INTRODUCTION TO TENSORFLOW for Deep Learning



Mara Sorella @maruscia

Web Algorithmics and Data Mining Lab

TensorFlow

TensorFlow[™] is an open source software library for numerical computation using **data flow graphs**.

Originally developed (for internal use) by researchers and engineers working on the **Google Brain Team** for the purposes of conducting **machine learning** and **deep neural networks research.**

It was released under the Apache 2.0 open source license on November 2015.





What is a Tensor?

Tensor

a typed multi-dimensional array.

i.e. a mini-batch of images

a 4-D array of floating point numbers with dimensions [batch_size, height, width, channels]

Operation

takes zero or more Tensors, performs some computation, and produces zero or more Tensors

Data flow graph

- Nodes in the graph represent operations
- Edges are directed and represent passing the result of an operation (a tensor) as an operand to another operation



h = ReLU(Wx + b)



TensorFlow architecture



TensorFlow consists of the following two components:

• a graph prototype

a computation specification.

• a runtime that executes the (distributed) graph



TensorFlow architecture



TensorFlow consists of the following two components:

• a graph prototype

a computation specification.

• a runtime that executes the (distributed) graph



MNIST: Handwritten digits recognition

Righteously identified as the Hello World of machine learning.



60K labeled examples

 28x28 pixel grayscale image



We know that every image in MNIST is a handwritten digit between zero and nine: 10 classes



We know that every image in MNIST is a handwritten digit between zero and nine: **10 classes**





We know that every image in MNIST is a handwritten digit between zero and nine: **10 classes**





We know that every image in MNIST is a handwritten digit between zero and nine: **10 classes**



0	1.	 5	6	 9	
0.02	0.01	0.5	0.45	0.01	



We know that every image in MNIST is a handwritten digit between zero and nine: **10 classes**





We know that every image in MNIST is a handwritten digit between zero and nine: 10 classes





Softmax regression for MLC



$$\mathbf{L}_{i} = \sum_{j} W_{i,j} x_{j} + b_{i} \qquad softmax(L_{i}) = \frac{e^{L_{i}}}{\|e^{L}\|}$$

Our promise: build a neural network that achieves 99% accuracy in this task



1-Layer simple NN for MLC



Activation function: a function (in this case softmax, yet also sigmoid, ReLu,...) that takes in the weighted sum of all of the inputs from the previous layer and then generates and passes an output value (typically nonlinear) to the next layer.



Working on a batch



$$L = X.W + b$$



Working on a batch





Working on a batch





Softmax over a batch of images







A loss function for MLC

- MNIST: 60K images train, 10K test 6 7 Labeled instances actual probabilities, "one-hot" encoded $-\sum Y'_i . log(Y_i)$ Cross Entropy = this is a "6" computed probabilities **NN output** 0.1||0.2||0.1||0.3||0.2||0.1/|0.9| 0.2 0.1 0.1 Training minimize CE, across
 - all training examples



Demo: mnist_1.0_softmax.py





0.2

0.0

0.4

0.8

1.0

0.6

Training images (current batch)







Demo: mnist_1.0_softmax.py





0.2

0.0

0.4

0.8

1.0

0.6

Training images (current batch)







TensorFlow: init





Training = computing variables W and b



Model

% 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))



Optimizer

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

Loss function

learning rate



TF: full python code





Kun

Go deep! 5-layer NN



Note: 5 layers are overkill for this task



K = 200

- L = 100
- M = 60

```
N = 30
```

```
W1 = tf.Variable(tf.truncated_normal([28*28, K] ,stddev=0.1))
B1 = tf.Variable(tf.zeros([K]))
```

```
weights initialised
with random values
```



W2 = tf.Variable(tf.truncated_normal([K, L], stddev=0.1)) B2 = tf.Variable(tf.zeros([L]))

```
W3 = tf.Variable(tf.truncated_normal([L, M], stddev=0.1))
B3 = tf.Variable(tf.zeros([M]))
W4 = tf.Variable(tf.truncated_normal([M, N], stddev=0.1))
B4 = tf.Variable(tf.zeros([N]))
W5 = tf.Variable(tf.truncated_normal([N, 10], stddev=0.1))
B5 = tf.Variable(tf.zeros([10]))
```



