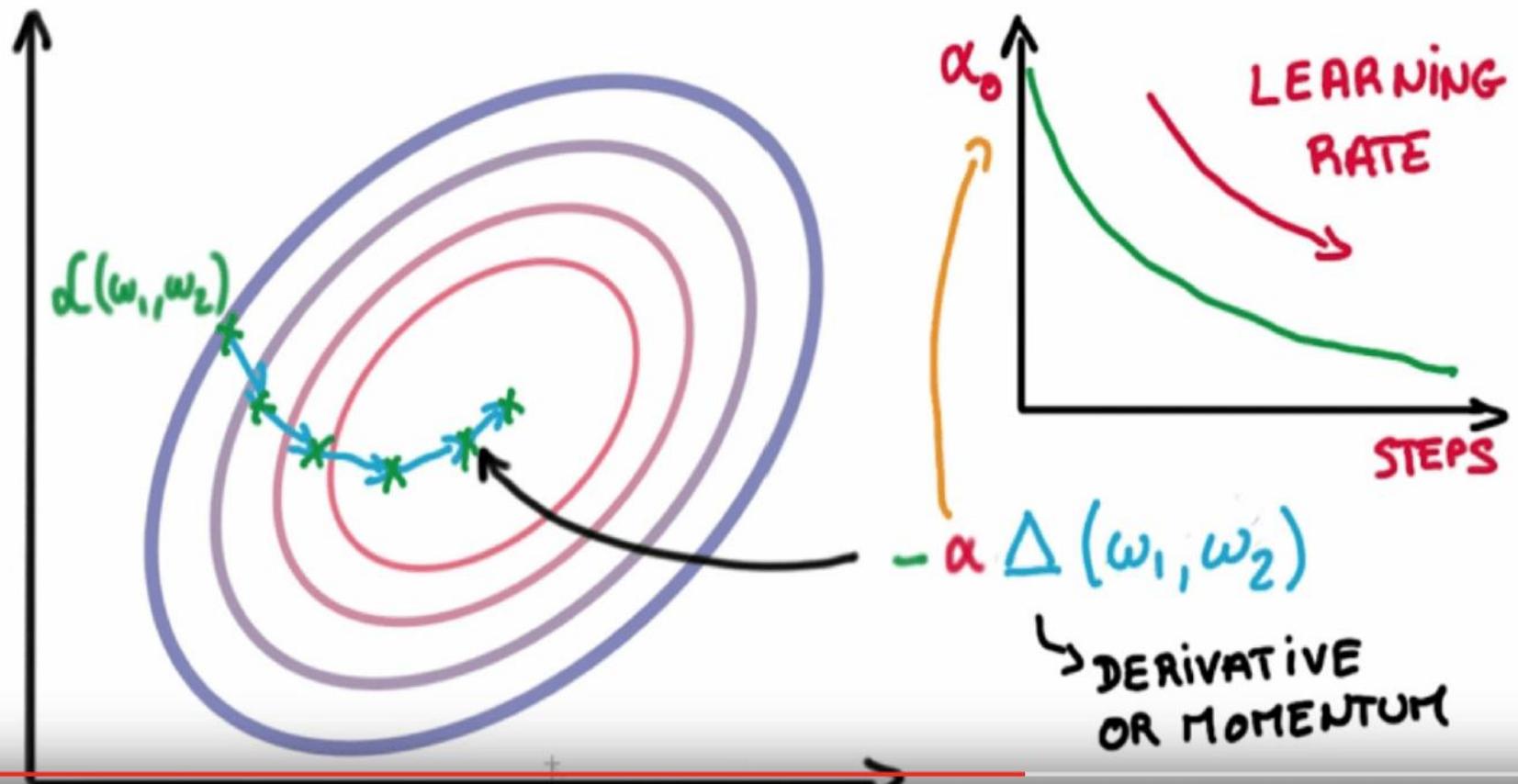
SDG: optimising with momentum



SDG: learning rate

Momentum and Learning Rate Decay

LEARNING RATE DECAY

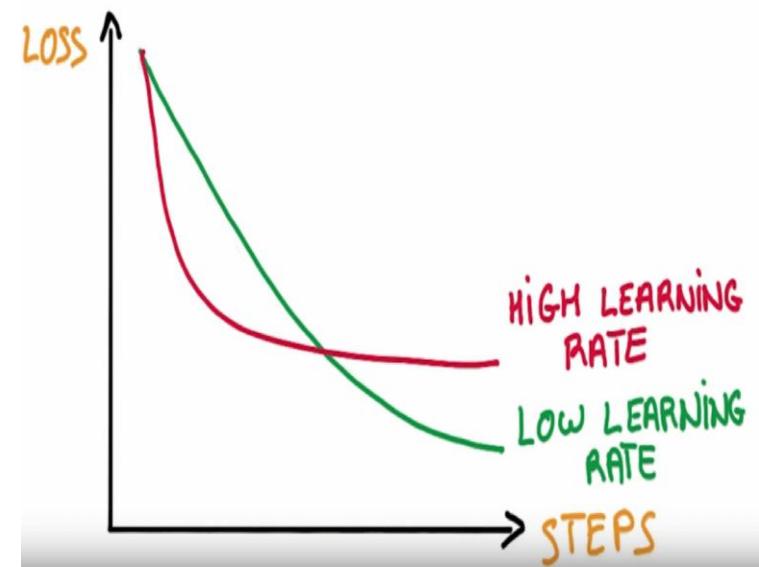


SDG: „black magic”

