

FASHION-MNIST

FULLY CONNECTED NEURAL NETWORKS



Stephanie & Fanny

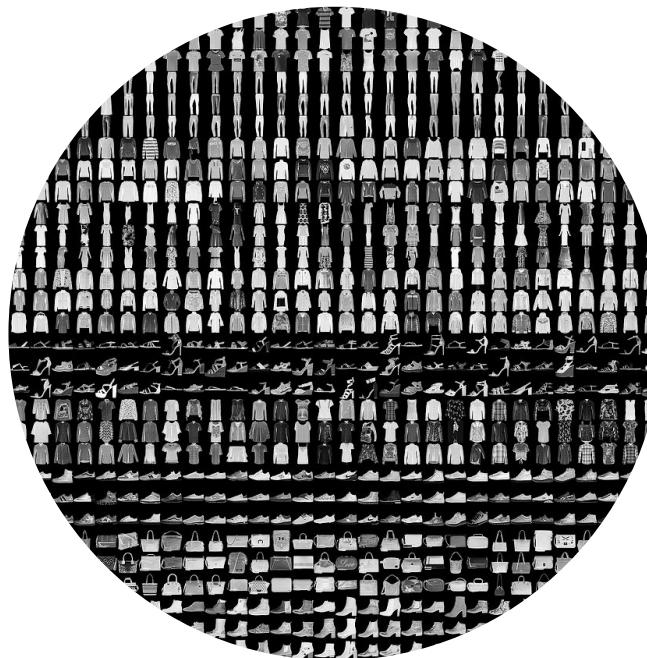
OBJECTIVES

1. Brief overview of Fashion-MNIST Data set
2. Fully connected neural network with 1 layer
3. 5 layers with sigmoid
4. 5 layers with relu
5. 5 layers relu & exponential decay
6. 5 layers relu exponential decay & dropout

FASHION-MNIST

Zalando Research created:

- 60,000 training images
- 10,000 test images
- 28x28 gray-scale image
 - 784 pixels





SHOES & FASHION ONLINE - ZALANDO IS AVAILABLE IN THE FOLLOWING COUNTRIES



ZALANDO.DE	ZALANDO.FR	ZALANDO.LU
ZALANDO.AT	ZALANDO.IT	ZALANDO.NO
ZALANDO.CH	ZALANDO.ES	ZALANDO.SE
ZALANDO.UK	ZALANDO.NL	ZALANDO.DK
ZALANDO.PL	ZALANDO.BE	ZALANDO.FI

[Legal notice](#)

Zalando SE is a German electronic commerce company seated in Berlin. The company maintains a cross-platform online store that sells shoes, clothing and other fashion items.

“overall purpose of wanting to deliver fashion for the good of all.”



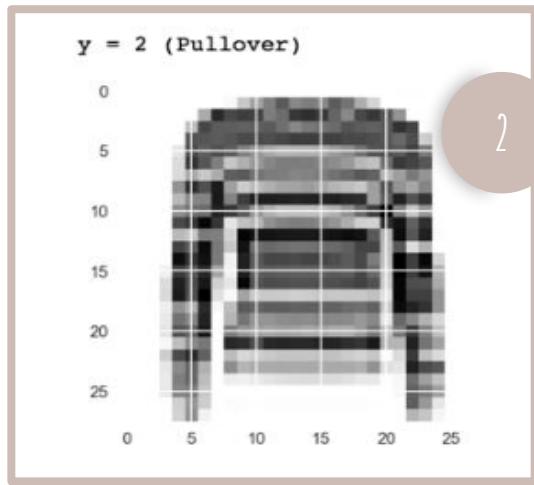
Zalando Launches Research Lab

by Dr. Reiner Kraft - 28 Sep 2016

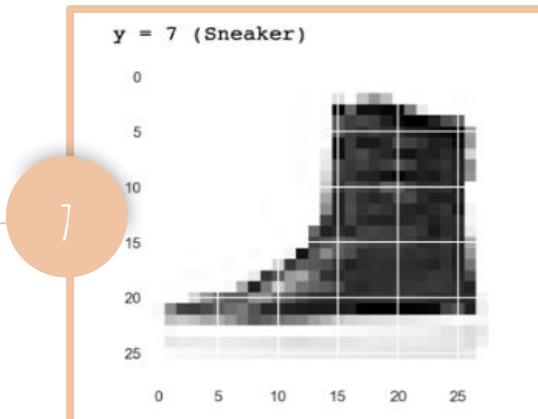
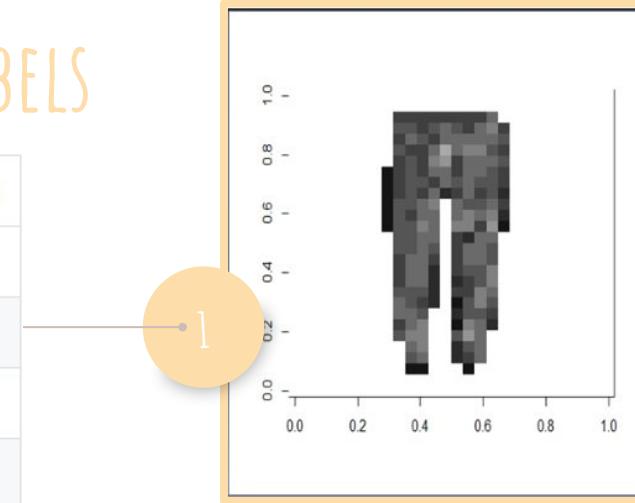
Zalando Tech is embarking on an exciting new chapter in its evolution. We've recently launched Zalando Research – an endeavour to place Zalando at the forefront of cutting-edge research, to complement our already strong foothold on technology.

Why are we doing this? We already boast great researchers at Zalando Tech, so we want to give them an outlet where they can be better organised and aligned. We also want to have a greater impact in the tech research community, especially in the fields of data science, machine learning, and artificial intelligence (AI). We want to carve out

10 DIFFERENT LABELS



Label	Description
0	T-shirt/top
1	Trouser
2	Pullover
3	Dress
4	Coat
5	Sandal
6	Shirt
7	Sneaker
8	Bag
9	Ankle boot

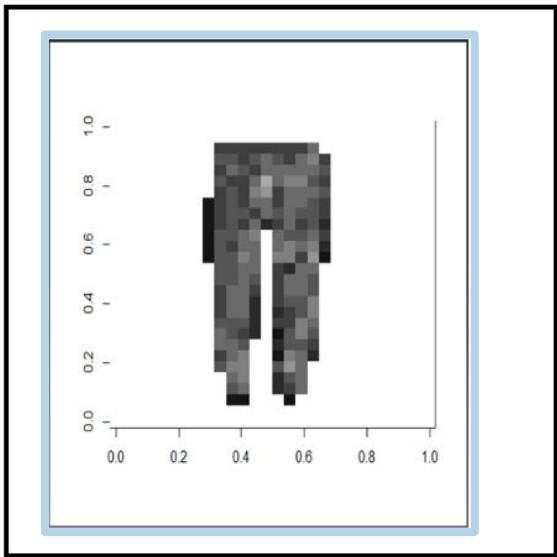


1. ONE LAYER FULLY-CONNECTED



Let's learn by explaining the foundations.

EVERY IMAGE CAN BE THOUGHT OF AS AN ARRAY OF NUMBERS DESCRIBING HOW DARK EACH PIXEL IS.



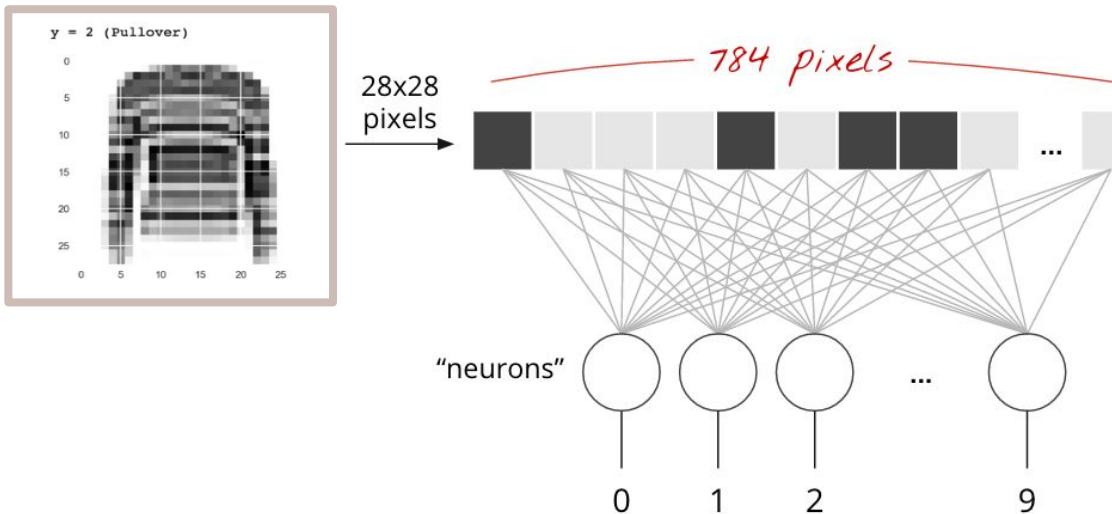
|?

