

# Deep Learning on Big Data Sets in the Cloud with Apache Spark and Google TensorFlow

December 9, 2016

Patrick GLAUNER and Radu STATE

first.last@uni.lu

SEDAN Lab,  
SnT - Interdisciplinary Centre for Security, Reliability and Trust,  
University of Luxembourg

- ▶ PhD Student at the University of Luxembourg
- ▶ Adjunct Lecturer of Artificial Intelligence at Karlsruhe University of Applied Sciences
- ▶ MSc in Machine Learning from Imperial College London
- ▶ Previously worked at CERN and SAP

## Definition (Artificial Intelligence)

"AI is the science of knowing what to do when you don't know what to do." (Peter Norvig)<sup>a</sup>

---

<sup>a</sup><http://www.youtube.com/watch?v=rtmQ3x1t-4A4m45>

## Definition (Machine Learning)

Machine Learning is the field of study that gives computers the ability to learn without being explicitly programmed.

# Motivation

Goal: recognition of characters

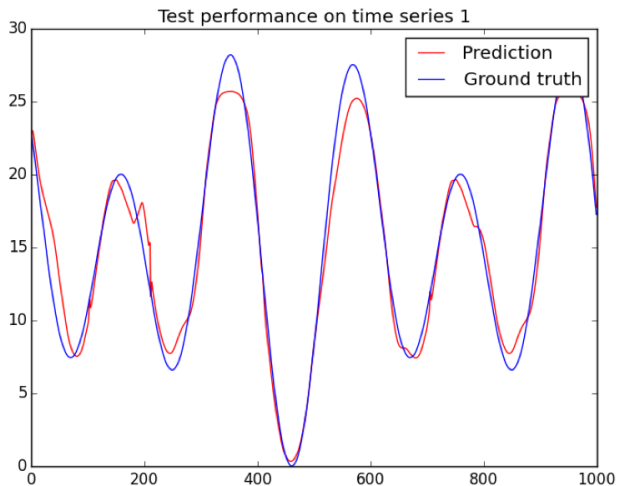


notMNIST examples<sup>1</sup>.

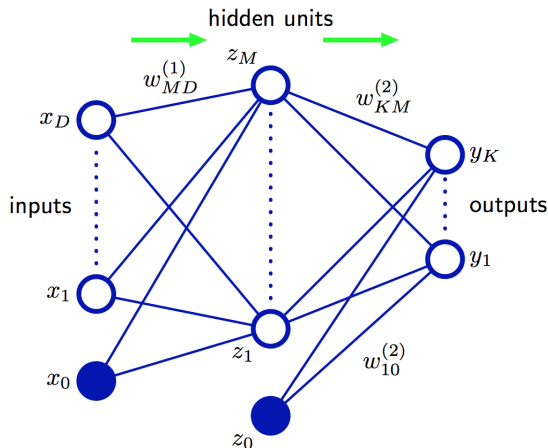
<sup>1</sup><http://yaroslavvb.blogspot.lu/2011/09/notmnist-dataset.html>

# Motivation

Goal: forecasting of time series

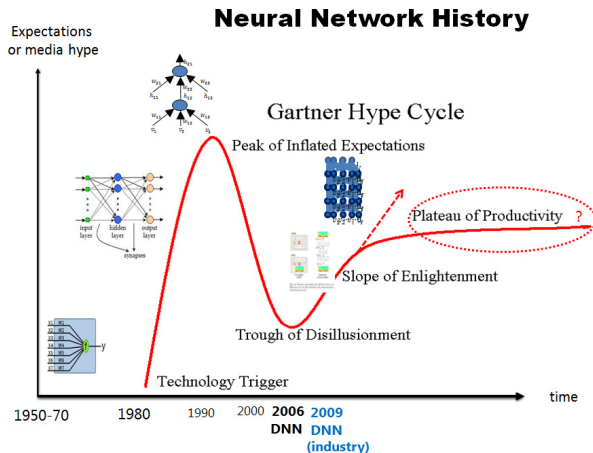


1. Neural networks
2. Deep Learning
3. TensorFlow
4. Distributed computing
5. Example: character recognition
6. Example: time series forecasting
7. Rise of the machines?
8. Conclusions and outreach



Neural network with two input and output units<sup>2</sup>.