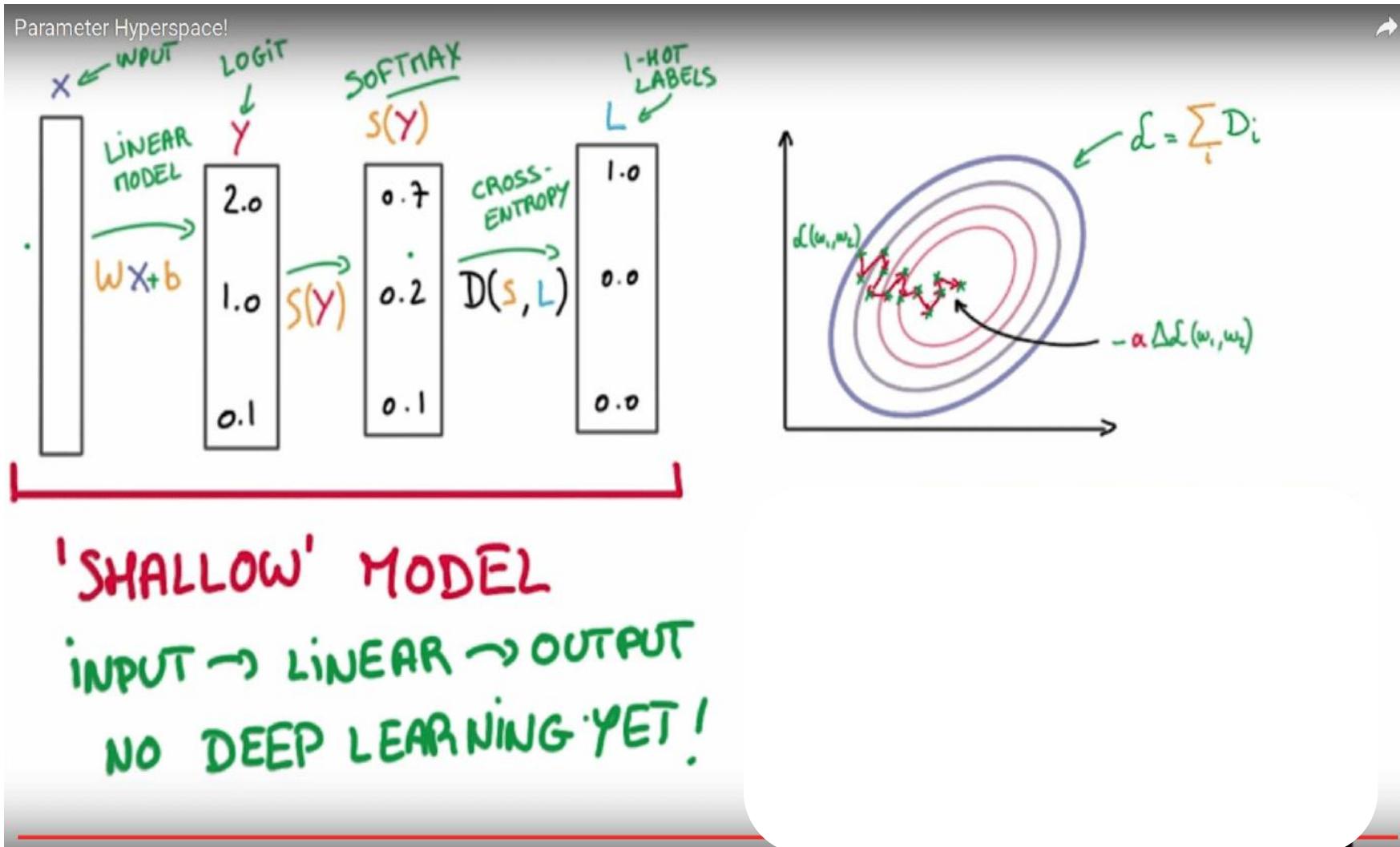
MANY HYPER-PARAMETERS

- INITIAL LEARNING RATE
- LEARNING RATE DECAY
- MOMENTUM
- BATCH SIZE
- WEIGHT INITIALIZATION

↑



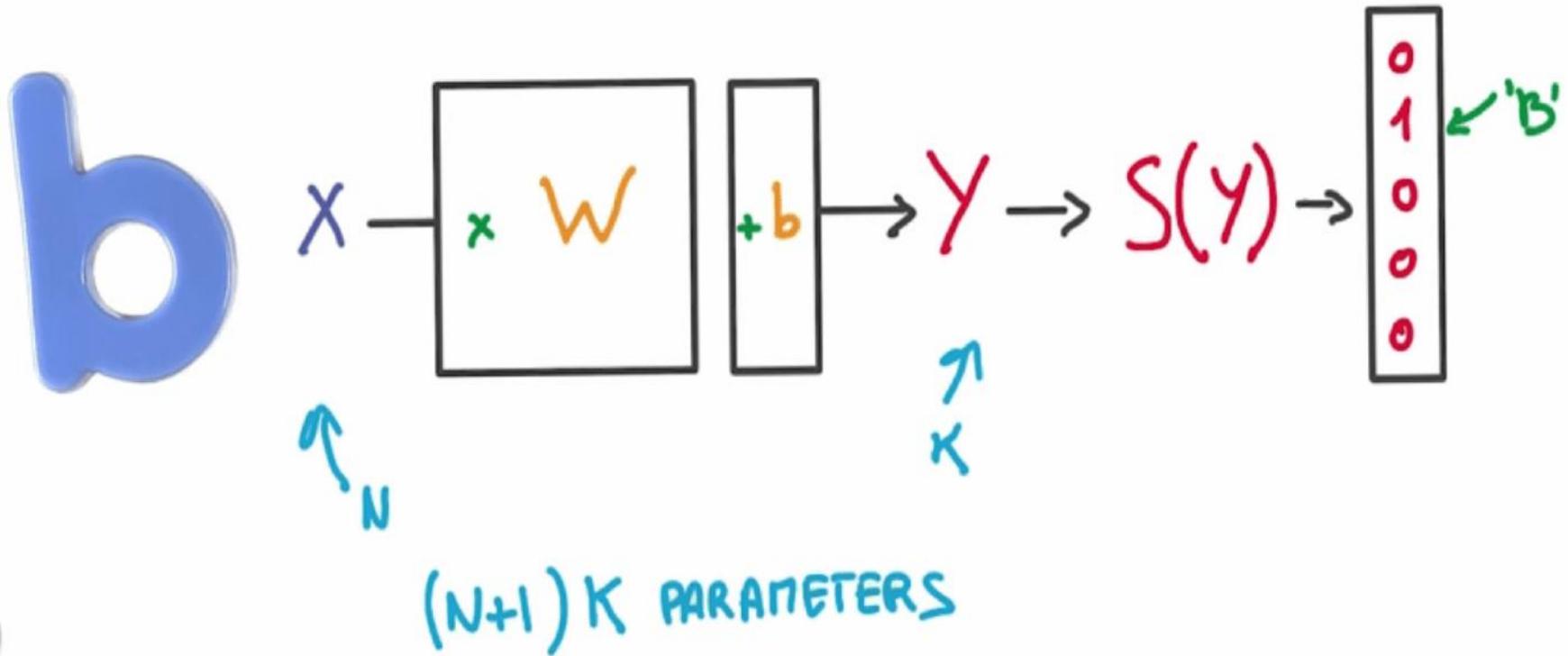
Input – linear - output



Linear models are linear

Linear Models are Limited

LINEAR MODEL COMPLEXITY



Linear models are stable

$$Y = W X \rightarrow \Delta Y \sim |W| \Delta X$$

↑
SMALL BOUNDED
↑
SMALL

$$\frac{\Delta Y}{\Delta X} = W^T$$

CONSTANTS

$$\frac{\Delta Y}{\Delta W} = X^T$$

Linear models are here to stay

- This is still linear

$$Y = \omega_1 \omega_2 \omega_3 X = w X$$

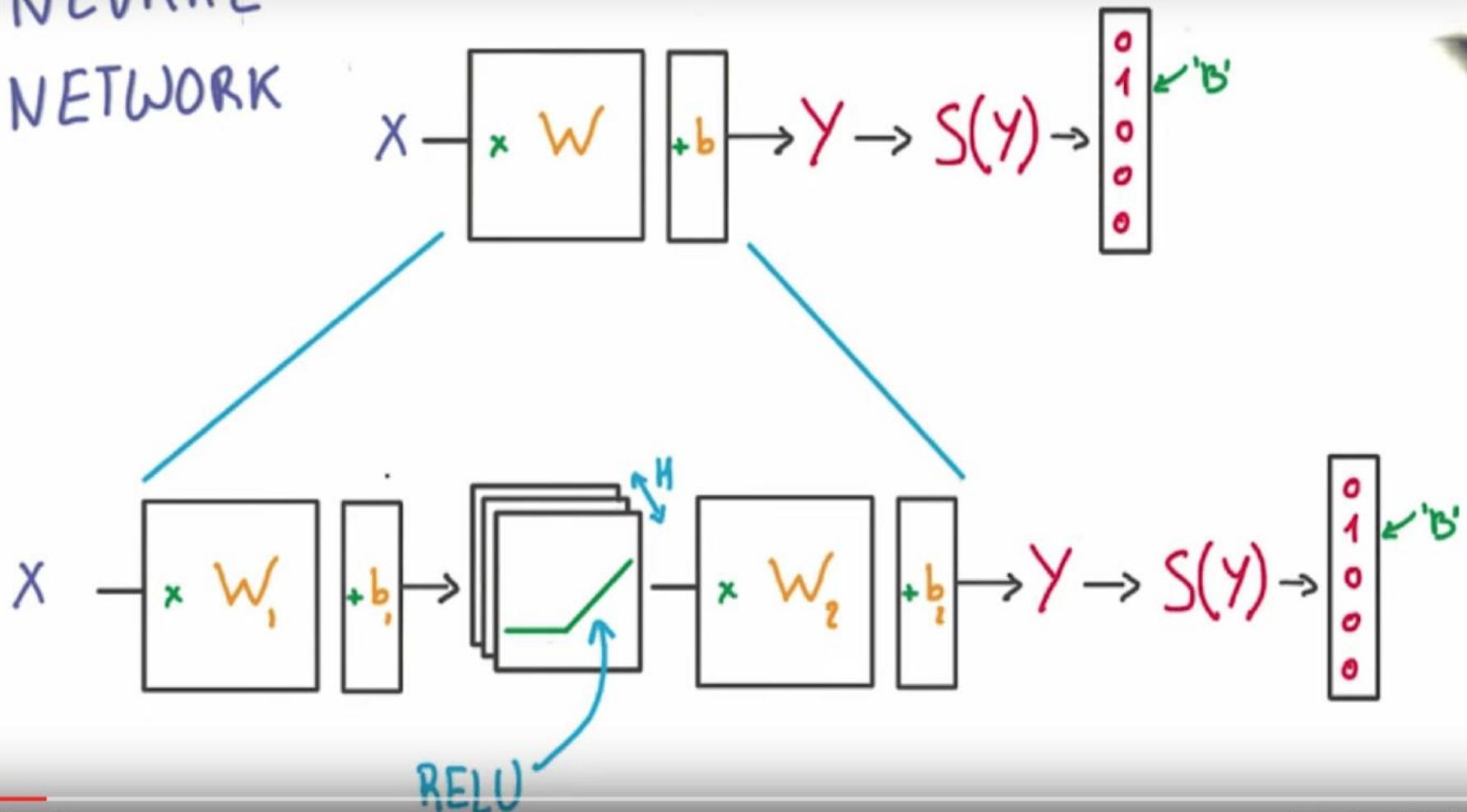
- Lets introduce non-linearity

$$Y = \omega_1 \omega_2 \omega_3 X = \textcolor{brown}{w} X$$

↑ ↑
NON - LINEARITIES

RELU: Rectified Linear Unit

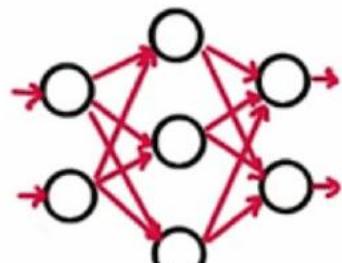
NEURAL
NETWORK



Networks of RELU

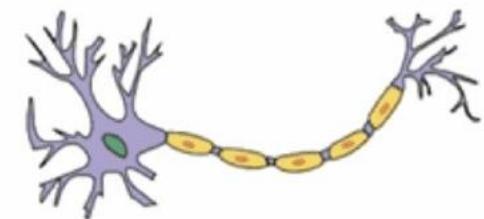
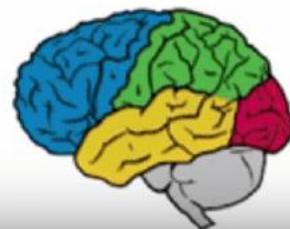
Network of ReLUs

NEURAL NETWORK



NEURONS?

HOW THE
BRAIN WORKS?



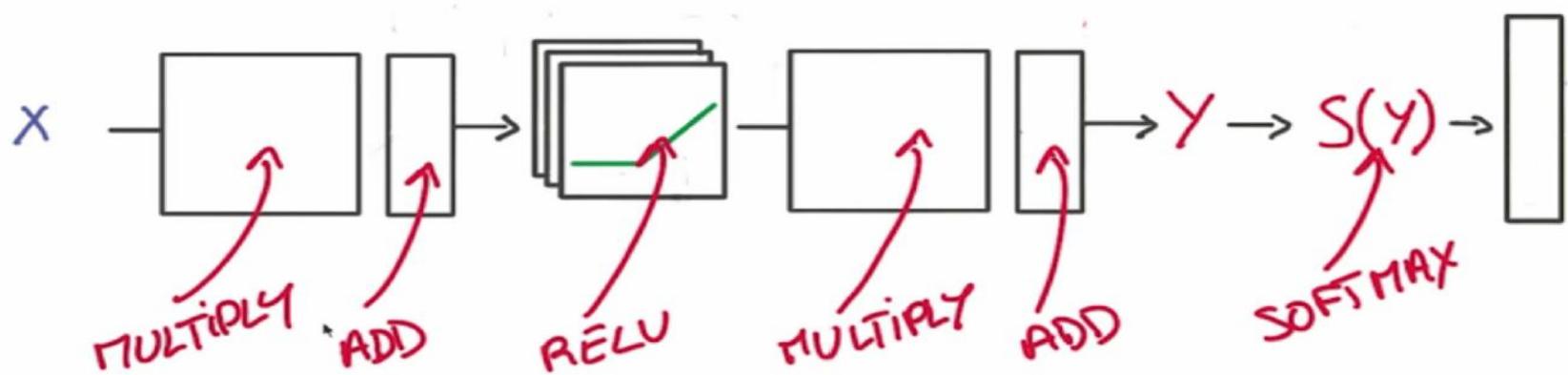
NEUROMORPHIC
ENGINEERING?