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	.6	.8	0	0	0	0	0	0	0
0	0	0	0	0	0	0	.7	.1	0	0	0	0	0	0	0
0	0	0	0	0	0	0	.7	.1	0	0	0	0	0	0	0
0	0	0	0	0	0	0	.5	.1	.4	0	0	0	0	0	0
0	0	0	0	0	0	0	.6	.1	.4	0	0	0	0	0	0
0	0	0	0	0	0	0	.6	.1	.4	0	0	0	0	0	0
0	0	0	0	0	0	0	.6	.1	.7	0	0	0	0	0	0
0	0	0	0	0	0	0	.8	.1	.1	0	0	0	0	0	0
0	0	0	0	0	0	0	.9	.1	.1	0	0	0	0	0	0
0	0	0	0	0	0	0	.3	.1	.1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

READ-IN DATA

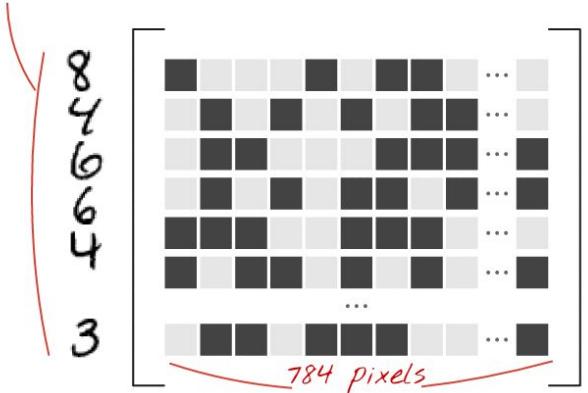
```
# Download images and labels into mnist.test and mnist.train
```

```
mnist = mnist_data.read_data_sets("data/fashion",  
                                  one_hot=True, reshape=False, validation_size=0)
```



How are weighted sums calculated?!

*X: 100 images,
one per line,
flattened*



10 columns									
$W_{0,0}$	$W_{0,1}$	$W_{0,2}$	$W_{0,3}$...	$W_{0,9}$				
$W_{1,0}$	$W_{1,1}$	$W_{1,2}$	$W_{1,3}$...	$W_{1,9}$				
$W_{2,0}$	$W_{2,1}$	$W_{2,2}$	$W_{2,3}$...	$W_{2,9}$				
$W_{3,0}$	$W_{3,1}$	$W_{3,2}$	$W_{3,3}$...	$W_{3,9}$				
$W_{4,0}$	$W_{4,1}$	$W_{4,2}$	$W_{4,3}$...	$W_{4,9}$				
$W_{5,0}$	$W_{5,1}$	$W_{5,2}$	$W_{5,3}$...	$W_{5,9}$				
$W_{6,0}$	$W_{6,1}$	$W_{6,2}$	$W_{6,3}$...	$W_{6,9}$				
$W_{7,0}$	$W_{7,1}$	$W_{7,2}$	$W_{7,3}$...	$W_{7,9}$				
$W_{8,0}$	$W_{8,1}$	$W_{8,2}$	$W_{8,3}$...	$W_{8,9}$				
...									
$W_{783,0}$	$W_{783,1}$	$W_{783,2}$	$W_{783,3}$...	$W_{783,9}$				

$$L = X \cdot W + b$$

Predictions

$Y[100, 10]$

Images

$X[100, 784]$

Weights

$W[784, 10]$

Biases

$b[10]$

$$Y = \text{softmax}(X \cdot W + b)$$

applied line
by line

matrix multiply

broadcast
on all lines

tensor shapes in []

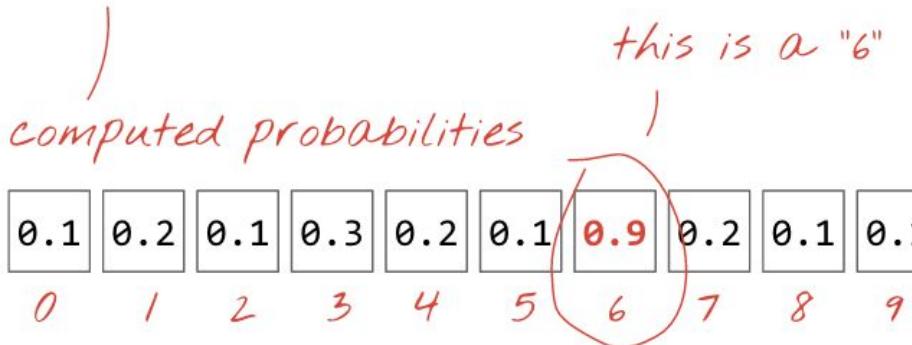
```
# The model
```

```
Y = tf.nn.softmax(tf.matmul(XX, W) + b)
```

0	1	2	3	4	5	6	7	8	9
0	0	0	0	0	0	1	0	0	0

actual probabilities, "one-hot" encoded

Cross entropy: $-\sum Y'_i \cdot \log(Y_i)$



TRAINING

Goal: Find weights and biases that minimize the distance between what system predicts and what we know to be true.

learning rate



```
optimizer = tf.train.GradientDescentOptimizer(0.003)  
train_step = optimizer.minimize(cross_entropy)
```