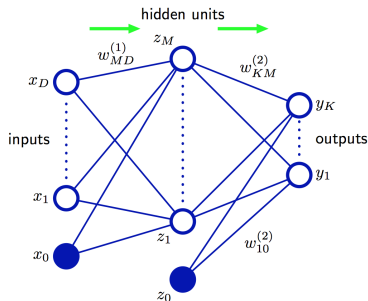
<sup>2</sup>Christopher M. Bishop, "Pattern Recognition and Machine Learning", Springer, 2007.



History of neural networks<sup>3</sup>.

<sup>3</sup>Li Deng and Dong Yu, "Deep Learning Methods and Applications", Foundations and Trends in Signal Processing, vol. 7 issues 3-4, pp. 197-387, 2014.



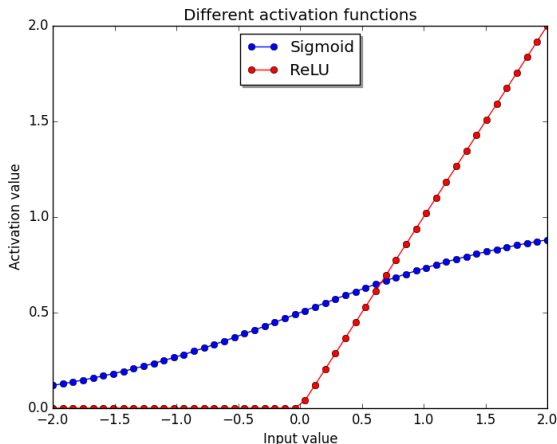


Neural network with two input and output units.

The activation of unit  $i$  of layer  $j + 1$  can be calculated:

$$z_i^{(j+1)} = \sum_{k=0}^{S_j} \Theta_{ik}^{(j)} x_k \quad (1)$$

$$a_i^{(j+1)} = g(z_i^{(j+1)}) \quad (2)$$



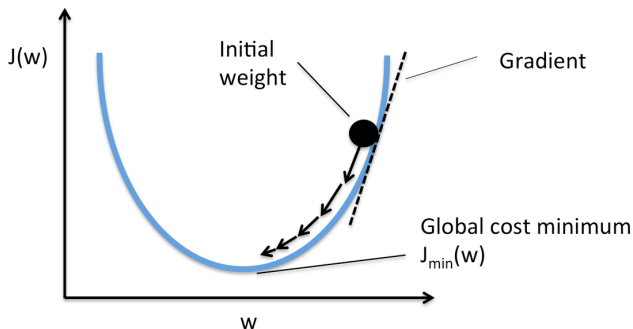
Sigmoid and rectified linear unit (ReLU) activation functions.

Cost function for  $m$  examples, hypothesis  $h_{\theta}$  and target values  $y^{(i)}$ :

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 \quad (3)$$

# Deep Learning: parameter optimization

How to optimize the weights?



Visualization for one parameter<sup>4</sup>.

<sup>4</sup><http://sebastianraschka.com/faq/docs/closed-form-vs-gd.html>

---

**Algorithm 1** Batch gradient descent: training size  $m$ , learning rate  $\alpha$

---

**repeat**

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \text{ (simultaneously for all } j)$$

**until** convergence

---

---

**Algorithm 2** Stochastic gradient descent: training size  $m$ , learning rate  $\alpha$ .

---

Randomly shuffle data set

**repeat**

**for**  $i = 1$  to  $m$  **do**

$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta, (x^{(i)}, y^{(i)}))$  (simultaneously for all  $j$ )

**end for**

**until** convergence

---

How to compute the partial derivatives?