The Chain Rule

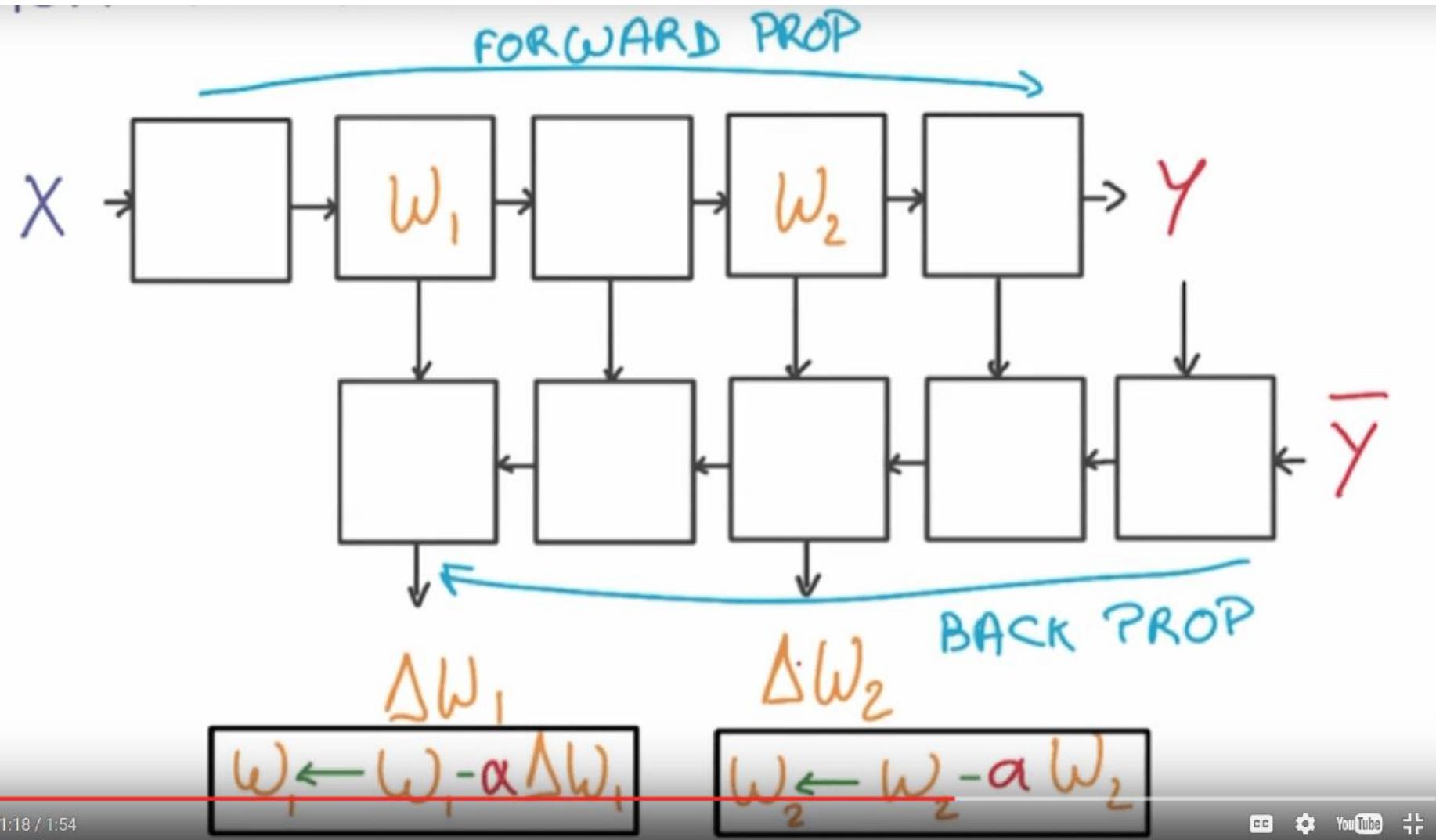
$$[\phi(f(x))]' = \phi'(f(x)) \cdot f'(x)$$

DERIVATIVE *PRODUCT*

STACKING UP SIMPLE OPERATIONS



Back - propagation



Optimisation tricks

REGULARIZATION

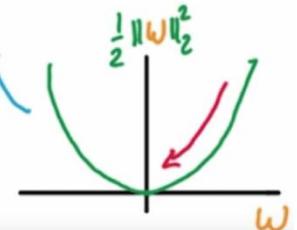


L₂ REGULARIZATION

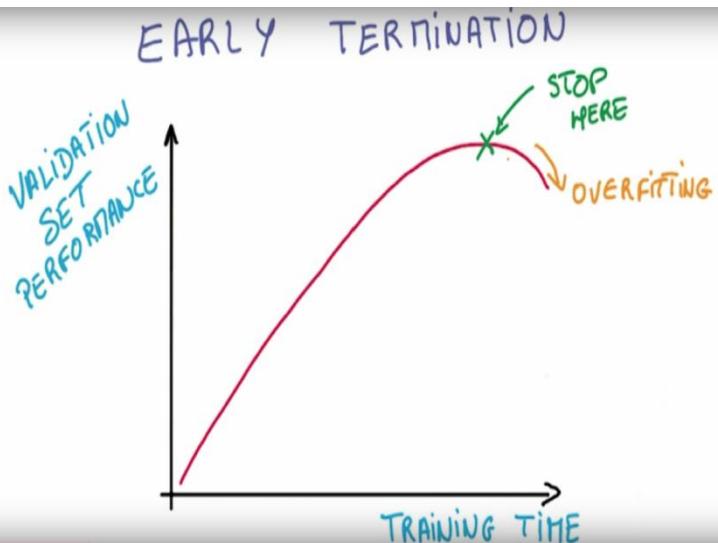
NEW LOSS

$$\mathcal{L}' = \mathcal{L} + \beta \frac{1}{2} \|W\|_2^2$$

LOSS



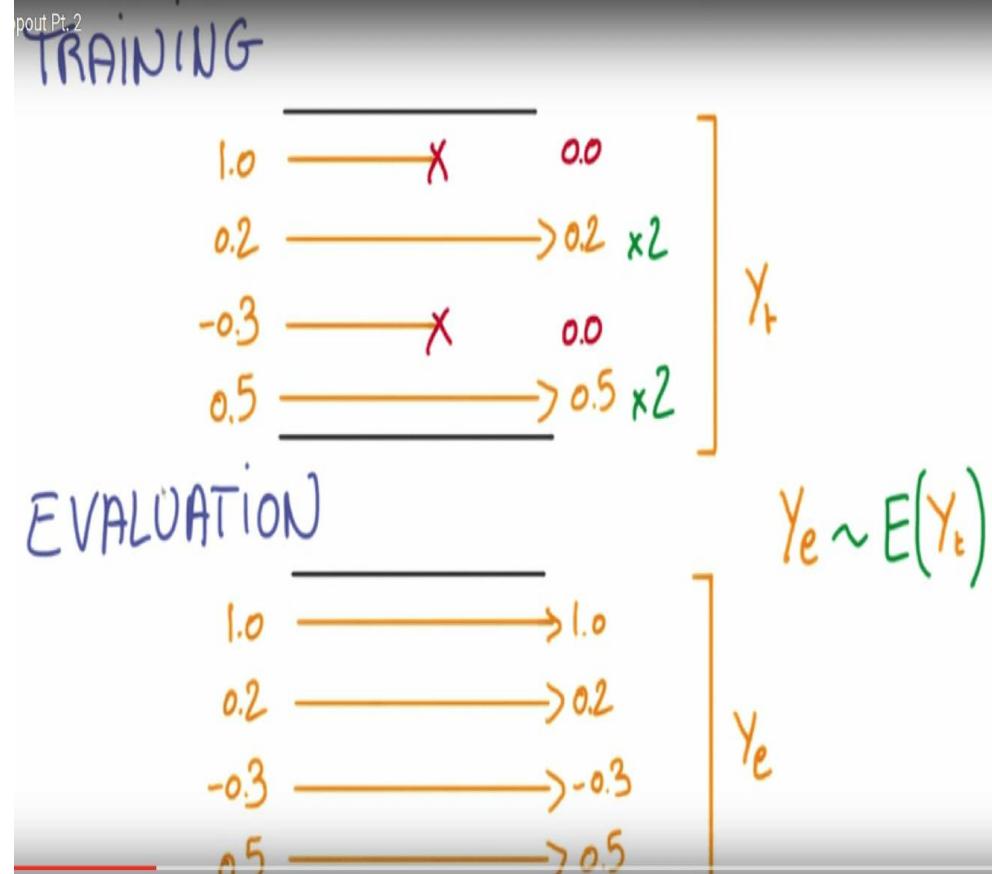
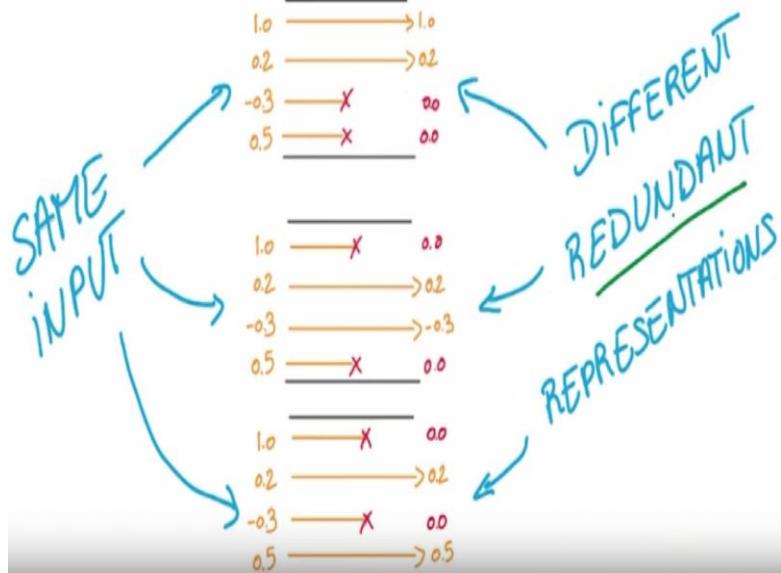
EARLY TERMINATION