loss function

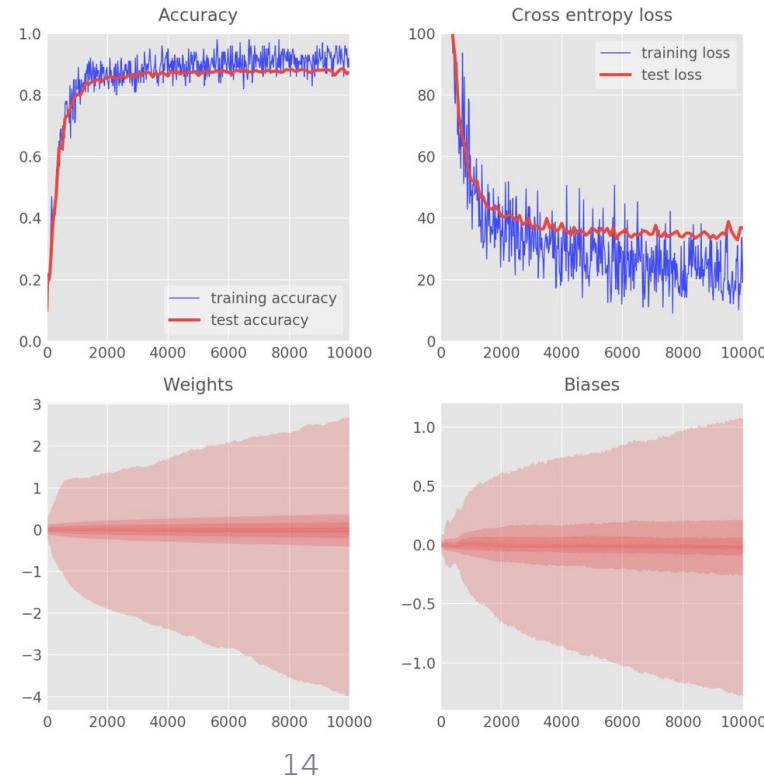


RESULTS

Softmax 1 layer

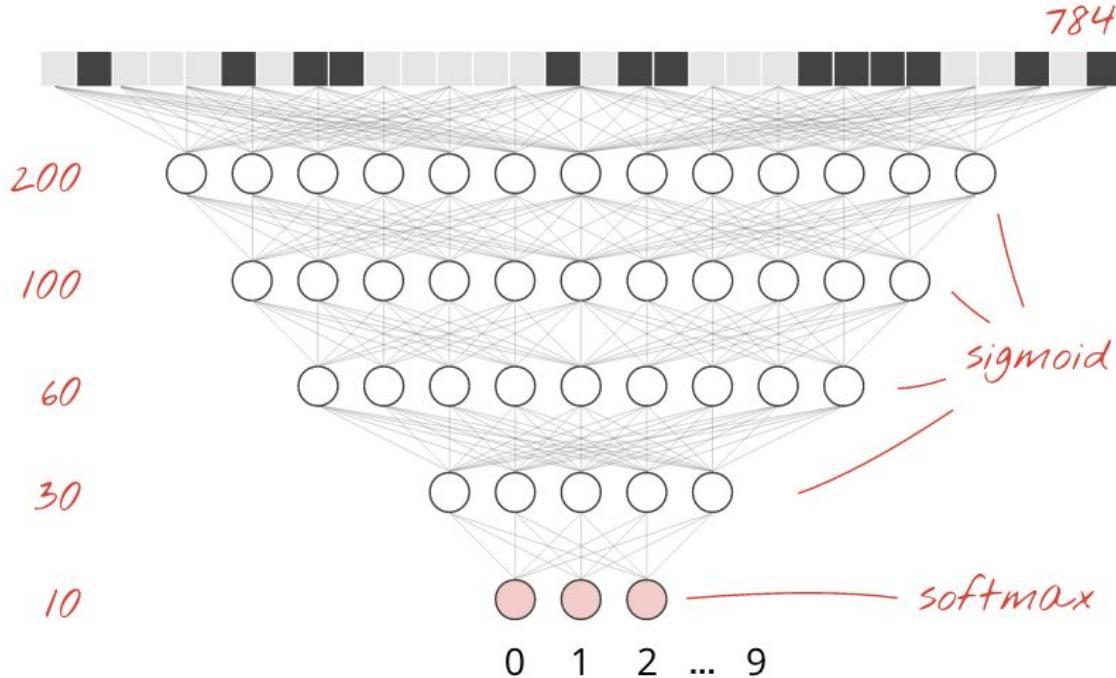
Fashion MNIST
Max test
accuracy: 83.13%

Digits Mnist
accuracy: 92%





5 LAYERS FULLY-CONNECTED WITH SIGMOID



Here is for example a 5-layer fully connected neural network

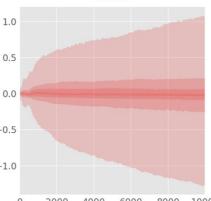
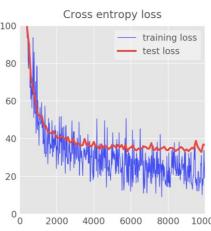
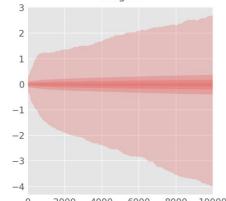
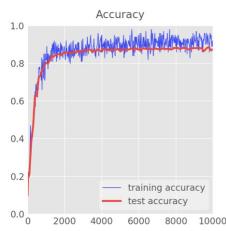
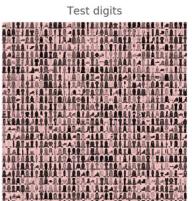
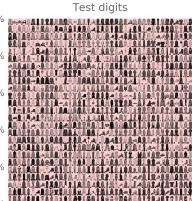
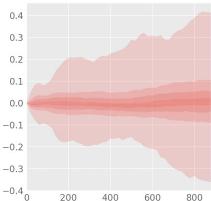
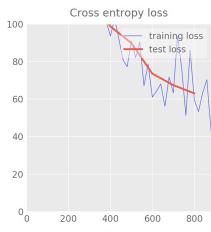
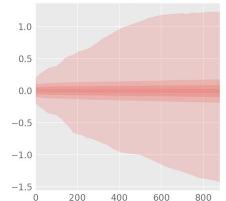
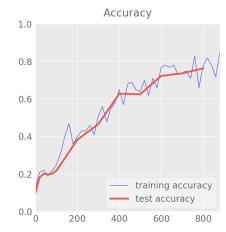
The neurons in the second layer compute weighted sums of neuron outputs from the previous layer.

Pros

- Max test accuracy: 88.45%
- Improvement compared to 1 layer softmax max test accuracy: 83.13%
- Adding layers improves accuracy

Cons

- Sigmoid activation function with many layers in initial iterations accuracy increases slowly

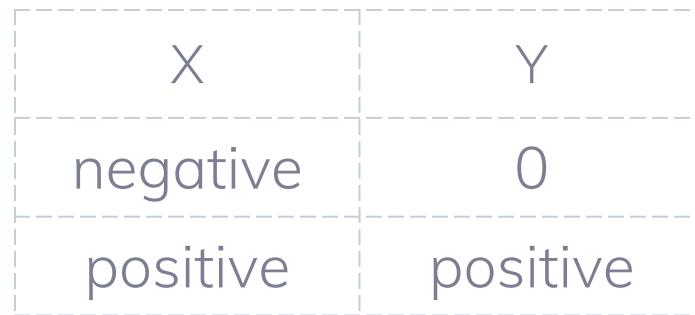
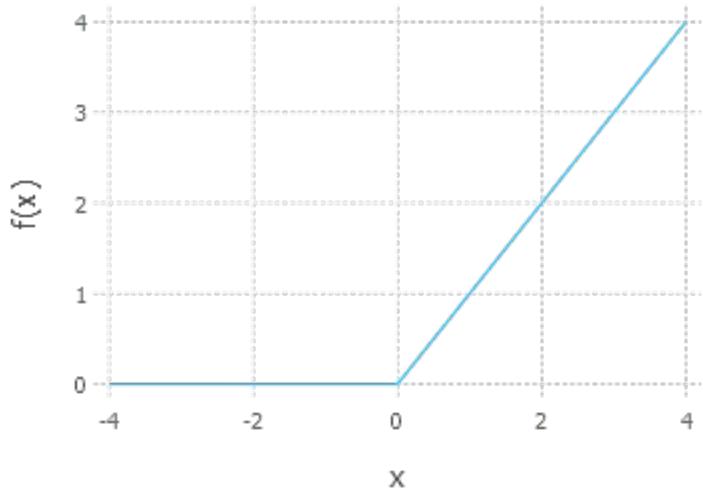


3.5 LAYER FULLY-CONNECTED WITH RELU



Sigmoid does not perform as well with more layers. Let's use the ReLu as our activation function instead.

REPLACE SIGMOID ACTIVATION FUNCTION WITH RELU (RECTIFIED LINEAR UNIT)



SWAP

```
# previous model
```

```
XX = tf.reshape(X, [-1, 784])  
Y1 = tf.nn.sigmoid(tf.matmul(XX, W1) + B1)  
Y2 = tf.nn.sigmoid(tf.matmul(Y1, W2) + B2)  
Y3 = tf.nn.sigmoid(tf.matmul(Y2, W3) + B3)  
Y4 = tf.nn.sigmoid(tf.matmul(Y3, W4) + B4)  
Ylogits = tf.matmul(Y4, W5) + B5  
Y = tf.nn.softmax(Ylogits)
```

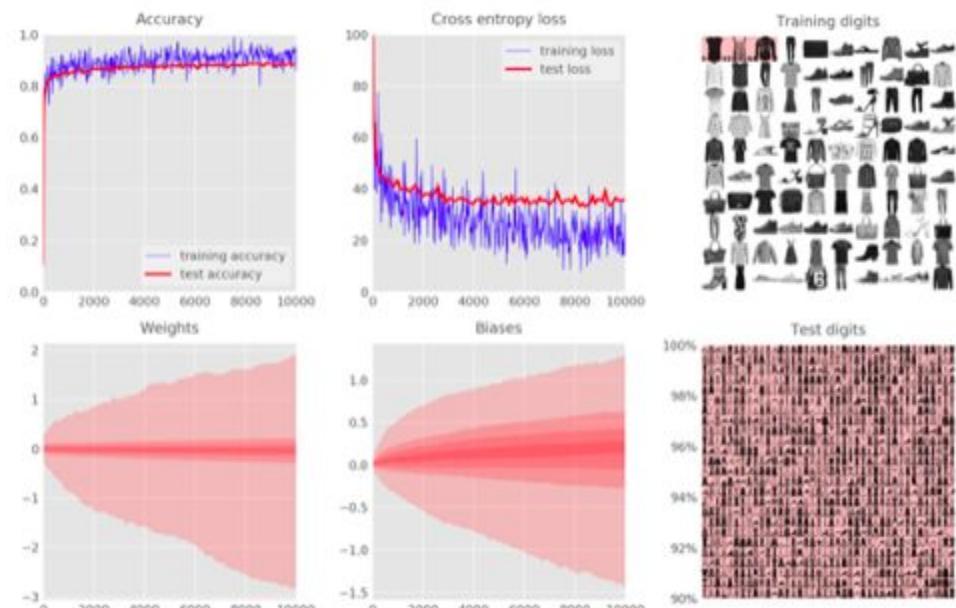
```
# new model
```

```
XX = tf.reshape(X, [-1, 784])  
Y1 = tf.nn.relu(tf.matmul(XX, W1) + B1)  
Y2 = tf.nn.relu(tf.matmul(Y1, W2) + B2)  
Y3 = tf.nn.relu(tf.matmul(Y2, W3) + B3)  
Y4 = tf.nn.relu(tf.matmul(Y3, W4) + B4)  
Ylogits = tf.matmul(Y4, W5) + B5  
Y = tf.nn.softmax(Ylogits)
```

5 LAYER FULLY-CONNECTED WITH RELU RESULTS

	Sigmoid	ReLU
Max test accuracy	0.8813	0.8893

Cons	Pros
Noisy accuracy curve	Faster training
Noisy cross entropy loss	Improved accuracy