---

## Algorithm 3 Backpropagation: training size $m$

---

$\Theta_{ij}^{(l)} \leftarrow \text{rand}(-\varepsilon, \varepsilon)$  (for all  $l, i, j$ )

$\Delta_{ij}^{(l)} \leftarrow 0$  (for all  $l, i, j$ )

**for**  $i = 1$  to  $m$  **do**

$\mathbf{a}^{(1)} \leftarrow \mathbf{x}^{(i)}$

    Perform forward propagation to compute  $\mathbf{a}^{(l)}$  for  $l = 2, 3, \dots, L$

    Using  $\mathbf{y}^{(i)}$ , compute  $\delta^{(L)} = \mathbf{a}^{(L)} - \mathbf{y}^{(i)}$  ▷ "error"

    Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ :  $\delta^{(l)} = (\Theta^{(l)})^T \delta^{(l+1)} \circ g'(z^{(l)})$

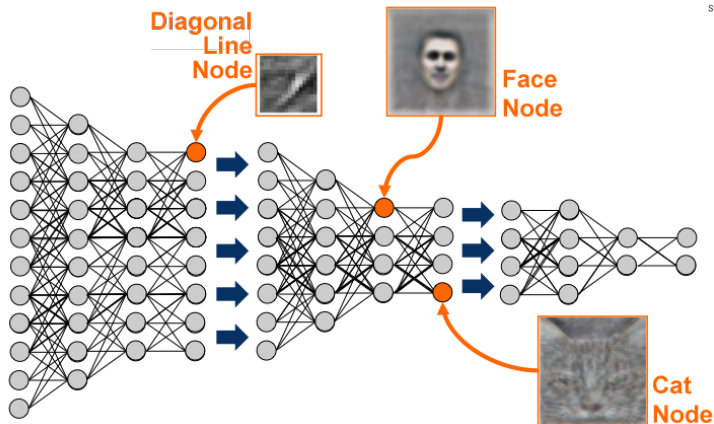
$\Delta^{(l)} \leftarrow \Delta^{(l)} + \delta^{(l+1)}(\mathbf{a}^{(l)})^T$  ▷ Matrix of errors for units of a layer

**end for**

$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) \leftarrow \frac{1}{m} \Delta_{ij}^{(l)}$

---





Deep neural network layers learning complex feature hierarchies<sup>5</sup>.

<sup>5</sup>The Analytics Store, "Deep Learning",  
<http://theanalyticsstore.com/deep-learning/>, retrieved: March 1, 2015.

# Deep Learning: DeepMind

- ▶ Founded in 2010 in London
- ▶ Created a neural network that learns how to play video games in a similar fashion to humans
- ▶ Acquired by Google in 2014, estimates range from USD 400 million to over GBP 500 million
- ▶ Now being used in Google's search engine
- ▶ AlphaGo played the game of Go at super-human performance



Google DeepMind<sup>6</sup>.

---

<sup>6</sup><http://deepmind.com/>, retrieved: March 2, 2016.

TensorFlow<sup>7</sup> is used by Google for most of its Deep Learning products:

- ▶ Offers neural networks (NN), convolutional neural networks (CNN), recurrent neural networks (RNN) and long-short term memories (LSTM)
- ▶ Computations are expressed as a data flow graph
- ▶ Can be used for research and production
- ▶ Python and C++ interfaces
- ▶ Code snippets available from Udacity class<sup>8</sup>

---

<sup>7</sup>J. Dean, R. Monga et al.: TensorFlow, "Large-Scale Machine Learning on Heterogeneous Distributed Systems", 2015.

<sup>8</sup><http://www.udacity.com/course/deep-learning--ud730>

# TensorFlow Playground

Let us experiment together with this playground for the next 20 minutes to get a better understanding of neural networks:

<http://playground.tensorflow.org>

The screenshot shows the TensorFlow Playground interface with the following settings and components:

- Iterations:** 000,126
- Learning rate:** 0.03
- Activation:** ReLU
- Regularization:** L2
- Regularization rate:** 0.01
- Problem type:** Classification

**DATA:** Which dataset do you want to use? (Icons for various datasets)

**FEATURES:** Which properties do you want to feed in? (Sliders for Ratio of training to test data: 50%, Noise: 25, Batch size: 11)

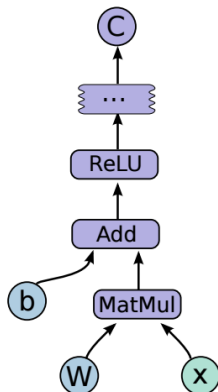
**2 HIDDEN LAYERS:** Each layer contains 4 neurons. The input features are  $X_1$ ,  $X_3$ ,  $X_1^2$ ,  $X_2^2$ , and  $X_1X_2$ . The output features are  $\sin(X_1)$  and  $\sin(X_2)$ .