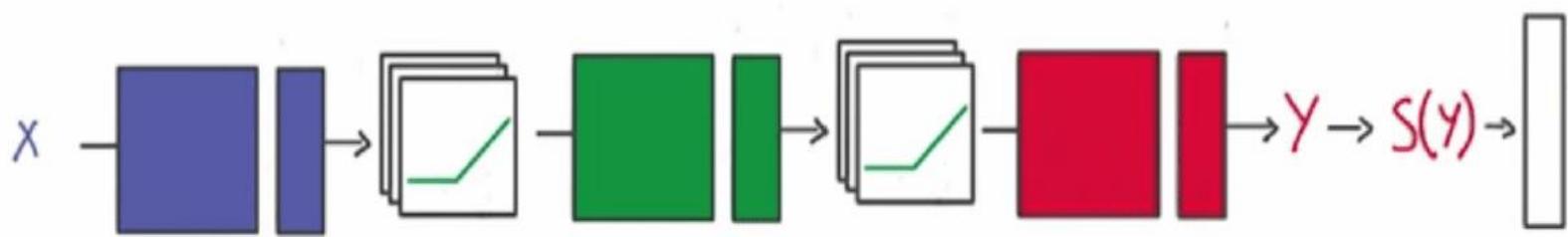
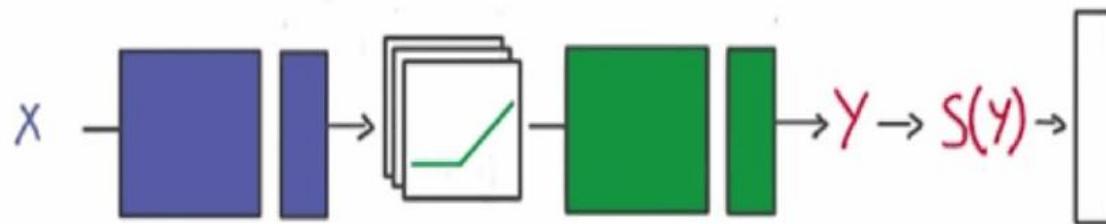
ACTIVATIONS

1.0	x	0.0
0.2		> 0.2
-0.3	x	0.0
0.5		> 0.5

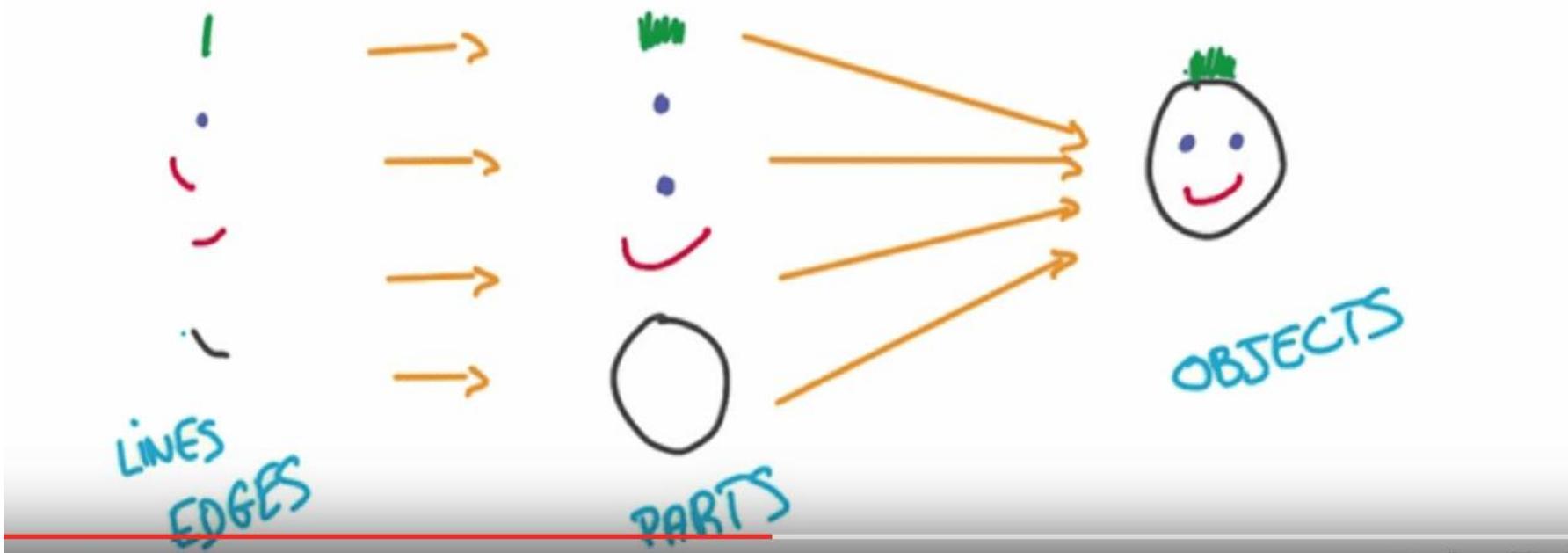
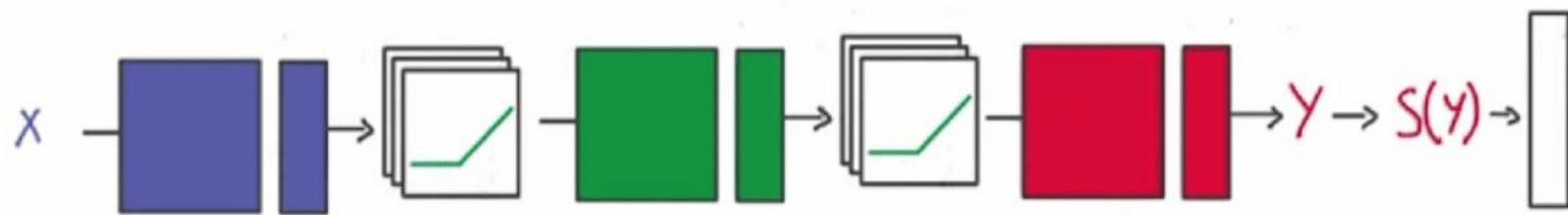
Optimisation trick: dropout



Deep networks



Deep networks



tensorflow.org/paper/whitepaper2015.pdf

TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems

(Preliminary White Paper, November 9, 2015)

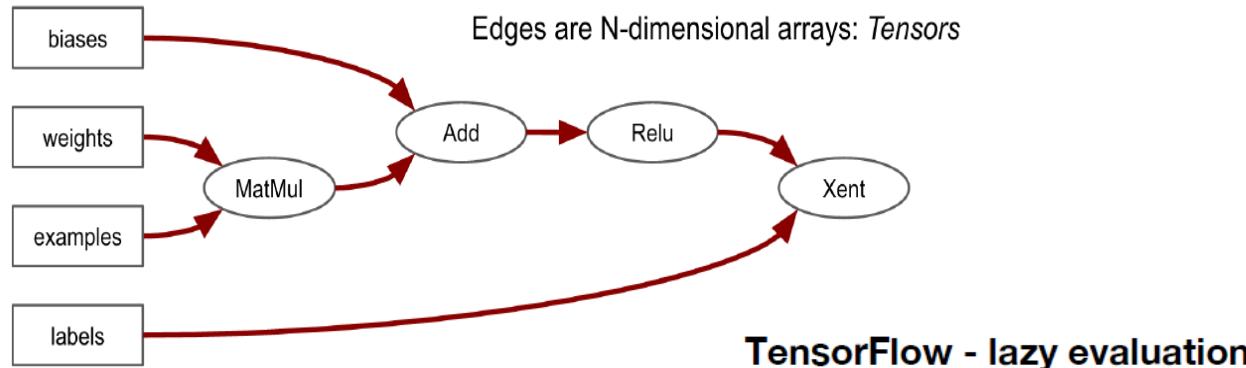
Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zhen^o
Google Research*

Abstract

TensorFlow [1] is an interface for expressing machine learning algorithms, and an implementation for executing such algorithms. A computation expressed using TensorFlow can be executed with little or no change on a wide variety of heterogeneous systems, ranging from mobile devices such as phones and tablets up to large-scale distributed systems of hundreds of machines and thousands of computational devices such as GPU cards. The system is flexible and can be used to express a wide variety of algorithms, including training and inference algorithms for deep neural network models, and it has been used for conducting research and for deploying machine learning systems into production across more than a dozen areas of computer science and other fields, including speech recognition, computer vision, robotics, information retrieval, natural language processing, geographic information extraction, and computational drug discovery. This paper describes the TensorFlow interface and an implementation of that interface that we have built at Google. The TensorFlow API and a reference implementation were released as an open-source package under the Apache 2.0 license in November, 2015 and are available at www.tensorflow.org.