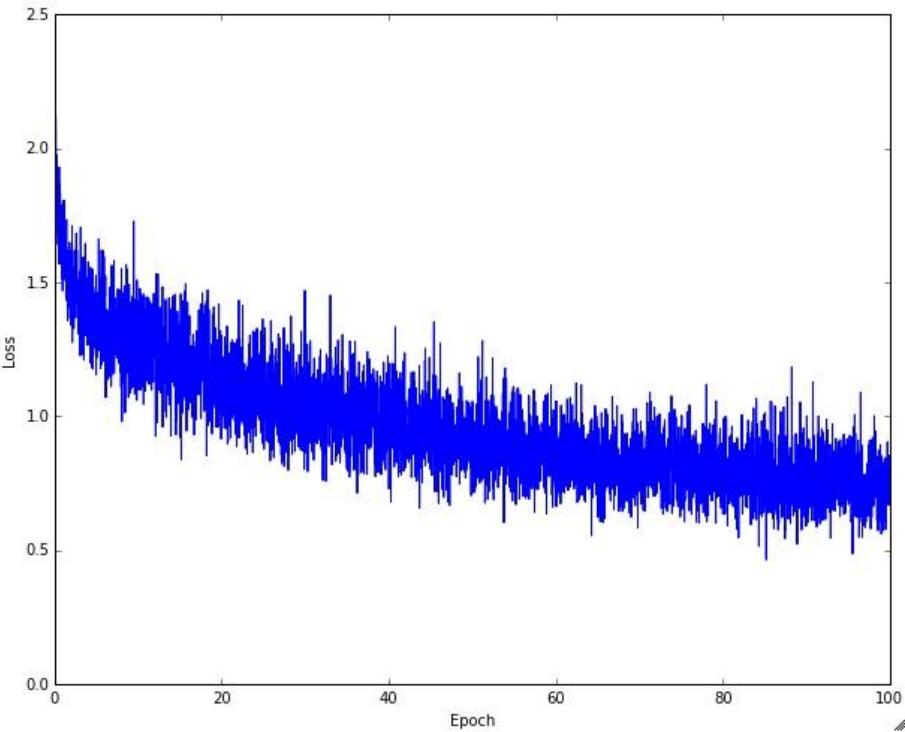
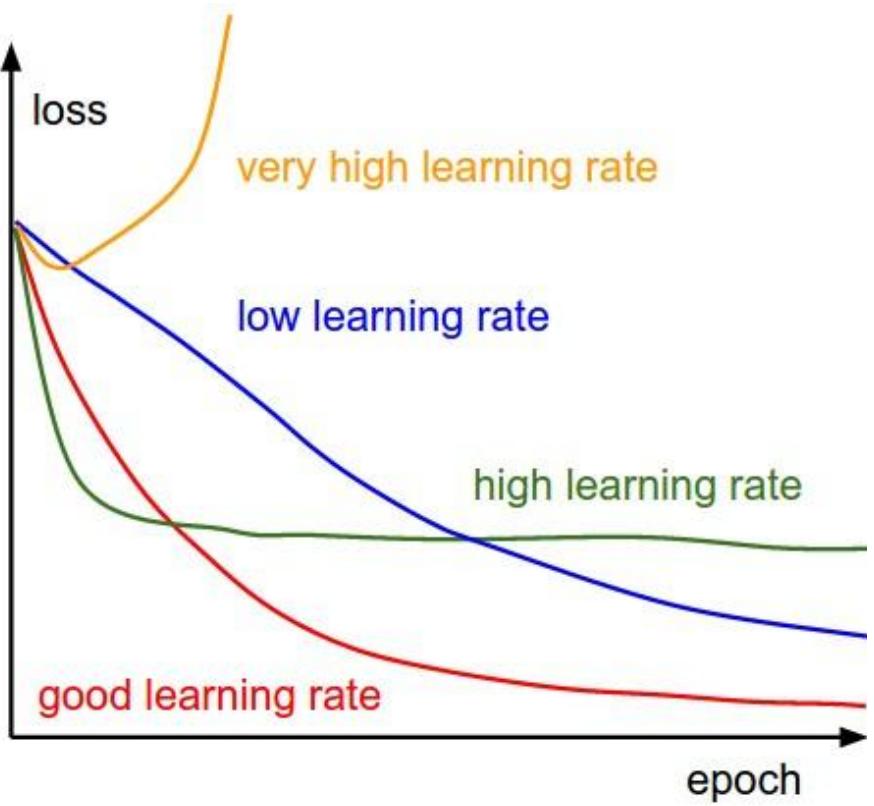


4. 5 LAYER FULLY-CONNECTED WITH RELU & DECAY LEARNING RATE



What happens when you add a decay learning rate?

EFFECTS OF DIFFERENT LEARNING RATES



CHANGING LEARNING RATE

CONSTANT LEARNING RATE

```
learning_rate = 0.003
```

```
train_step =
```

```
tf.train.AdamOptimizer(learning_rate).minimize(cross_entropy)
```

DECAY LEARNING RATE

```
# the learning rate is: 0.0001 + 0.003 * (1/e)^(step/2000)
```

```
# i.e. exponential decay from 0.003, then dropping exponentially to 0.0001
```

```
lr = 0.0001 + tf.train.exponential_decay(0.003, step, 2000, 1/math.e)
```

```
train_step = tf.train.AdamOptimizer(lr).minimize(cross_entropy)
```

FULLY-CONNECTED NEURAL NETWORK: RELU DECAY LEARNING RATE

Without
decay
learning rate

Max test
accuracy

0.8893

Decay
learning rate

0.8949

Cons

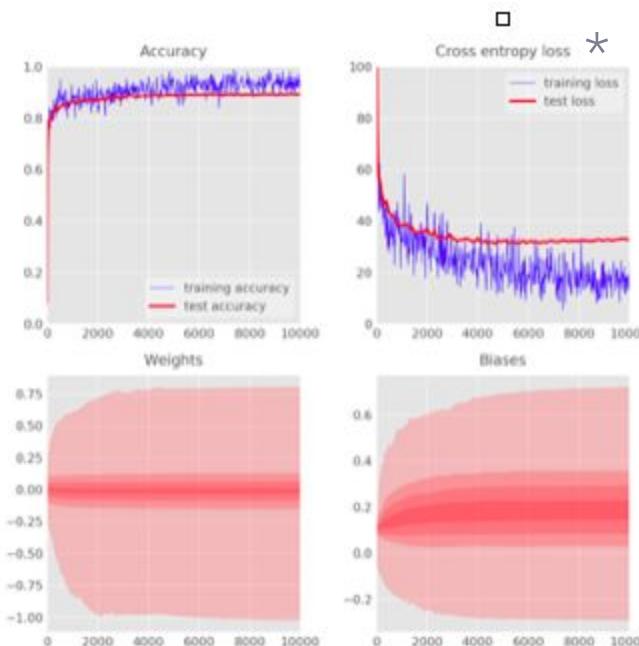
Cross-entropy
loss function
shifting up

Risk of
overfitting*

Pros

Faster
training

Improved
accuracy



5. 5 LAYER FULLY-CONNECTED WITH RELU + DECAY LEARNING + RATE + DROPOUT



What happens when you shoot 25% of your neurons?

WHAT IS DROPOUT?

- Dropout refers to ignoring units (i.e. neurons) during the training phase of certain set of neurons which is chosen at random.
- These units are not considered during a particular forward or backward pass.
- At each training stage, nodes are:
 - dropped with probability = $1 - p$
 - Kept with probability = p

WHY DROPOUT?

“To prevent overfitting”

CODE THE DROPOUT

```
# the previous model  
XX = tf.reshape(X, [-1, 784])  
Y1 = tf.nn.relu(tf.matmul(XX, W1) + B1)  
Y2 = tf.nn.relu(tf.matmul(Y1, W2) + B2)  
Y3 = tf.nn.relu(tf.matmul(Y2, W3) + B3)  
Y4 = tf.nn.relu(tf.matmul(Y3, W4) + B4)  
Ylogits = tf.matmul(Y4, W5) + B5  
Y = tf.nn.softmax(Ylogits)
```

```
# The model, with dropout at each layer
```

```
XX = tf.reshape(X, [-1, 28*28])
```

```
Y1 = tf.nn.relu(tf.matmul(XX, W1) + B1)
```

```
Y1d = tf.nn.dropout(Y1, pkeep)
```

```
Y2 = tf.nn.relu(tf.matmul(Y1d, W2) + B2)
```

```
Y2d = tf.nn.dropout(Y2, pkeep)
```

```
Y3 = tf.nn.relu(tf.matmul(Y2d, W3) + B3)
```

```
Y3d = tf.nn.dropout(Y3, pkeep)
```

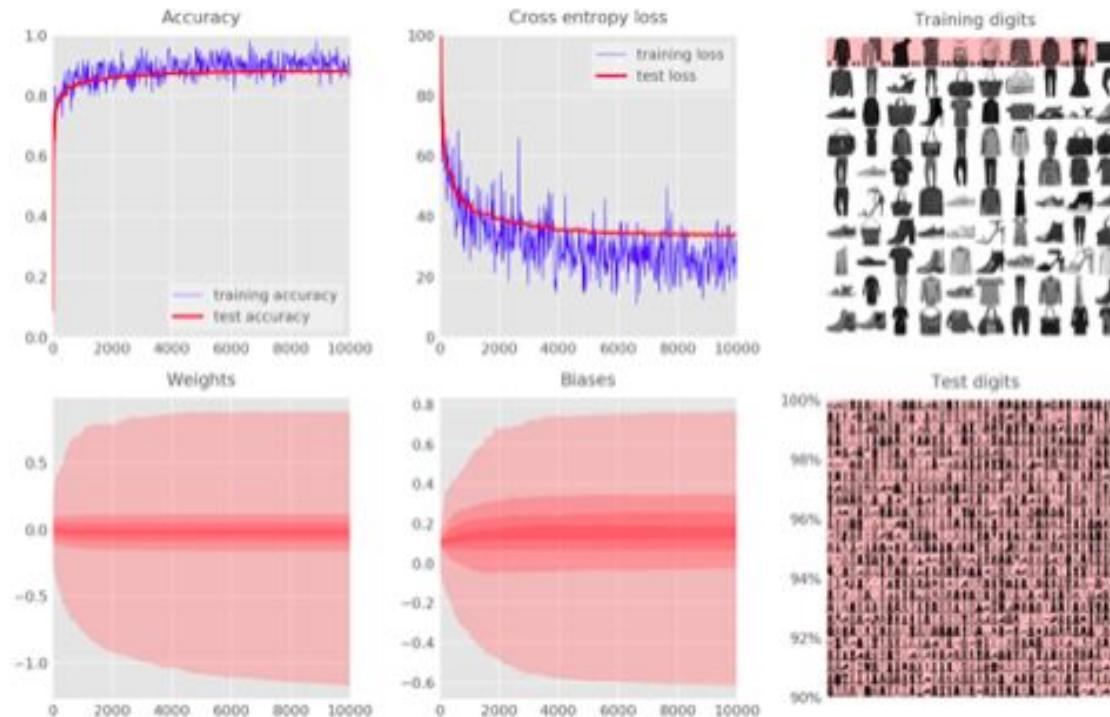
```
Y4 = tf.nn.relu(tf.matmul(Y3d, W4) + B4)
```

```
Y4d = tf.nn.dropout(Y4, pkeep)
```

```
Ylogits = tf.matmul(Y4d, W5) + B5
```

```
Y = tf.nn.softmax(Ylogits)
```

FULLY-CONNECTED NEURAL NETWORK WITH RELU + DECAY LEARNING RATE + DROPOUT



max test accuracy: 0.882

QUESTION

What do you think will happen when we change
the learning rate and the number of layers?