**OUTPUT:** Test loss 0.136, Training loss 0.113. A 2D scatter plot shows data points colored by neuron and weight values. A legend indicates that colors show data, neuron and weight values, with a color scale from -1 to 1. There are checkboxes for "Show test data" and "Discretize output".

Annotations on the neural network diagram:

- "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."

- ▶ A Tensor is a typed multi-dimensional array
- ▶ Nodes in the graph are called ops
- ▶ An op takes zero or more Tensors, performs some computation, and produces zero or more Tensors
- ▶ Two phases:
  - ▶ Construction phase, that assembles a graph
  - ▶ Execution phase that uses a session to execute ops in the graph
- ▶ Auto-differentiation of the graph to compute partial derivatives used in stochastic gradient descent (SGD)



Sample computation graph<sup>9</sup>.

<sup>9</sup>J. Dean, R. Monga et al., "TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems", 2015.

Great documentation<sup>10</sup>.

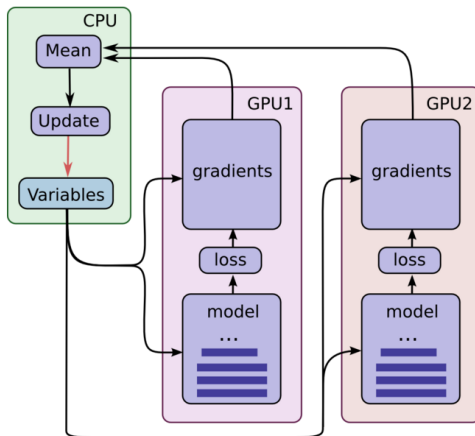
```
# Anaconda
$ sudo conda install \
  -c http://conda.anaconda.org/jjhelmus tensorflow
```

Support for Linux and Mac platforms, virtuelenv and Docker<sup>11</sup>. : time series

---

<sup>10</sup>[http://www.tensorflow.org/versions/0.6.0/get\\_started](http://www.tensorflow.org/versions/0.6.0/get_started)

<sup>11</sup>[http://www.tensorflow.org/versions/0.6.0/get\\_started/os\\_setup.html#pip\\_install](http://www.tensorflow.org/versions/0.6.0/get_started/os_setup.html#pip_install)



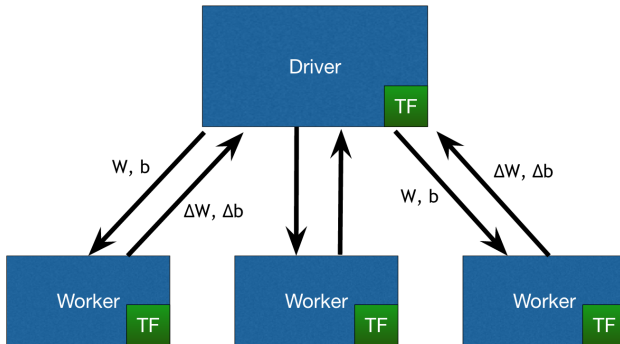
Parallel execution on multiple units<sup>12</sup>.

<sup>12</sup>J. Dean, R. Monga et al., "TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems", 2015.



# Distributed computing: Spark

Use of Spark for distributed computation of gradients:



Distributed computation of gradients<sup>13</sup>.