Model

weights and biases
X = tf.reshape(X, [-1, 28*28])
...yet with a different activation function
Y1 = tf.nn.sigmoid(tf.matmul(X, 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)
Y = tf.nn.softmax(tf.matmul(Y4, W5) + B5)



Demo: mnist_2.0_five_layers_sigmoid.py





Training images (current batch)



Test images 100% 7.2.4.5.5.3.3.4.5.6.4.3.3.4.7.2.7.5.7.4.4.5.1.4.4.6.3.5.5.6.4.7.8.7.8.1.8 ₩721#24682#7361#+65#248739+443575566838 #62388655562486596548484253521\6#8635566838 6759382113287714559714559735444.57712855834115 144582242424247293673512348472537863+8595 98% 55-215056372544337662428525442855137 2244-51363120649513131554663127584 \$P\$237.4233.3557.8436571/3.8344642545942 2433-2223593251-615299456439995731293 96% 1295549145545246433269263462542 * CH2+584242552134584584531652634F 672729324356733535526K92**8**3522562 \$5E: 11, \$247, 33, 33, 34, 55, 84, 59, 67, 68, 321, 8. 644546633266931029289472639423 94% 1,2*2,4,2,1,2,4,4*,2,2,2,3,5,5,5,5,5,7,4,4,9,5,5,5,5,2,3,8,3,9 \$374A35A7X)&85233361A547A7545454848444464 19245137224/6242517636267372362935 \$8523158668725892389384458759756 92% 476655723533775433623363384635673746356758 \$5954.8539217342255654+22554424554455447374 90%



Demo: mnist_2.0_five_layers_sigmoid.py





Training images (current batch)



Test images 100% 7.2.4.5.5.3.3.4.5.6.4.3.3.4.7.2.7.5.7.4.4.5.1.4.4.6.3.5.5.6.4.7.8.7.8.1.8 ₩721#24682#7361#+65#248739+443575566838 #62388655562486596548484253521\6#8635566838 6759382113287714559714559735444.57712855834115 144582242424247293673512348472537863+8595 98% 55-215056372544337662428525442855137 2244-51363120649513131554663127584 \$P\$237.4233.3557.8436571/3.8344642545942 2433-2223593251-615299456439995731293 96% 1295549145545246433269263462542 * CH2+584242552134584584531652634F 672729324356733535526K92**8**3522562 \$5E: 11, \$247, 33, 33, 34, 55, 84, 59, 67, 68, 321, 8. 644546633266931029289472639423 94% 1,2*2,4,2,1,2,4,4*,2,2,2,3,5,5,5,5,5,7,4,4,9,5,5,5,5,2,3,8,3,9 \$374A35A7X)&85233361A547A7545454848444464 19245137224/6242517636267372362935 \$8523158668725892389384458759756 92% 476655723533775433623363384635673746356758 \$5954.8539217342255654+22554424554455447374 90%



Slow convergence



Vanishing gradients problem

- 1. If the gradient at each step is too small, then greater repetitions will be needed until convergence, because the weight is not changing enough at each iteration.
- 2. In **deep** networks, computing these gradients can involve taking the product of many small terms. and when you do so repeatedly, neuron outputs and their gradients can vanish entirely starting from the first layers (close to input).

The **ReLU** activation function can help prevent vanishing gradients.







How to use ReLU

Simply swap tf.nn.sigmoid with tf.nn.relu in code.

Biases initialization with small positive weights (just for ReLUs)

W = tf.Variable(tf.truncated_normal([K,L], stddev=0.1))

B = tf.Variable(tf.ones([L])/10)



In very high dimensional spaces like here - we have in the order of 10K weights and biases - "saddle points" are frequent.

These are points that are not local minima but where the gradient is nevertheless zero and the gradient descent optimizer stays stuck there.

TensorFlow has a **full array** of available optimizers, including **some that work with an amount of inertia and will safely sail past saddle points**.

e.g.,

You can replace tf.train.GradientDescentOptimiser with a tf.train.AdamOptimizer



From sigmoid to ReLU

First 300 iterations





From sigmoid to ReLU

First 300 iterations





300
Noisy accuracy







Learning rate decay

If we reduce the learning rate, say by a factor of 10, we might slow down our learning by a factor of 10.

Right approach

Start fast, then decay.



Learning rate decay

If we reduce the learning rate, say by a factor of 10, we might slow down our learning by a factor of 10.

Right approach

Start fast, then decay.