COMPARING ACCURACY RATES FOR 5 LAYER RELU W/ DROPOUT

	Number of Layers	Number of Layers	Number of Layers
Exponential Decay	3	5	7
From .03 to .0001	0.814	0.7706	.1116
Default From .003 to 0.0001	0.8832	0.882	0.8886
From .0003 to .0001	0.8692	0.8731	0.8891
	32		

“

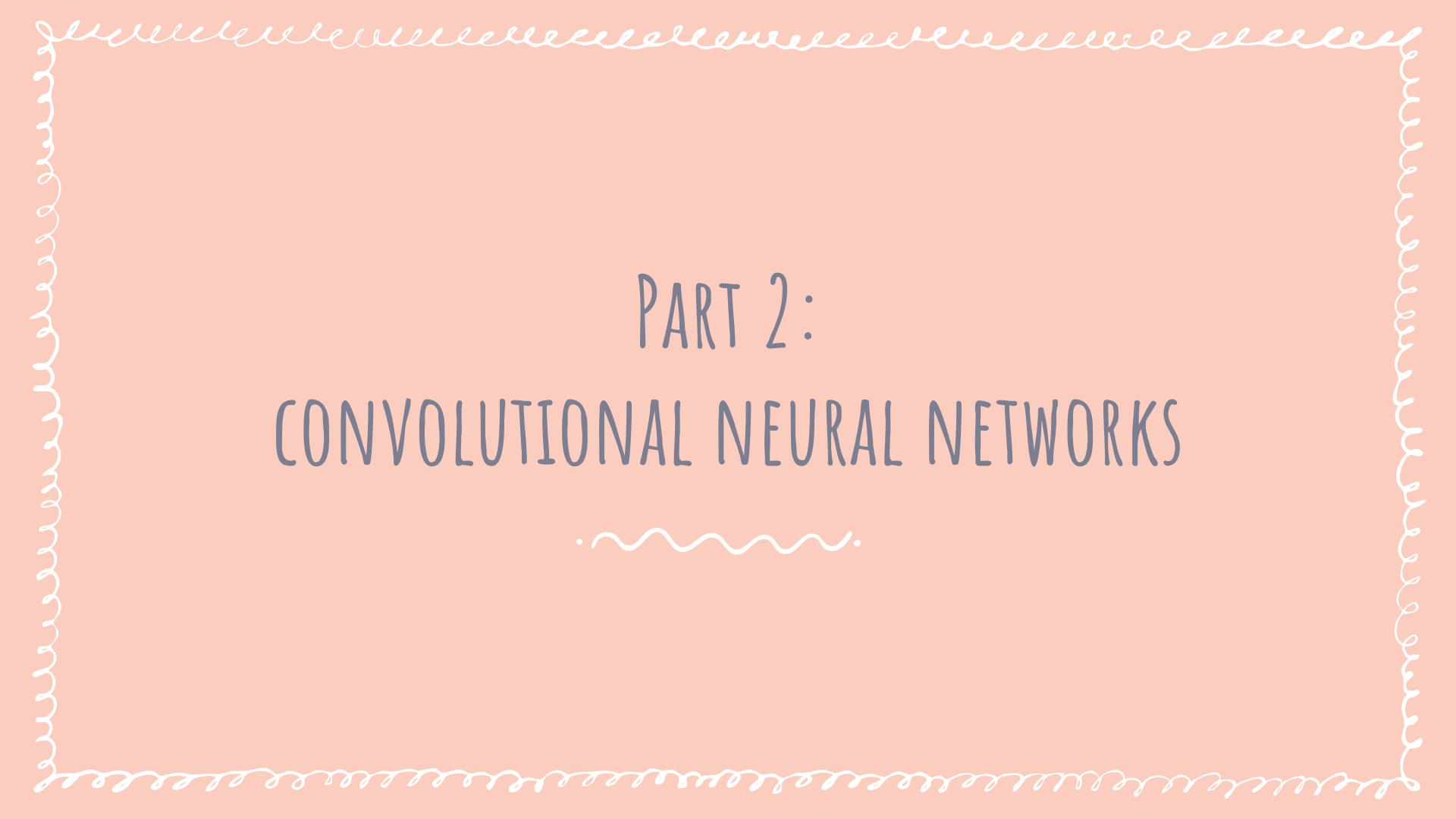
88.9% test accuracy at most?!
We can do better than that!

—••—

CNNs

What happens if we maintain the 2-D structure of images?

Next time, we'll discuss Convolutional Neural Networks!



PART 2: CONVOLUTIONAL NEURAL NETWORKS



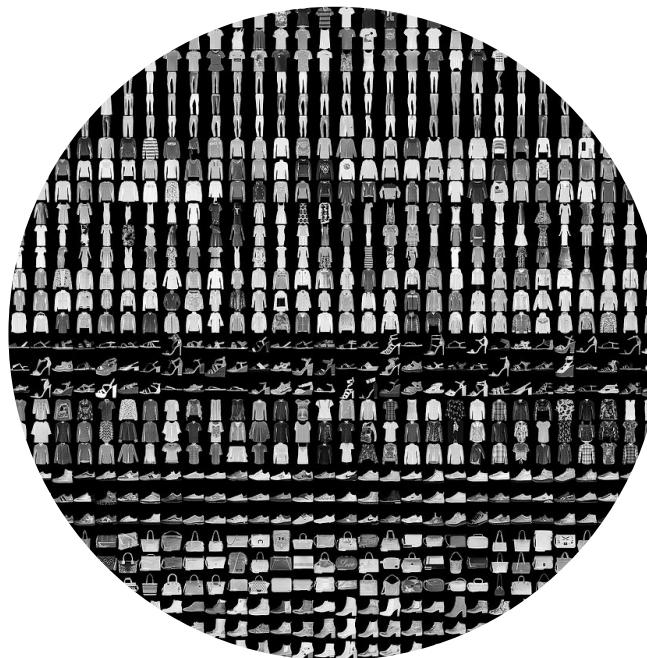
OBJECTIVES

1. Brief review of Fashion-MNIST data set
2. History of CNNs
3. Add layers! More degrees of freedom!
4. Add dropout (regularization)
 - a. $P = 0.75$
 - b. Change that up!
5. CNN with Batch Normalization

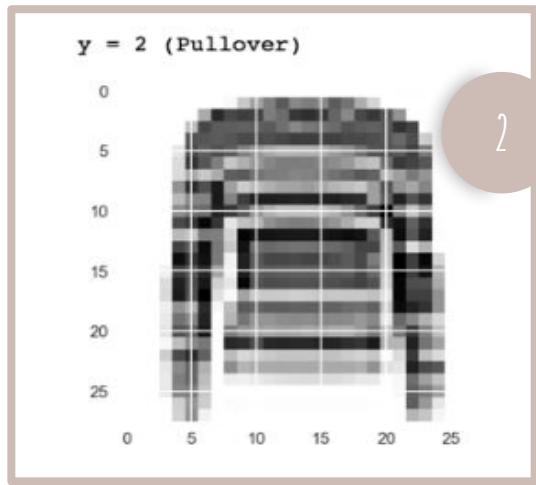
FASHION-MNIST

Zalando Research created:

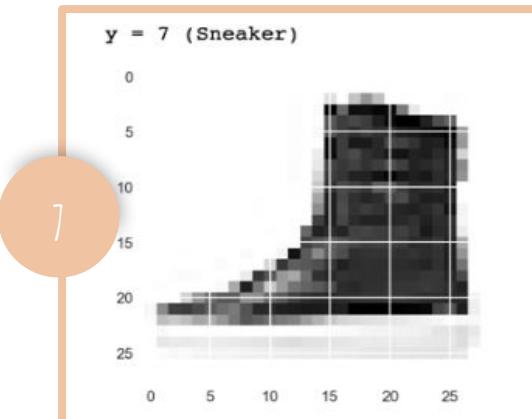
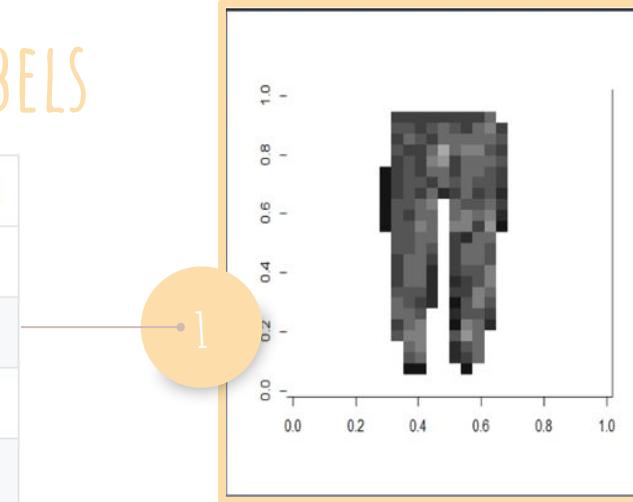
- 60,000 training images
- 10,000 test images
- 28x28 gray-scale image
 - 784 pixels



10 DIFFERENT LABELS



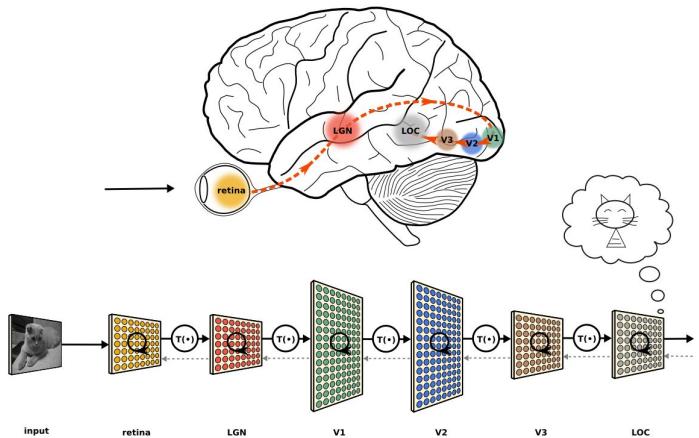
Label	Description
0	T-shirt/top
1	Trouser
2	Pullover
3	Dress
4	Coat
5	Sandal
6	Shirt
7	Sneaker
8	Bag
9	Ankle boot



2. HISTORY OF CONVOLUTIONAL NEURAL NETWORKS



BIOLOGICAL BACKGROUND



Research paper written by D.H. Hubel and T.N. Wiesel in 1962 on mammalian visual cortex.