TensorFlow is an open source software library for numerical computation using data flow graphs.



What is a Data Flow Graph?

Data flow graphs describe mathematical computation with a directed graph of nodes & edges. Nodes typically implement mathematical operations, but can also represent endpoints to feed in data, push out results, or read/write persistent variables. Edges describe the input/output relationships between nodes. These data edges carry dynamically-sized multidimensional data arrays, or tensors. The flow of tensors through the graph is where TensorFlow gets its name. Nodes are assigned to computational devices and execute asynchronously and in parallel once all the tensors on their incoming edges becomes available.

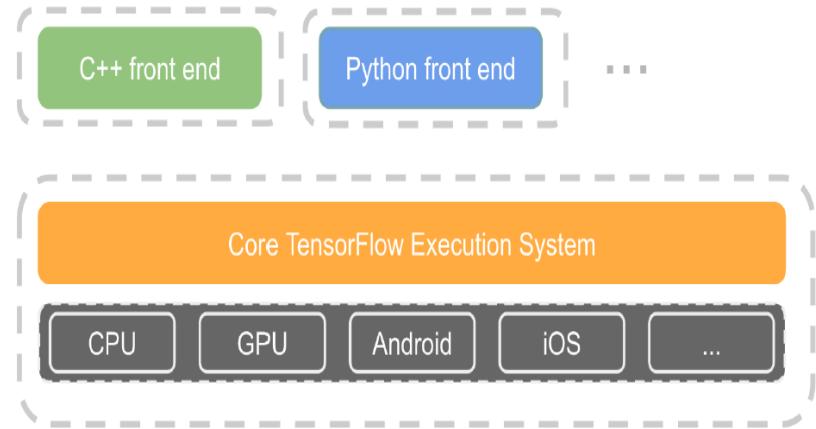
Key Features

- Deep Flexibility
- True Portability
- Connect Research and Production
- Auto-Differentiation
- Language Options
- Maximize Performance



Parallel Execution

- Launch graph in a Session
- Request output of some Ops with Run API
- TensorFlow computes set of Ops that must run to compute the requested outputs
- Ops execute, in parallel, as soon as their inputs are available



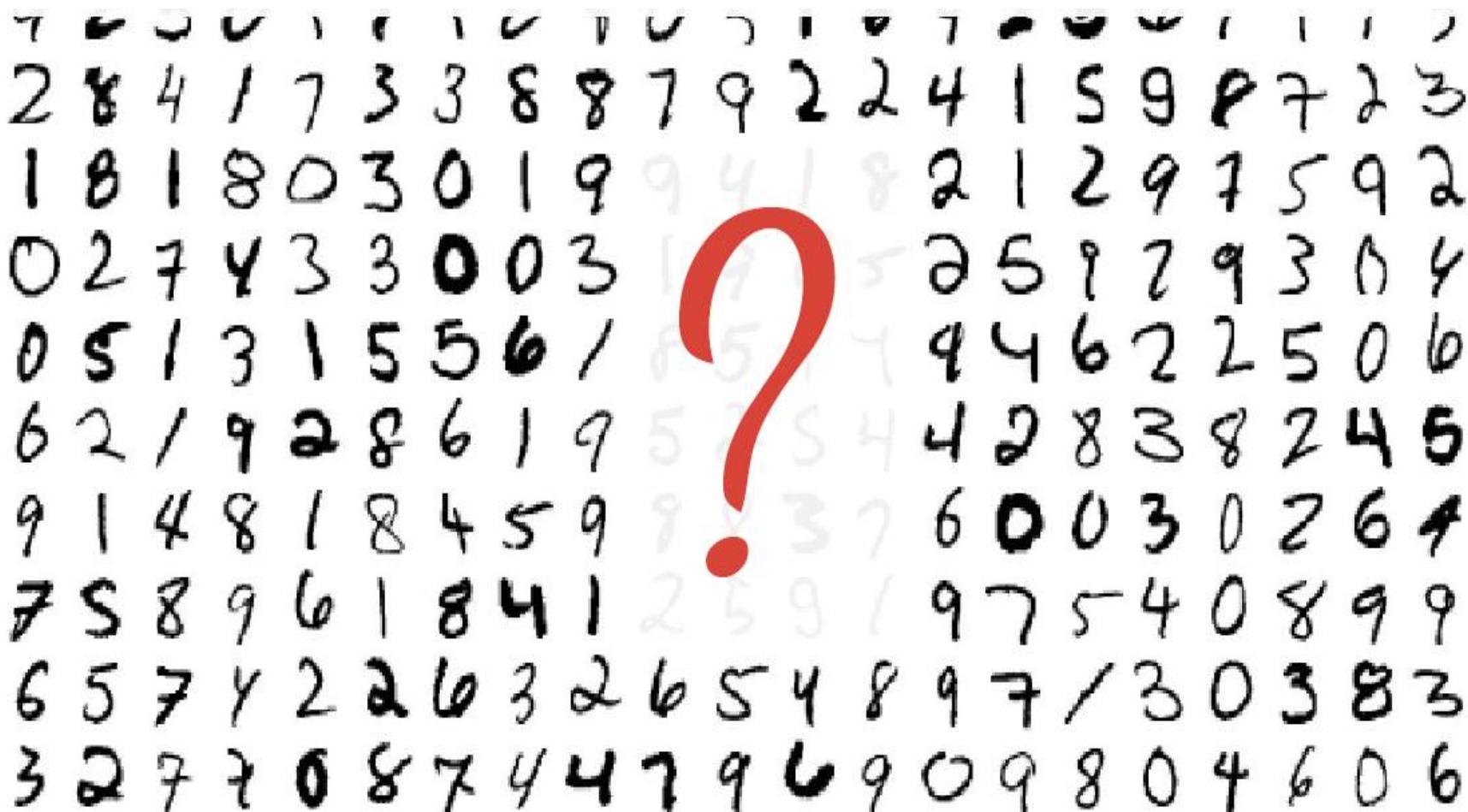
There are hundreds of predefined Ops (and easy to add more)

- Basics: constant, random, placeholder, cast, shape
- Variables: assign, assign_sub, assign_add
- Queues: enqueue, enqueue_batch, dequeue, blocking or not.
- Logical: equal, greater, less, where, min, max, argmin, argmax.
- Tensor computations: all math ops, matmul, determinant, inverse, cholesky.
- Images: encode, decode, crop, pad, resize, color spaces, random perturbations.
- Sparse tensors: represented as 3 tensors.
- ...

And many neural net specific Ops

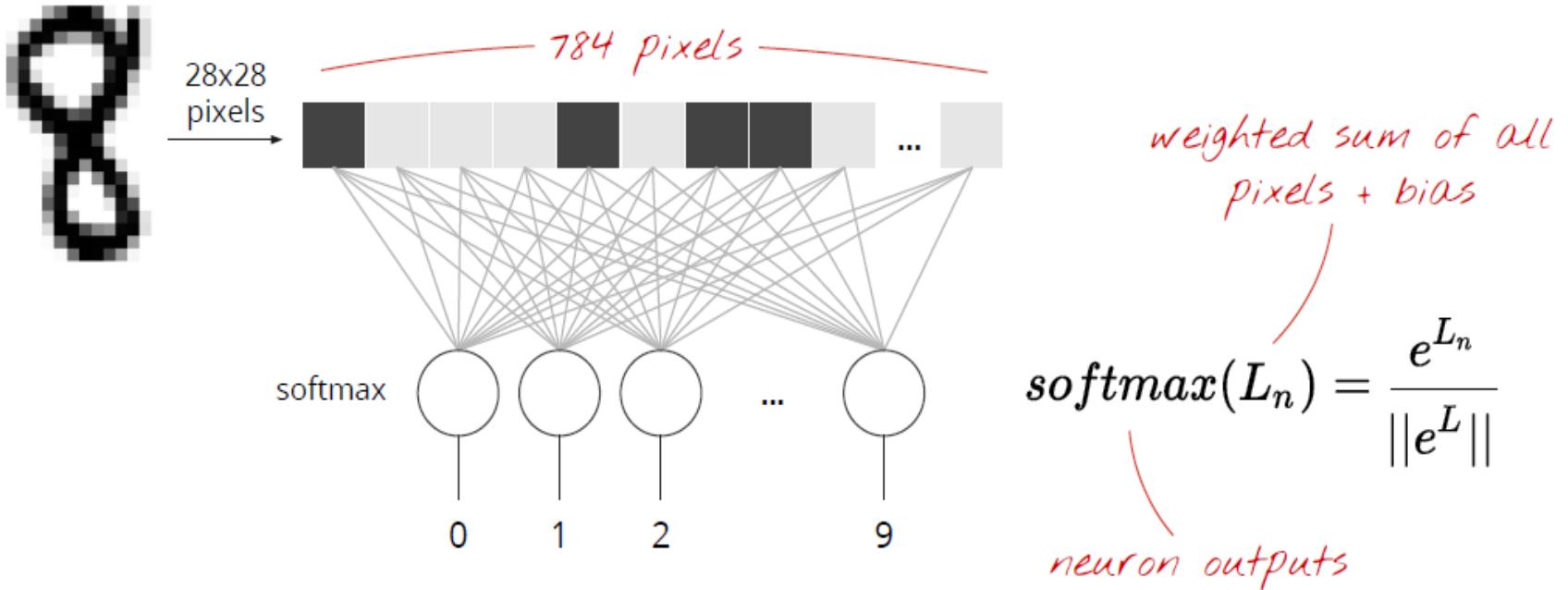
- Activations: sigmoid, tanh, relu, dropout, ...
- Pooling: avg, max.
- Convolutions: with many options.
- Normalization: local, batch, moving averages.
- Classification: softmax, softmax loss, cross entropy loss, topk.
- Embeddings: lookups/gather, scatter/updates.
- Sampling: candidate sampler (various options), sampling softmax.
- Updates: "fused ops" to speed-up optimizer updates (Adagrad, Momentum.)
- Summaries: Capture information for visualization.