<sup>13</sup><http://arimo.com/machine-learning/deep-learning/2016/arimo-distributed-tensorflow-on-spark/>

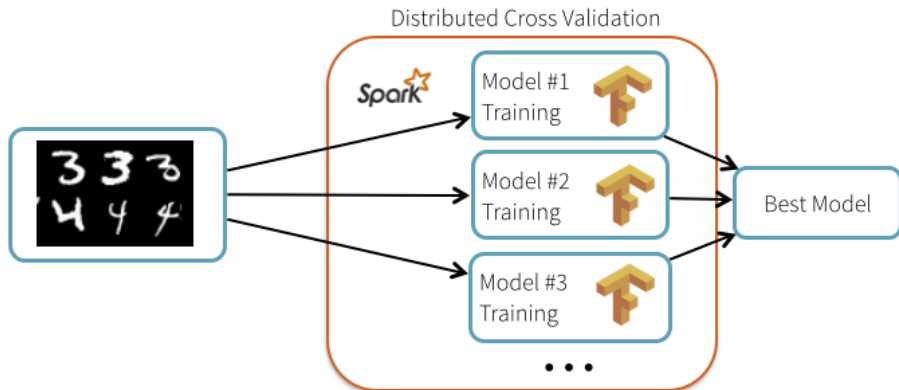
## Model selection

The process of optimizing various hyper parameters, including:

- ▶ Number of layers
- ▶ Size of a layer
- ▶ Learning rate
- ▶ Regularization
- ▶ ...

# Distributed computing: Spark

Use of Spark for distributed computation of model selection:



Distributed model selection on a single node<sup>14</sup>.

<sup>14</sup><http://databricks.com/blog/2016/01/25/deep-learning-with-apache-spark-and-tensorflow.html>

Google Cloud Machine Learning: <https://cloud.google.com/ml/>

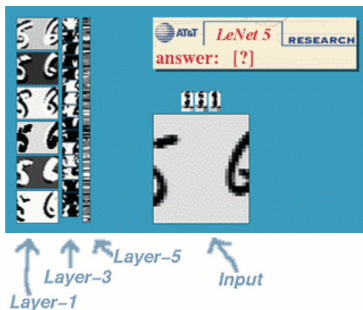


**CLOUD MACHINE LEARNING** BETA

Machine Learning on any data, of any size

# Example: character recognition

MNIST:



Hand-written digit recognition learned by a neural network<sup>15</sup>.

<sup>15</sup>Yann LeCun et al.: LeNet-5, convolutional neural networks.  
<http://yann.lecun.com/exdb/lenet/>. Retrieved: April 22, 2015.

# Example: character recognition

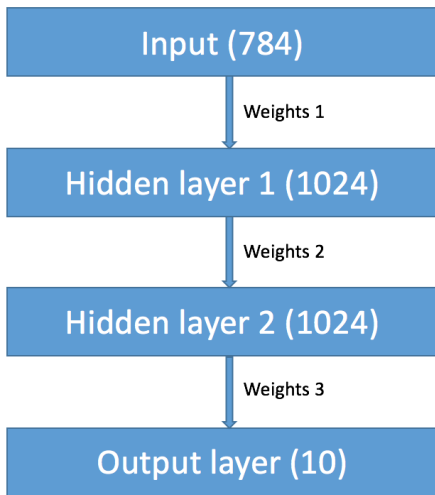
notMNIST: letters A-J.



notMNIST examples<sup>16</sup>.

<sup>16</sup><http://yaroslavvb.blogspot.lu/2011/09/notmnist-dataset.html>

# Example: character recognition



Architecture of network (biases omitted).

# Example: character recognition

- ▶ Source code: [http://github.com/pglauner/UCC\\_2016\\_Tutorial](http://github.com/pglauner/UCC_2016_Tutorial)
- ▶ Run `create_notmnist.py` once to get and convert the data
- ▶ Run `notmnist_classifier.py` for the experiments



# Example: character recognition

```
weights1 = tf.Variable(  
    tf.truncated_normal([image_size * image_size, 1024])  
biases1 = tf.Variable(tf.zeros([1024]))  
  
weights2 = tf.Variable(  
    tf.truncated_normal([1024, 1024]))  
biases2 = tf.Variable(tf.zeros([1024]))  
  
weights3 = tf.Variable(  
    tf.truncated_normal([1024, num_labels]))  
biases3 = tf.Variable(tf.zeros([num_labels]))
```