Fixed learning rate Ir = 0.003 (no decay)



Learning rate decay

If we reduce the learning rate, say by a factor of 10, we might slow down our learning by a factor of 10.

Right approach

Start fast, then decay.



Start with Ir=0.003 then drop exponentially to 0.0001

lr = lrmin+(lrmax-lrmin)*exp(-i/2000)



Overfitting?



The system is getting better and better on training data but performance on test data is not improving anymore *ML handbook* solution: **regularization**



Dropout: shooting neurons



At each iteration:

drop fraction **(1-pkeep)** of neurons (with all their weights and biases)

(boosting the output of remaining neurons)



When testing, all neurons are put back in

effect: those weights will not change for that iteration

TF code

```
pkeep = tf.Placeholder(float32)
Yf = tf.nn.relu(tf.matmul(X, W) + B)
Y = tf.nn.dropout(Yf, pkeep)
```

applied after each layer





Sigmoid, learning rate = 0.003





RELU, learning rate = 0.003





RELU, decaying learning rate 0.003 -> 0.0001





RELV, decaying learning rate 0.003 -> 0.0001 and dropout 0.75





Overfitting: the system is getting better and better on training data but performance on test data is not improving anymore.

At its core overfitting happens when you have **too many degrees of freedom** (weights and biases) for the problem at hand.

Questions:

- 1. too few data? (no)
- 2. too many layers? (maybe)
- 3. faulty design?

At the present state, apparently by some faulty design our NN is not capable to constrain the d.o.f. and **extract the structure** we need. In other words it cannot learn *categories* in a sufficient way to generalize well to data it has never seen before.





Overfitting: the system is getting better and better on training data but performance on test data is not improving anymore.

At its core overfitting happens when you have **too many degrees of freedom** (weights and biases) for the problem at hand.

Questions:

- 1. too few data? (no)
- 2. too many layers? (maybe)
- 3. faulty design?

At the present state, apparently by some faulty design our NN is not capable to constrain the d.o.f. and **extract the structure** we need. In other words it cannot learn *categories* in a sufficient way to generalize well to data it has never seen before.

Why?





Overfitting: the system is getting better and better on training data but performance on test data is not improving anymore.

At its core overfitting happens when you have **too many degrees of freedom** (weights and biases) for the problem at hand.

Questions:

- 1. too few data? (no)
- 2. too many layers? (maybe)
- 3. faulty design?

At the present state, apparently by some faulty design our NN is not capable to constrain the d.o.f. and **extract the structure** we need. In other words it cannot learn *categories* in a sufficient way to generalize well to data it has never seen before.

Why?

Hint: We did something very stupid





Overfitting: the system is getting better and better on training data but performance on test data is not improving anymore.

At its core overfitting happens when you have **too many degrees of freedom** (weights and biases) for the problem at hand.

Questions:

- 1. too few data? (no)
- 2. too many layers? (maybe)
- 3. faulty design?

At the present state, apparently by some faulty design our NN is not capable to constrain the d.o.f. and **extract the structure** we need. In other words it cannot learn *categories* in a sufficient way to generalize well to data it has never seen before.

Why?

Hint: We did something very stupid with the input.



Remember how we are using our images, all pixels flattened into a single long vector ? That was **a really bad idea**.

Handwritten digits are made of **shapes** and we discarded the shape information when we flattened the pixels. However, there is a type of neural network that can take advantage of shape information: **convolutional networks**.

These network are specifically designed for > 1D datasets (images) where shape information (locality) is important.



















Convolutional Neural Networks





Convolutional Neural Networks





Convolutional Neural Networks





Piping information down



Traditional approach: "max pooling"



subsample the neuron outputs by keeping only outputs where the signal is strongest

An alternative way that allows using only convolutional layers:Directly slide the patches across the image with a stride of2 pixels instead of 1: this allows to retain less outputs

This approach has proven just as effective and today's convolutional networks use convolutional layers only (see next slide)



Piping information down



Traditional approach: "max pooling"



subsample the neuron outputs by keeping only outputs where the signal is strongest

An alternative way that allows using only convolutional layers:Directly slide the patches across the image with a stride of2 pixels instead of 1: this allows to retain less outputs

This approach has proven just as effective and today's convolutional networks use convolutional layers only (see next slide)



A CNN for MNIST (all convolutional)