Hand-written digits: MNIST



MNIST = Mixed National Institute of Standards and Technology - Download the dataset at <http://yann.lecun.com/exdb/mnist/>

Simple linear model



Slides from M. Gorner tutorial
<http://www.youtube.com/watch?v=vq2nnJ4g6NO>

TensorFlow full python code

```
import tensorflow as tf  
  
X = tf.placeholder(tf.float32, [None, 28, 28, 1])  
W = tf.Variable(tf.zeros([784, 10]))  
b = tf.Variable(tf.zeros([10]))  
init = tf.initialize_all_variables()  
  
# model  
Y = tf.nn.softmax(tf.matmul(tf.reshape(X, [-1, 784]), W) + b)  
  
# placeholder for correct answers  
Y_ = tf.placeholder(tf.float32, [None, 10])  
  
# Loss function  
cross_entropy = -tf.reduce_sum(Y_ * tf.log(Y))  
  
# % of correct answers found in batch  
is_correct = tf.equal(tf.argmax(Y, 1), tf.argmax(Y_, 1))  
accuracy = tf.reduce_mean(tf.cast(is_correct, tf.float32))
```

initialisation

model

SUCCESS metrics

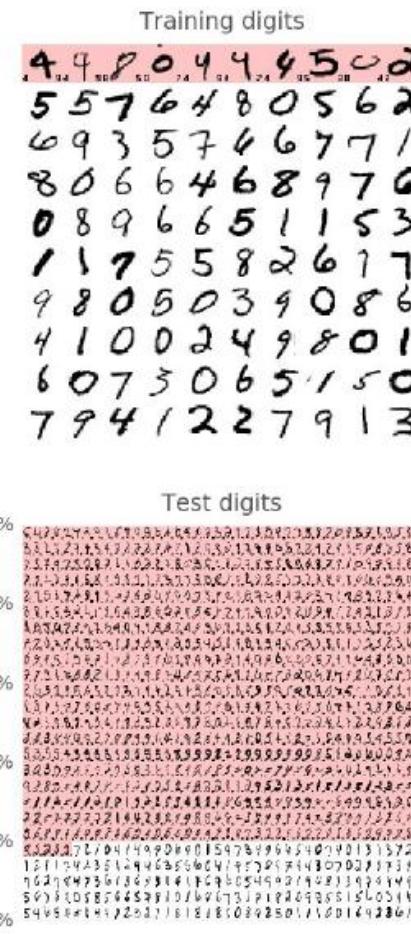
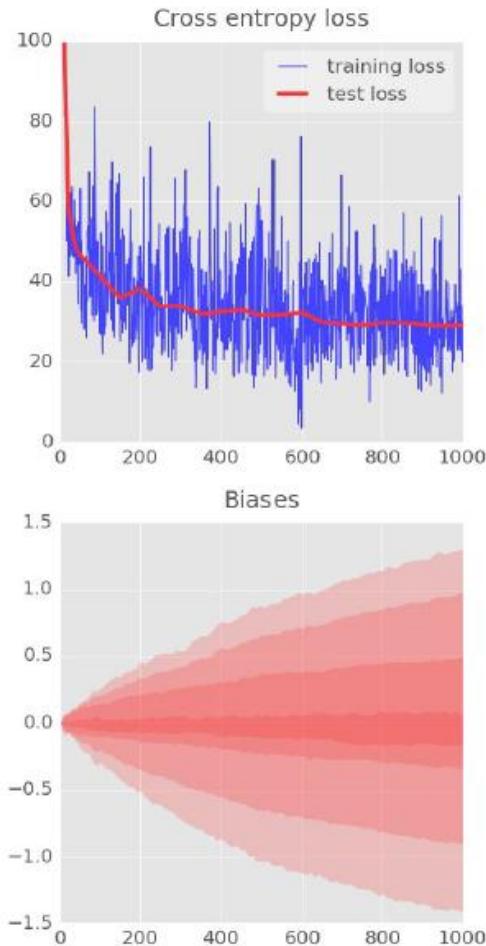
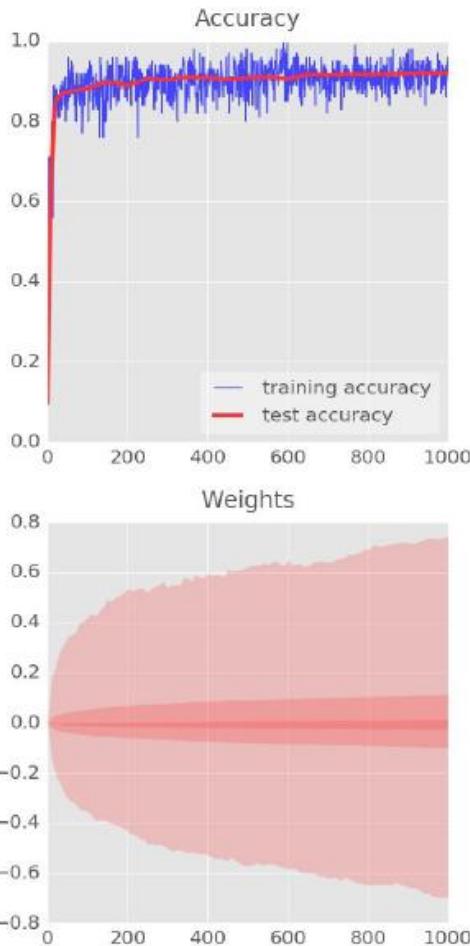
```
optimizer = tf.train.GradientDescentOptimizer(0.003)  
train_step = optimizer.minimize(cross_entropy)  
  
sess = tf.Session()  
sess.run(init)  
  
for i in range(10000):  
    # Load batch of images and correct answers  
    batch_X, batch_Y = mnist.train.next_batch(100)  
    train_data = {X: batch_X, Y_: batch_Y}  
  
    # train  
    sess.run(train_step, feed_dict=train_data)  
  
    # success ? add code to print it  
    a, c = sess.run([accuracy, cross_entropy], feed=train_data)  
  
    # success on test data ?  
    test_data = {X: mnist.test.images, Y_: mnist.test.labels}  
    a, c = sess.run([accuracy, cross_entropy], feed=test_data)
```

training step

Run

Slides from M. Gorner@youtube

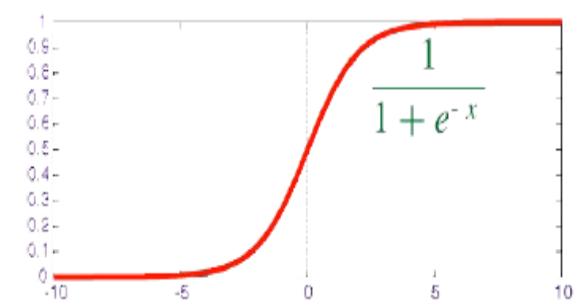
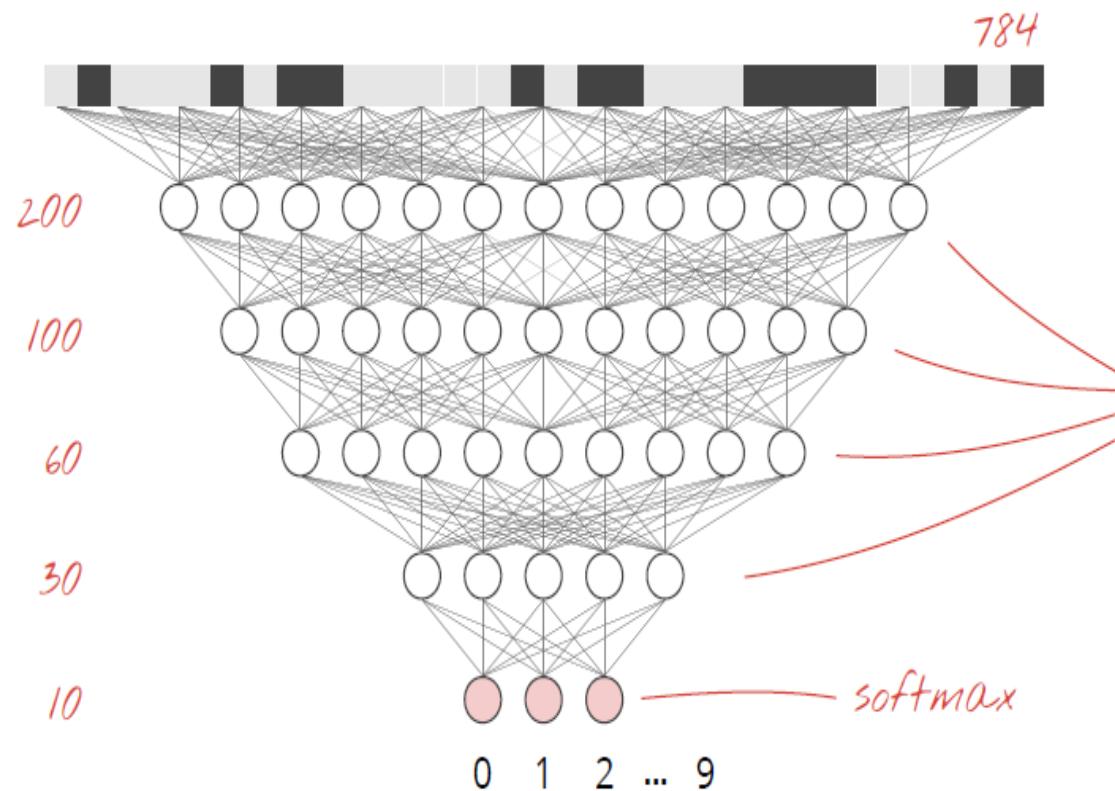
Simple linear model