They proposed an explanation for the way in which mammals visually perceive the world around them using a layered architecture of neurons in the brain.

This inspired engineers to attempt to develop similar pattern recognition mechanisms in computer vision.

3. CONVOLUTIONAL NEURAL NETWORKS



Let's utilize the 2-D structure of images.

CONVOLUTION

1 x1	1 x0	1 x1	0	0
0 x0	1 x1	1 x0	1	0
0 x1	0 x0	1 x1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

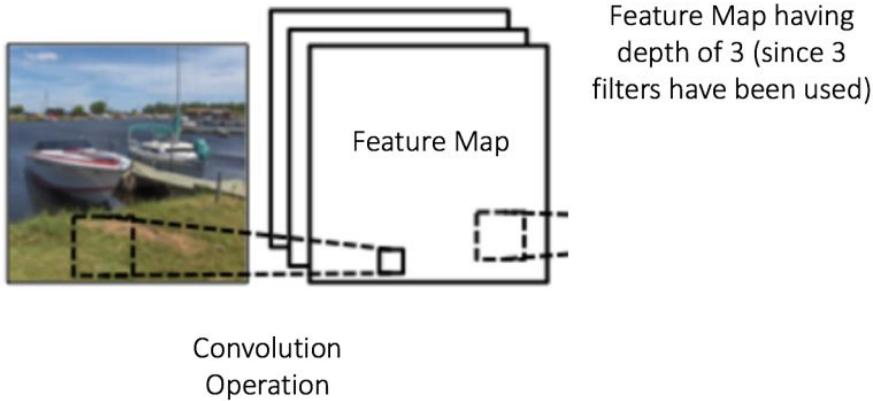
Convolved
Feature

- Preserves the spatial relationship between pixels
- Learn features using small squares of input data

4. ADD LAYERS- DEGREES OF FREEDOM



ADD LAYERS



- **Depth** corresponds to the number of filters we use for the convolution operation.
- **Stride** is the number of pixels by which we slide our filter matrix over the input matrix. When the stride is 1 then we move the filters one pixel at a time.

The more number of filters we have, the more image features get extracted and the better our network becomes at recognizing patterns in unseen images.

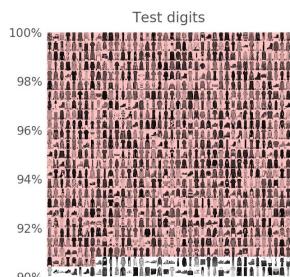
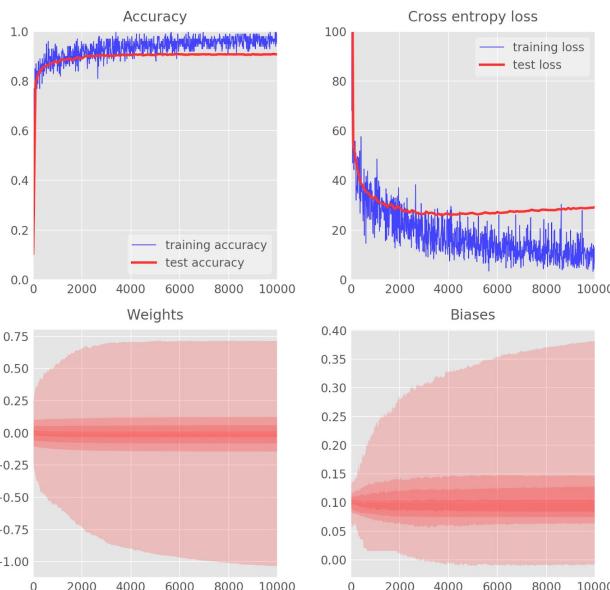
3 CONVOLUTIONAL LAYERS 2 FULLY CONNECTED

```
# The model
stride = 1 # output is 28x28
Y1 = tf.nn.relu(tf.nn.conv2d(X, W1, strides=[1, stride, stride, 1], padding='SAME') + B1)
stride = 2 # output is 14x14
Y2 = tf.nn.relu(tf.nn.conv2d(Y1, W2, strides=[1, stride, stride, 1], padding='SAME') + B2)
stride = 2 # output is 7x7
Y3 = tf.nn.relu(tf.nn.conv2d(Y2, W3, strides=[1, stride, stride, 1], padding='SAME') + B3)

# reshape the output from the third convolution for the fully connected layer
YY = tf.reshape(Y3, shape=[-1, 7 * 7 * M])

Y4 = tf.nn.relu(tf.matmul(YY, W4) + B4)
Ylogits = tf.matmul(Y4, W5) + B5
Y = tf.nn.softmax(Ylogits)
```

RESULTS

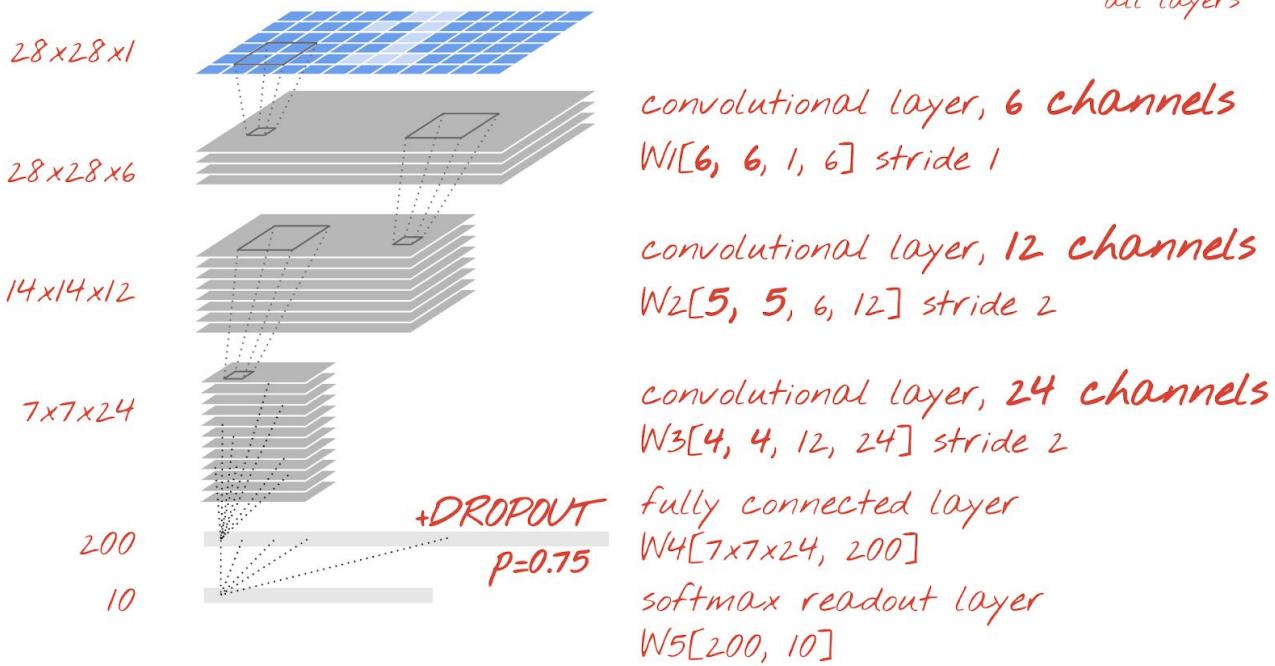


Model	Max Accuracy
Fully Connected 7 Layers + Relu+Dropout	88.91%
Convolutional Network	90.91%

5. CONVOLUTIONAL NEURAL NETWORKS WITH DROPOUT

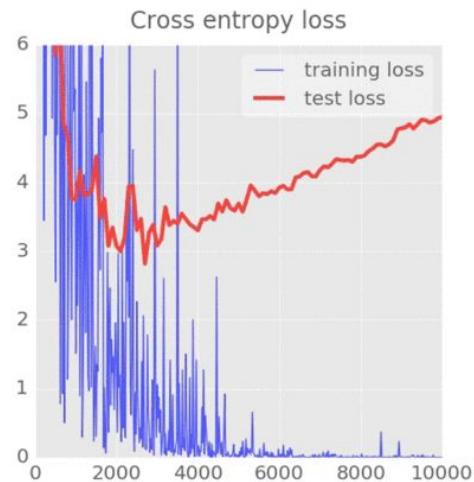
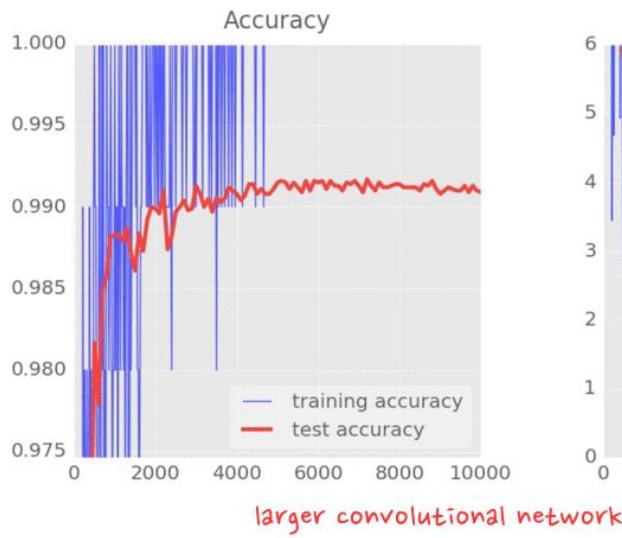


Fix Cross Entropy Loss



Increase number of channels to increase degrees of freedom but also add dropout on fully connected layer so increase in degrees of freedom doesn't result in overfitting.