+ biases on all layers

28×28×1

28x28x4

14×14×8

7x7x12

200

10



convolutional layer, 4 channels WI[5, 5, 1) 4] stride D

convolutional layer, 8 channels W2[4, 4, 4, 8] stride 2

convolutional layer, 12 channels W3[4, 4, 8, 12] stride 2

fully connected layer W4[7x7x12, 200] softmax readout layer W5[200, 10]



CNN in TF (Init)

	filter	input	output	= number of patches we	
K=4	Size	channels	channels	produce in output	
L=8					
M=12			/		
			/		
<pre>W1 = tf.Variable(tf.truncated_normal([5, 5, 1, K],stddev=0.1))</pre>					
<pre>B1 = tf.Variable(tf.ones(</pre>	[K])/10)				
W2 = tf.Variable(tf.trunc	ated_normal	([5, 5, K, L],stddev=0	.1))	
<pre>B2 = tf.Variable(tf.ones(</pre>	[L])/10)				
W3 = tf. <pre>Variable(tf.trunc</pre>	ated_normal	([4, 4, L, M],stddev=0	.1))	
B3 = tf. <pre>Variable(tf.ones(</pre>	[M])/10)				
				weights initialised	
N=200				with random value	5
W4 = tf Variable(tf trunc)	ated normal([7*7*Μ N] c	tddev=0 1))		
B4 = tf.Variable(tf.cruck)	$\left[N \right] \left(10 \right)$	כן באו נויו לי לב			
$U_{\rm T} = t_{\rm T} \cdot v_{\rm all} t_{\rm able}(t_{\rm T} \cdot t_{\rm bles})$	// TO/		$\left(0, -0, 1 \right)$		
WS = TT.Variable(TT.Trunca	<pre>iced_normal(</pre>	[N, LO] ,STOC	ev=0.1))		
B5 = tf.Variable(tf.zeros)	([10])/10)				



CNN in TF (Model)

conv2d: you give it batch of input values (images), a set of weights, it will scan input in both directions and produce the weighted sums





Demo: mnist_3.0_convolutional.py





Biases range

6000

10000

8000

Training images (current batch)



Test images



100% 65346139054518434303351113506+523998/722348 473496454014013134727131.544351244135560 419178137-243070281732971227847361369314 98% 1749:054512/146813911-41159:6710585665381 616457319182549858156034465485301447282 1818 8508 230111090 169186/ 139529159399 36551347:28417338870224159P730424145772 1208577918180301994192124155264134252040 96% 523+412402443:005(9(105))+34420+11715359 18601381051915800/851194622505868311 1802376162/92861452544283424603(77371719 3 929204114818+59868750031256/8332391768 05666782275896 841299/97599991052378910 94% 1395213136577226 + 5474303231434464215 154: \$40023232308984114901807:063549331332 780211065-3869638099686957860240221/4751 01042719,018229273541802054137(1:25803460 92% 12611437617517547691383361258511443103-01444855408213460+0613269264314625120521 1341054311779918402451184919124455363(45 6894153803211283440233317359632613667217 1424217961124817)45058/31077035527644283 62754082758688274720463213227360578347602 90%



Results -CNN





Results -CNN





WT H?



Cross entropy loss training loss test loss

Solutions?



Go bigger

28×28×1

28×28×6

14×14×12

7×7×24

200

10

сопvolutional WI[6, 6, 1, 6] patch size convolutional W2[5, 5, 6, 1, convolutional W3[4, 4, 12, 2

convolutional layer, 6 channels WI[6, 6, 1, 6] stride 1

convolutional layer, 12 channels W2[5, 5, 6, 12] stride 2

convolutional layer, 24 channels W3[4, 4, 12, 24] stride 2

fully connected layer W4[7x7x24, 200] softmax readout layer W5[200, 10]

why not apply dropout to previous layers?

.....



+ biases on

all layers

Go bigger

+ biases on all layers

28×28×6

14×14×12

7x7x24



convolutional layer, 6 channels WI[6, 6, 1, 6] stride 1 patch size

convolutional layer, 12 channels W2[5, 5, 6, 12] stride 2

convolutional layer, 24 channels W3[4, 4, 12, 24] stride 2

W4[7x7x24, 200] fully connected layer W5[200, 10] softmax readout layer

why not apply dropout to previous layers?