Slides from M. Gorner@youtube

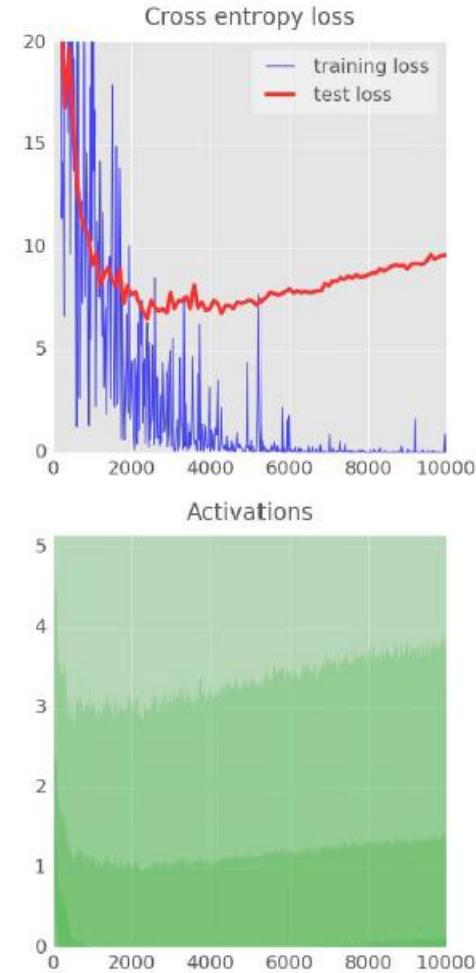
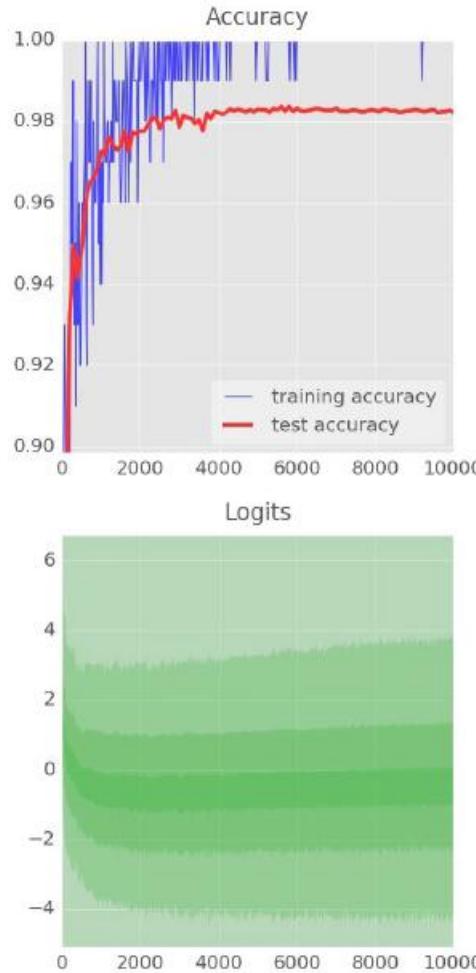
Multi-layer connected network

overKill



Slides from M. Gorner@youtube

Multi-layer connected network



Training digits

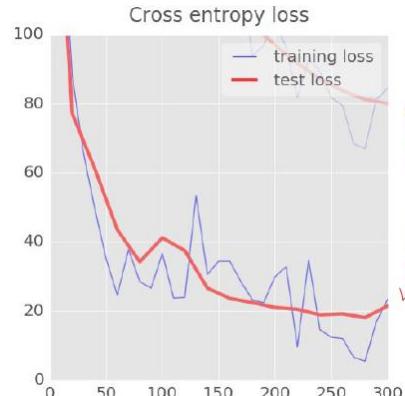
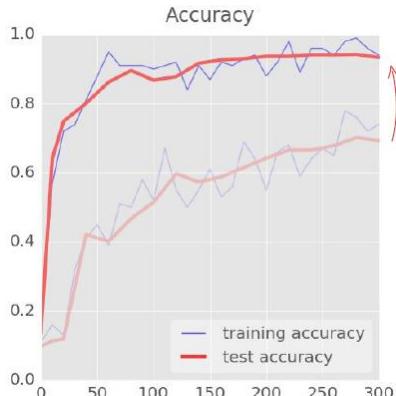
Test digits

98%
98%
98%
98%
98%
94%
94%
90%

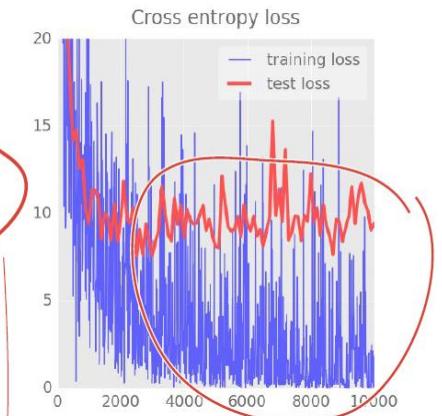
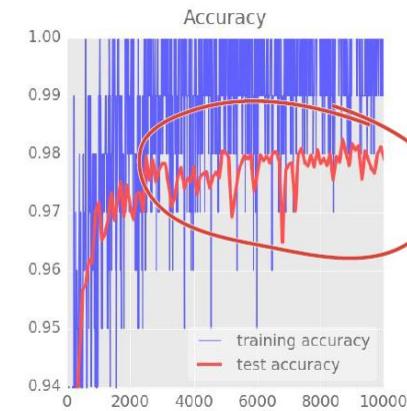
98!.

All tricks count

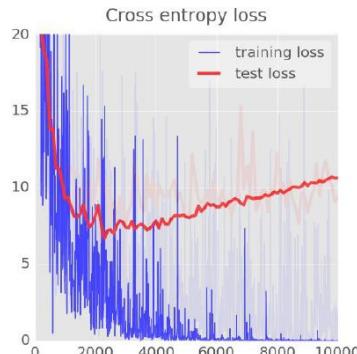
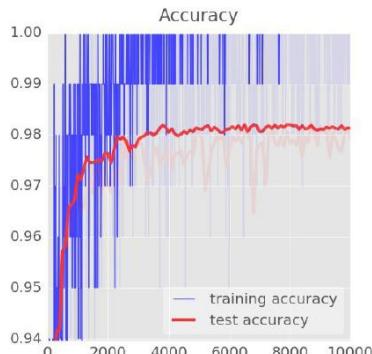
Use RELU



But noisy accuracy

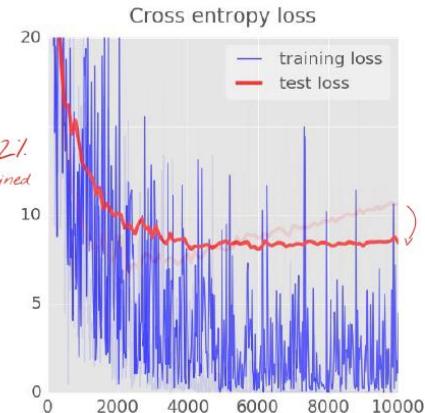
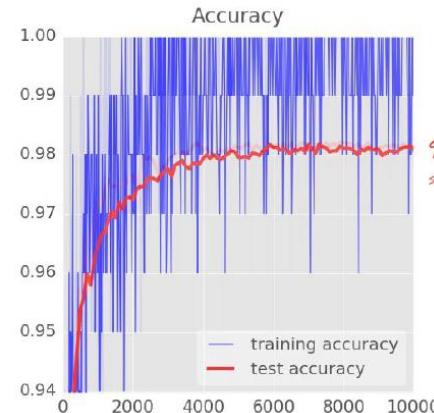


Exponentially reduce learning rates



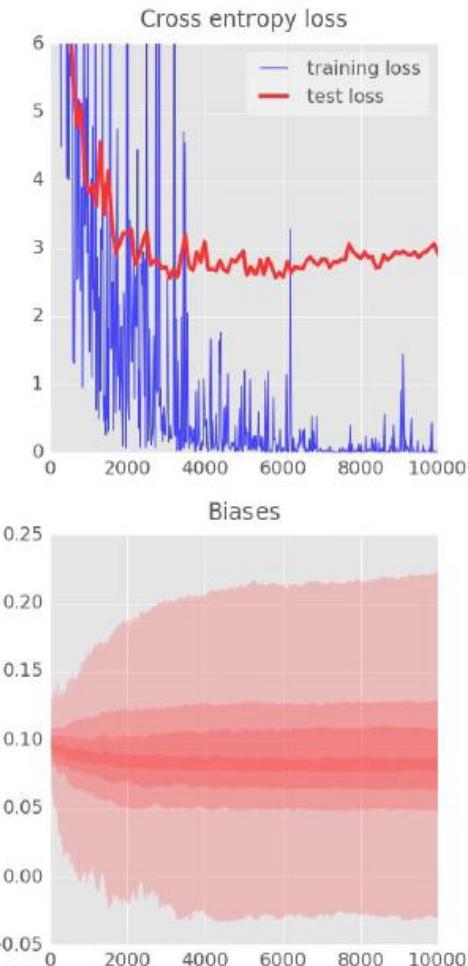
Learning rate 0.003 at start then dropping exponentially to 0.0001