# Example: character recognition

```
[...]  
def model(data, train=False):  
    hidden1 = tf.nn.relu(  
        tf.matmul(data, weights1) + biases1)  
    if train:  
        hidden1 = tf.nn.dropout(hidden1, 0.7, seed=SEED)  
  
    hidden2 = tf.nn.relu(  
        tf.matmul(hidden1, weights2) + biases2)  
    if train:  
        hidden2 = tf.nn.dropout(hidden2, 0.7, seed=SEED)  
  
    return tf.matmul(hidden2, weights3) + biases3  
[...]
```

# Example: character recognition

```
logits = model(tf_train_dataset, True)

loss = tf.reduce_mean(
    tf.nn.softmax_cross_entropy_with_logits(
        logits, tf_train_labels))

# L2 regularization for the fully connected parameters
regularizers = (tf.nn.l2_loss(weights1)
                + tf.nn.l2_loss(biases1)
                + tf.nn.l2_loss(weights2)
                + tf.nn.l2_loss(biases2)
                + tf.nn.l2_loss(weights3)
                + tf.nn.l2_loss(biases3))

loss += 5e-4 * regularizers
```

# Example: character recognition

```
Training set (200000, 784) (200000, 10)  
Validation set (10000, 784) (10000, 10)  
Test set (10000, 784) (10000, 10)
```

Initialized

```
Minibatch loss at step 0: 13926.021484
```

```
Minibatch accuracy: 7.8%
```

```
Validation accuracy: 25.4%
```

```
Minibatch loss at step 500: 839.786133
```

```
Minibatch accuracy: 76.6%
```

```
Validation accuracy: 81.2%
```

```
[...]
```

```
Minibatch loss at step 2500: 515.079651
```

```
Minibatch accuracy: 78.9%
```

```
Validation accuracy: 80.4%
```

```
Minibatch loss at step 3000: 503.497894
```

```
Minibatch accuracy: 66.4%
```

```
Validation accuracy: 80.1%
```

```
Test accuracy: 87.2%
```

# Example: character recognition

Goal: become invariant to translation and rotation

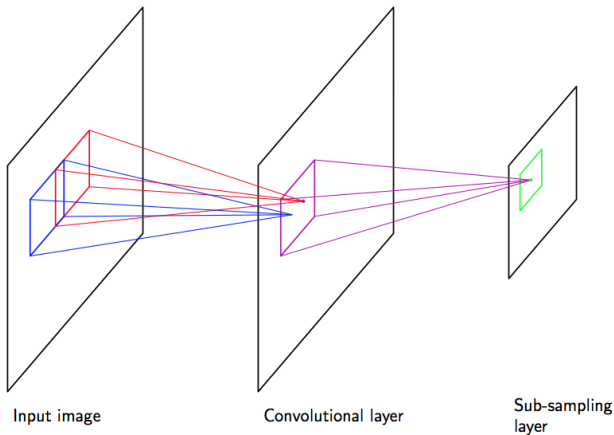


Illustration of a Convolutional Neural Network (CNN)<sup>17</sup>.

<sup>17</sup>C. M. Bishop, "Pattern Recognition and Machine Learning", Springer, 2007.

# Example: character recognition

- ▶ Source code: [http://github.com/pglauner/UCC\\_2016\\_Tutorial](http://github.com/pglauner/UCC_2016_Tutorial)
- ▶ Run `notminst_classifier_CNN.py` for the experiments

# Example: character recognition

Training set (200000, 28, 28, 1) (200000, 10)

Validation set (10000, 28, 28, 1) (10000, 10)

Test set (10000, 28, 28, 1) (10000, 10)

Initialized

Minibatch loss at step 0: 5.747538

Minibatch accuracy: 6.2%

Validation accuracy: 10.0%

Minibatch loss at step 500: 0.642069

Minibatch accuracy: 87.5%

Validation accuracy: 81.9%

[...]