Go bigger

+ biases on all layers

28×28×1 28×28×6 14×14×12 7x7x24 +DROPOU 200 P=0.75 10

convolutional layer, 6 channels WI[6, 6, 1, 6] stride 1 patch size

convolutional layer, 12 channels W2[5, 5, 6, 12] stride 2

convolutional layer, 24 channels W3[4, 4, 12, 24] stride 2

fully connected layer W4[7x7x24, 200] softmax readout layer W5[200, 10]

why not apply dropout to previous layers?



Demo: mnist_3.1_convolutional_bigger_dropout.py





Improvements



larger convolutional network


Improvements



larger convolutional network





Our model now misses only 72 out of the 10,000 test digits, with 99.3% accuracy

The world record, which you can find on the MNIST website is around **99.7%**. We are only 0.4 percentage points away from it with our model built with 100 lines of Python / TensorFlow.

Another technique that we didn't mention is **batch normalization**, that can lead us up to 99.5 accuracy.



TF.Learn, TF High level APIs

Tensorflow higher-level APIs too called <u>tf.learn</u>.

Simple but constrained.

TensorFlow ™	Install	Develop	API r1.4	Deploy	Extend	Community	Versions	TFRC
 tf.contrib.learn 		Class DNNRegressor						
Overview								
BaseEstimator		Inherits From: Estimator Defined in tensorflow/contrib/learn/python/learn/estimators/dnn.py.						
binary_svm_head								
build_parsing_serving_input_	fn							
DNNClassifier		See the guide: Learn (contrib) > Estimators						
DNNEstimator								
DNNLinearCombinedClassif	er	A regressor for TensorFlow DNN models.						
DNNLinearCombinedEstima	tor							
DNNLinearCombinedRegres	sor	Example:						
DNNRegressor				_				
DynamicRnnEstimator		<pre>sparse_feature_a = sparse_column_with_hash_bucket() sparse_feature_b = sparse_column_with_hash_bucket()</pre>						
Estimator								
Evaluable	<pre>sparse_feature_a_emb = embedding_column(sparse_id_column=sparse_feature_a,</pre>							
evaluate)							
Experiment		<pre>sparse_feature_b_emb = embedding_column(sparse_id_column=sparse_feature_b,</pre>						
ExportStrategy					•••)			
extract_dask_data		estimator =	DNNRegress	or(
extract_dask_labels		feature_columns=[sparse_feature_a, sparse_feature_b], hidden_units=[1024, 512, 256])						
extract_pandas_data								
extract_pandas_labels		# Or estima		- Duraudina I	A de euro d'On trà		and the second second	
extract_pandas_matrix		# Ur estima	ator using tr vation	ne Proximal	LAdagradupti	lmizer optimize	r with	
Head		estimator =	DNNRegress	or(
infer		feature	e_columns=[s	parse_featu	ure_a, spars	se_feature_b],		
infer_real_valued_columns_f	al_valued_columns_fro hidden_units=[1024, 512, 256],							
infer_real_valued_columns_f	ro	optimizer=tf.train.ProximalAdagradUptimizer(
InputFnOps	11 regularization strength=0.001							
KMeansClustering))						
LinearClassifier								
LinearEstimator		# Input but	ilders					
LinearRegressor	LinearRegressor def input_fn_train: # returns x, y							
LogisticRegressor		estimator f	fit(input fn:	=input fn t	rain)			
make_export_strategy					,			

Example: DNNRegressor



References and pointers

Code

git clone <u>https://github.com/maruscia/tensorflow-mnist-tutorial</u>

- . largely (and aggressively) borrowed by work by Martin Görner
 - . https://cloud.google.com/blog/big-data/2017/01/learn-tensorflow-and-deep-learning-without-a-phd
- Deep Learning: Advanced machine learning course on neural networks, with extensive coverage of image and text models
- Rules of ML: Best practices for machine learning engineering
- deeplearn.js: Open-source toolkit for interactive model training and inference in the browser
- distill.pub: Web journal presenting recent research in ML & AI in a divulgative way (reactive diagrams, visualizations, mind diagrams)
- colah.github.io: Blog by Chris Olah, (Google Brain team, also editor at distill)
 - In particular check out this blog post for (quite old, yet still very informative) visualizations of the MNIST

TensorFlow resources

- TensorFlow Programmer's Guide: In-depth guides to key TensorFlow features, including variables, threading, and debugging
- TensorFlow Dev Summit 2017: Series of tech talks and demos highlighting TensorFlow APIs and real-world applications