Add drop-out



RELU, decaying learning rate 0.003 → 0.0001 and dropout 0.75

Can do better with conv network



References

<http://www.deeplearning.book.org>

<http://download.tensorflow.org/paper/whitepaper2015.pdf>

<https://www.tensorflow.org/>

<http://www.youtube.com/watch?v=vq2nnJ4g6NO>

<https://www.udacity.com/course/deep-learning--ud730>

TensorFlow application

PHYSICAL REVIEW D 94, 093001 (2016)

Potential for optimizing the Higgs boson CP measurement in $H \rightarrow \tau\tau$ decays at the LHC including machine learning techniques

R. Józefowicz,¹ E. Richter-Was,² and Z. Was³

¹*Open AI, San Francisco, California 94110, USA*

²*Institute of Physics, Jagellonian University, Lojasiewicza 11, 30-348 Krakow, Poland*

³*Institute of Nuclear Physics Polish Academy of Sciences, PL-31342 Krakow, Poland*

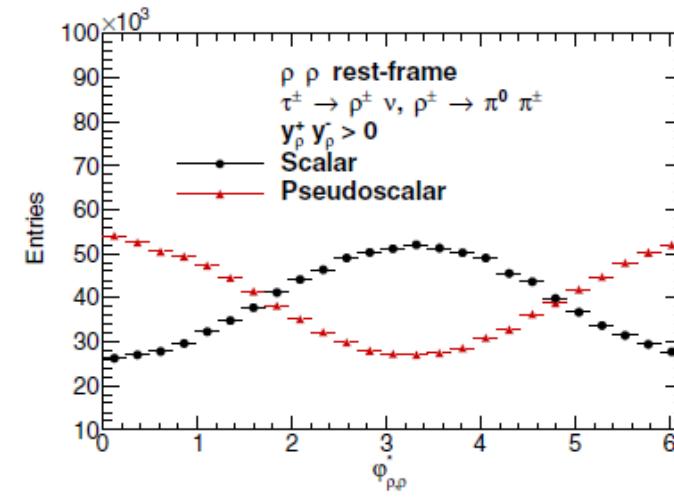
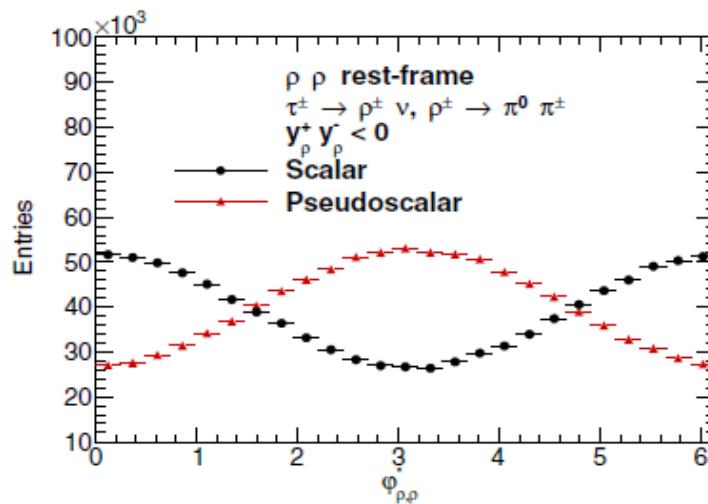
(Received 15 August 2016; published 7 November 2016)

- Multi-particle final state (cascade decays): 4 vectors
- CP information in correlations between decay planes and angles
- Physics intuition allowed to 1D optimal observables, but is there more to be explored.

Defining the problem

TABLE II. Dimensionality of the features which may be used in each discussed configuration of the decay modes. Note that in principle y_i^\pm, y_k^\mp may be calculated in the rest frame of the resonance pair used to define $\varphi_{i,k}^*$ planes, but in practice, a choice of the frames is of no numerically significant effect. We do not distinguish such variants.

Features/variables	Decay mode: $\rho^\pm - \rho^\mp$		Decay mode: $a_1^\pm - \rho^\mp$	Decay mode: $a_1^\pm - a_1^\mp$
	$\rho^\pm \rightarrow \pi^0 \pi^\pm$	$a_1^\pm \rightarrow \rho^0 \pi^\pm, \rho^0 \rightarrow \pi^+ \pi^-$ $\rho^\mp \rightarrow \pi^0 \pi^\mp$		
$\varphi_{i,k}^*$	1		4	16
$\varphi_{i,k}^*$ and y_i, y_k	3		9	24
$\varphi_{i,k}^*$, 4-vectors	25		36	64
$\varphi_{i,k}^*, y_i, y_k$ and m_i, m_k	5		13	30
$\varphi_{i,k}^*, y_i, y_k, m_i, m_k$ and 4-vectors	29		45	78
$\varphi_{i,k}^*, y_i, y_k, m_i, m_k$ and 4-vectors	29		45	78



Defining NN model

- 6 hidden layers, each 300 nodes, with RELU activation function

$$D \rightarrow 300 \rightarrow 300 \rightarrow 300 \rightarrow 300 \rightarrow 300 \rightarrow 300 \rightarrow 1.$$

- Sigmoid function on last layer $[\text{sigmoid}(x) = 1/(1 + \exp(-x))]$
- Metric: negative log-likelihood of the true targets under Bernoulli distribution
 - $-\log p(y|y_h) = -(y == 0) * \log(y_h) - (y == 1) * \log(1 - y_h),$
- Optimisation: SDG Adam algorithm
- Optimisation: batch normalisation, dropout
- Final score: weighted AUC

Defining NN model

```
# Linearly transforms X of shape [batch_size, size1] into [batch_size, size2].  
# Applies X -> XW+b, where W and b are trainable parameters.  
  
def linear(x, name, size, bias=True):  
    w = tf.get_variable(name + "/W", [x.get_shape()[1], size])  
    b = tf.get_variable(name + "/b", [1, size],  
                       initializer=tf.zeros_initializer)  
    return tf.matmul(x, w) + b  
  
  
# Applies batch normalization trick from https://arxiv.org/abs/1502.03167  
# by normalizing each feature in a batch.  
  
def batch_norm(x, name):  
    mean, var = tf.nn.moments(x, [0])  
    normalized_x = (x - mean) * tf.rsqrt(var + 1e-8)  
    gamma = tf.get_variable(name + "/gamma", [x.get_shape()[-1]],  
                           initializer=tf.constant_initializer(1.0))  
    beta = tf.get_variable(name + "/beta", [x.get_shape()[-1]])  
    return gamma * normalized_x + beta
```

```
class NeuralNetwork(object):  
  
    def __init__(self, num_features, batch_size, num_layers=6, size=300, lr=1e-3):  
        # Each input x is represented by a given number of features  
        # and corresponding weights for target distributions A and B.  
        self.x = x = tf.placeholder(tf.float32, [batch_size, num_features])  
        self.wa = wa = tf.placeholder(tf.float32, [batch_size])  
        self.wb = wb = tf.placeholder(tf.float32, [batch_size])  
        # The model will predict a single number, which is a probability of input x  
        # belonging to class A. That probability is equal to wa / (wa+wb).  
        y = wa / (wa+wb)  
        y = tf.reshape(y, [-1, 1])  
  
        # We apply multiple layers of transformations where each layer consists of  
        # linearly transforming the features, followed by batch normalization  
        # (described above)  
        # and ReLU nonlinearity (which is an elementwise operation: x -> max(x, 0))  
        for i in range(num_layers):  
            x = tf.nn.relu(batch_norm(linear(x, "linear %d" % i, size), "bn %d" % i))  
  
        # Finally, the output is transformed into a single number.  
        # After applying sigmoid nonlinearity (x -> 1 / (1 + exp(-x))) we'll interpret  
        # that number  
        # as a probability of x belonging to class A.  
        x = linear(x, "regression", 1)  
        self.p = tf.nn.sigmoid(x)  
  
        # The objective to optimize is negative log likelihood under Bernoulli  
        # distribution:  
        # loss = - (p(y==A) * log p(y==A|x) + p(y==B) * log p(y==B|x))  
        self.loss = loss = tf.reduce_mean(tf.nn.sigmoid_cross_entropy_with_logits  
                                         (x, y))  
        # The model parameters are optimized using gradient-based Adam optimizer  
        # (https://arxiv.org/abs/1412.6980) to minimize the loss on the training data  
        self.train_op = tf.train.AdamOptimizer(lr).minimize(loss)
```

Results

TABLE III. Average probability p_i (calculated as explained in Sec. III B) that a model predicts correctly event x_i to be of a type A (scalar), with training being performed for separation between types A and B (pseudoscalar).

Features/variables	Decay mode: $\rho^\pm - \rho^\mp$ $\rho^\pm \rightarrow \pi^0 \pi^\pm$		Decay mode: $a_1^\pm - \rho^\mp$ $a_1^\pm \rightarrow \rho^0 \pi^\mp, \rho^0 \rightarrow \pi^+ \pi^-$ $\rho^\mp \rightarrow \pi^0 \pi^\mp$	Decay mode: $a_1^\pm - a_1^\mp$ $a_1^\pm \rightarrow \rho^0 \pi^\pm, \rho^0 \rightarrow \pi^+ \pi^-$
True classification	0.782		0.782	0.782
$\varphi_{i,k}^*$	0.500		0.500	0.500
$\varphi_{i,k}^*$ and y_i, y_k	0.624		0.569	0.536
4-vectors	0.638		0.590	0.557
$\varphi_{i,k}^*$, 4-vectors	0.638		0.594	0.573
$\varphi_{i,k}^*, y_i, y_k$ and m_i^2, m_k^2	0.626		0.578	0.548
$\varphi_{i,k}^*, y_i, y_k, m_i^2, m_k^2$ and 4-vectors	0.639		0.596	0.573

- NN can capture „optimal variables” but can do better given simple of 4-vectors.
- Given „simple” and „higher level” features can still improve in more complicated case.
- Will try now on the experimental data.