RESULTS



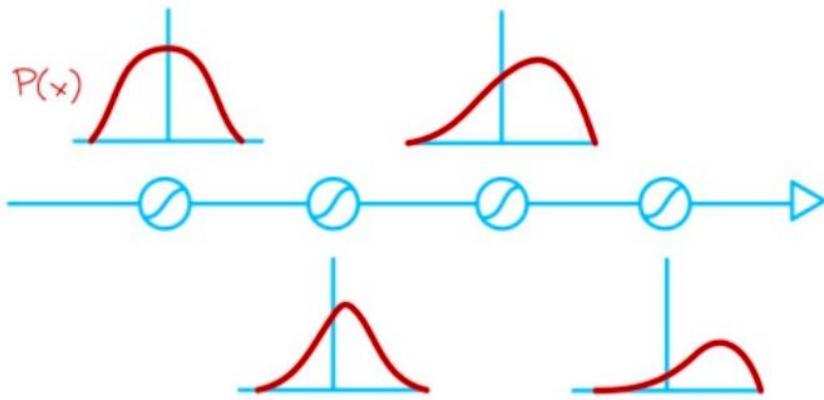
Model	Max Accuracy
Convolutional Network	90.91%
Convolutional Network + Dropout	91.16%

6. BATCH NORMALIZATION



Speed up training!

INTERNAL COVARIATE SHIFT

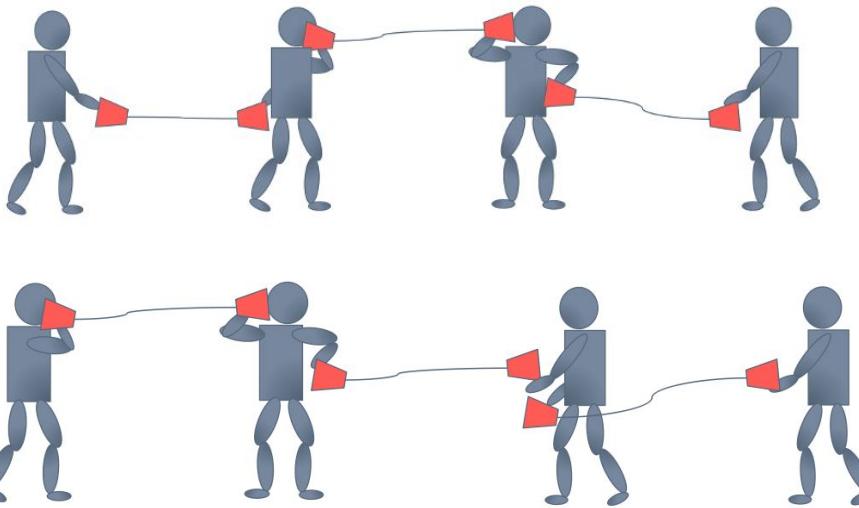


As learning progresses, the distribution of layer inputs changes due to parameter updates.

The weights in our model become updated over each epoch during training.

What if during training one of the weights ends up becoming drastically larger than the other weights?

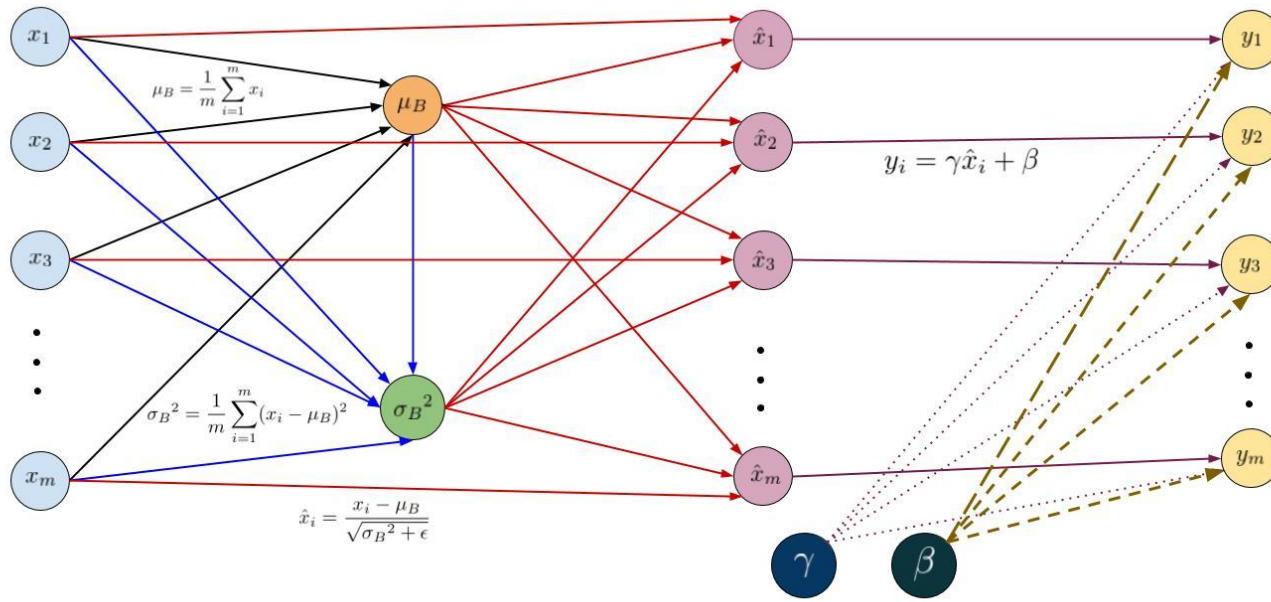
PROBLEM: INTERNAL COVARIATE SHIFT



Well this large weight will cause the output from its corresponding neuron to be extremely large.

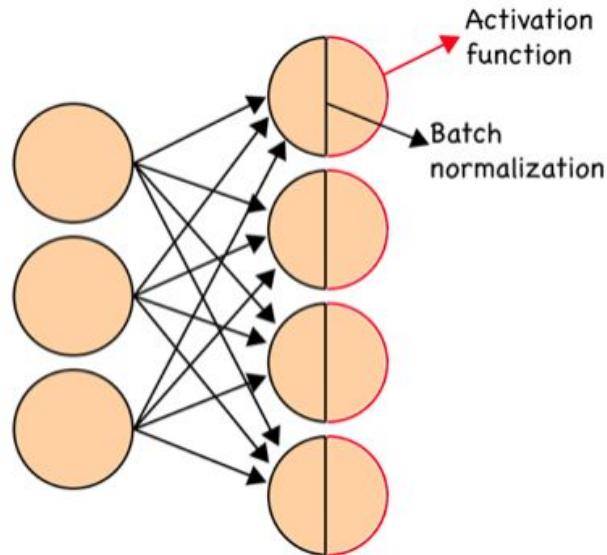
The imbalance will cascade through the neural network causing instability that slows down convergence.

SOLUTION: BATCH NORMALIZATION



1. Compute average & variance on mini-batch

BATCH NORMALIZATION



$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)$$

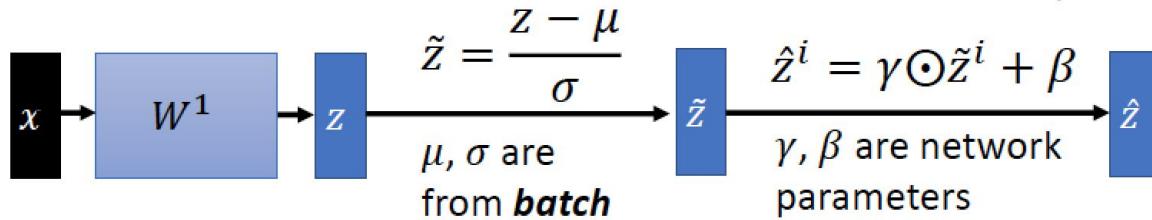
2. Center & rescale logits before the activation function
3. Add learnable scale & offset

WHY USE BATCH NORMALIZATION?

- Helps maximize ability to take advantage of activation function
- Makes training faster
- Unbalanced non-normalized data can make it harder to train data
- Regularization is a side-benefit!

Batch normalization

- At testing stage:



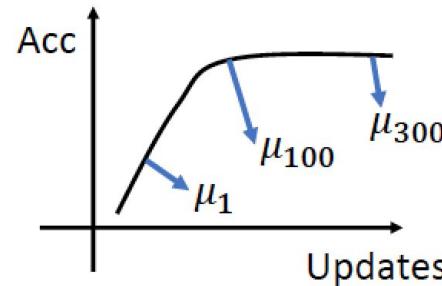
We do not have batch at testing stage.

Ideal solution:

Computing μ and σ using the whole training dataset.

Practical solution:

Computing the moving average of μ and σ of the batches during training.



CODE: BATCH NORMALIZATION

```
def batchnorm(Ylogits, is_test, iteration, offset, convolutional=False):
    # adding the iteration prevents from averaging across non-existing iterations
    exp_moving_avg = tf.train.ExponentialMovingAverage(0.999, iteration)
    bnepsilon = 1e-5
    if convolutional:
        mean, variance = tf.nn.moments(Ylogits, [0, 1, 2])
    else:
        mean, variance = tf.nn.moments(Ylogits, [0])
    update_moving_averages = exp_moving_avg.apply([mean, variance])
    m = tf.cond(is_test, lambda: exp_moving_avg.average(mean), lambda: mean)
    v = tf.cond(is_test, lambda: exp_moving_avg.average(variance), lambda: variance)
    Ybn = tf.nn.batch_normalization(Ylogits, m, v, offset, None, bnepsilon)
    return Ybn, update_moving_averages
```

BATCH NORMALIZATION MODEL

```
# batch norm offsets are used instead of biases
```

```
stride = 1 # output is 28x28
```

```
Y1l = tf.nn.conv2d(X, W1, strides=[1, stride, stride, 1], padding='SAME')
```

```
Y1bn, update_ema1 = batchnorm(Y1l, tst, iter, B1, convolutional=True)
```

```
Y1r = tf.nn.relu(Y1bn)
```

```
Y1 = tf.nn.dropout(Y1r, pkeep_conv, compatible_convolutional_noise_shape(Y1r))
```

```
stride = 2 # output is 14x14
```

```
Y2l = tf.nn.conv2d(Y1, W2, strides=[1, stride, stride, 1], padding='SAME')
```

```
Y2bn, update_ema2 = batchnorm(Y2l, tst, iter, B2, convolutional=True)
```

```
Y2r = tf.nn.relu(Y2bn)
```

```
Y2 = tf.nn.dropout(Y2r, pkeep_conv, compatible_convolutional_noise_shape(Y2r))
```

```
stride = 2 # output is 7x7
```

```
Y3l = tf.nn.conv2d(Y2, W3, strides=[1, stride, stride, 1], padding='SAME')
```

```
Y3bn, update_ema3 = batchnorm(Y3l, tst, iter, B3, convolutional=True)
```

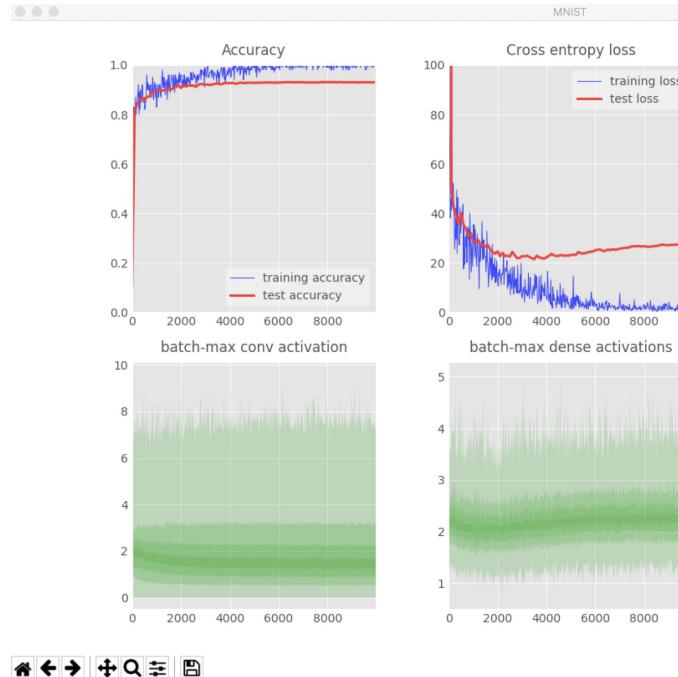
```
Y3r = tf.nn.relu(Y3bn)
```

```
Y3 = tf.nn.dropout(Y3r, pkeep_conv, compatible_convolutional_noise_shape(Y3r))
```

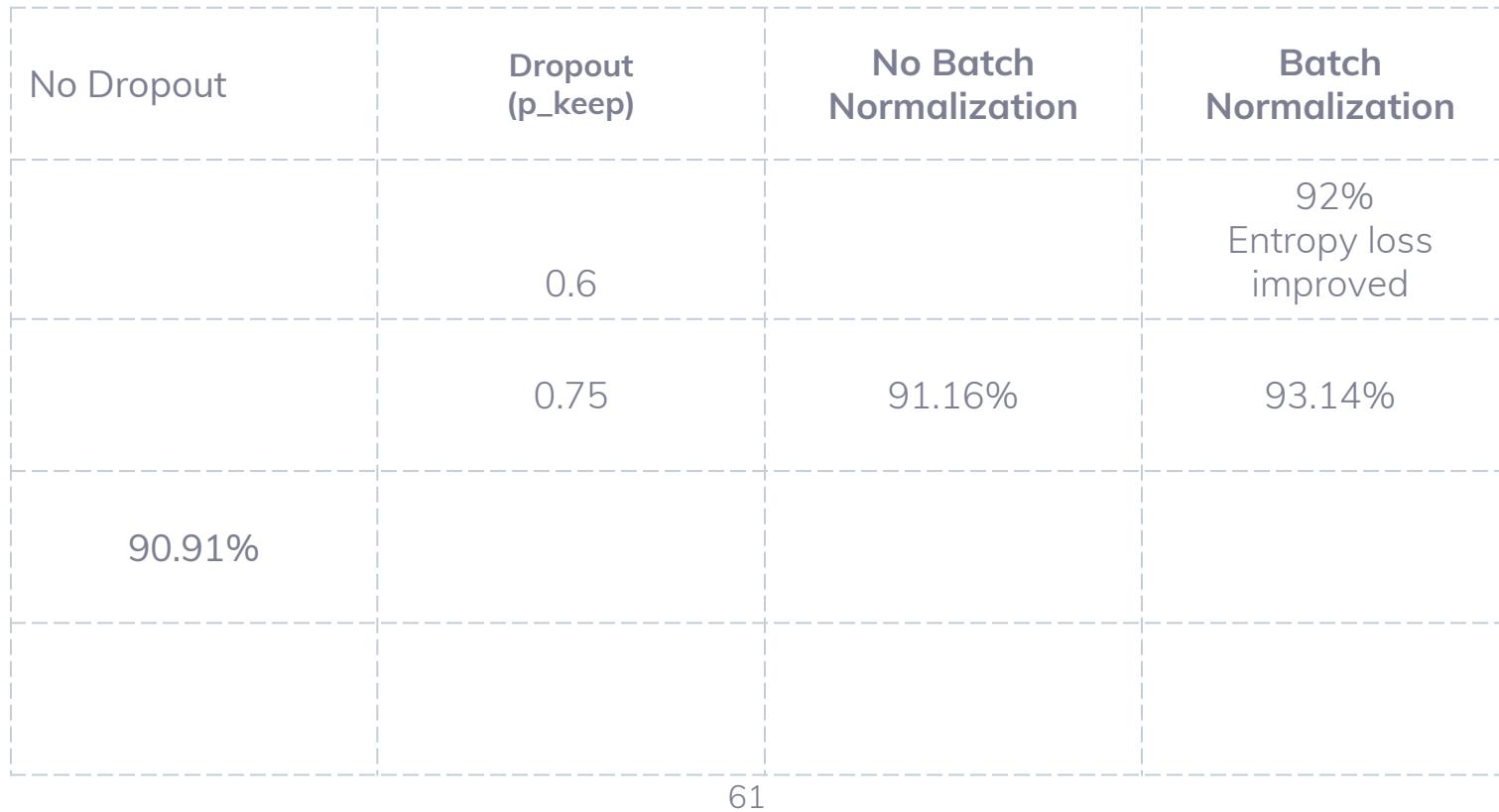
5 LAYER CNN WITH BATCH NORMALIZATION RESULTS

Dropout	Batch Normalization
Max test accuracy	91.16%

Cons	Pros
Noisy accuracy curve	Faster training
Noisy cross entropy loss	Improved accuracy



COMPARING ACCURACY RATES OF MODELS WE TESTED



HOW DO OUR MODELS COMPARE TO FASHION-MNIST BENCHMARKS?

Model	Accuracy Rate	Author
CNN w/ Batch Norm	93.14%	Stephanie & Fanny
Xgboost	95.3%	@anktplwl91
Resnet18	94.9%	<u>Kyriakos Efthymiadis</u>
Alexnet with triple loss	89.99%	<u>@Cenk Bircanoğlu</u>

THANKS! QUESTIONS?



SOURCES

- [https://www.tensorflow.org/versions/r1.1/get_st
arted/mnist/beginners](https://www.tensorflow.org/versions/r1.1/get_started/mnist/beginners)
- [http://colah.github.io/posts/2014-10-Visualizing-
MNIST/](http://colah.github.io/posts/2014-10-Visualizing-
MNIST/)
- [https://www.youtube.com/watch?v=qyvlt7kiQoI
&feature=youtu.be](https://www.youtube.com/watch?v=qyvlt7kiQoI
&feature=youtu.be)
- [https://cloud.google.com/blog/big-data/2017/01/
learn-tensorflow-and-deep-learning-without-a-
phd](https://cloud.google.com/blog/big-data/2017/01/
learn-tensorflow-and-deep-learning-without-a-
phd)