Minibatch loss at step 2500: 0.721265

Minibatch accuracy: 75.0%

Validation accuracy: 86.1%

Minibatch loss at step 3000: 0.646058

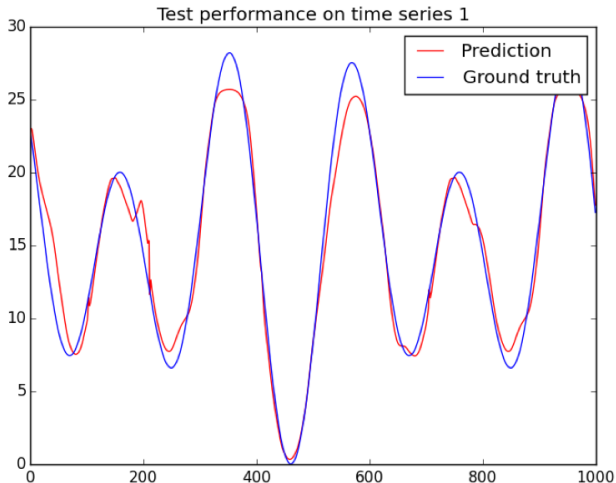
Minibatch accuracy: 87.5%

Validation accuracy: 86.5%

Test accuracy: 93.2%

# Example: time series forecasting

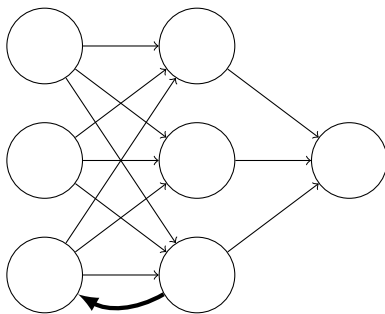
Goal: predict time series of electricity load





## Example: time series forecasting

- ▶ Feed-forward networks lack the ability to handle temporal data
- ▶ Recurrent neural networks (RNNs) have cycles in the graph structure, allowing them to keep temporal information



Simple RNN, current connection in **bold**.

- ▶ A long short-term memory (LSTM)<sup>18</sup> is a modular recurrent neural network composed of LSTM cell
- ▶ LSTM cells can be put together in a modular structure to build complex recurrent neural networks
- ▶ LSTMs have been reported to outperform regular RNNs and Hidden Markov Models in classification and time series prediction tasks<sup>19</sup>

---

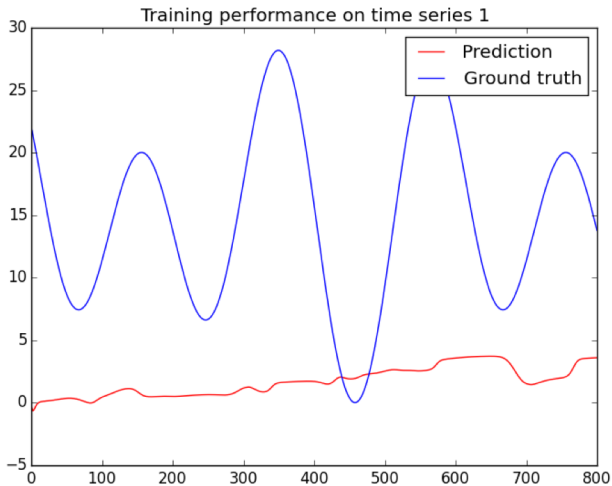
<sup>18</sup>S. Hochreiter and J. Schmidhuber, "Long short-term memory", Neural Computation, vol. 9, issue 8, pp. 1735-1780, 1997.

<sup>19</sup>N. Srivastava, E. Mansimov and R. Salakhutdinov, "Unsupervised Learning of Video Representations using LSTMs", University of Toronto, 2015.

- ▶ Source code: [http://github.com/pglauner/UCC\\_2016\\_Tutorial](http://github.com/pglauner/UCC_2016_Tutorial)
- ▶ Run LSTM.py for the experiments
- ▶ Simplified example, as time series is synthetic and harmonic
- ▶ More complex task will follow later

- ▶ Training on two time series at the same time
- ▶ Input values of each time series: value, derivative, second-order derivative
- ▶ Training data must be sufficiently long

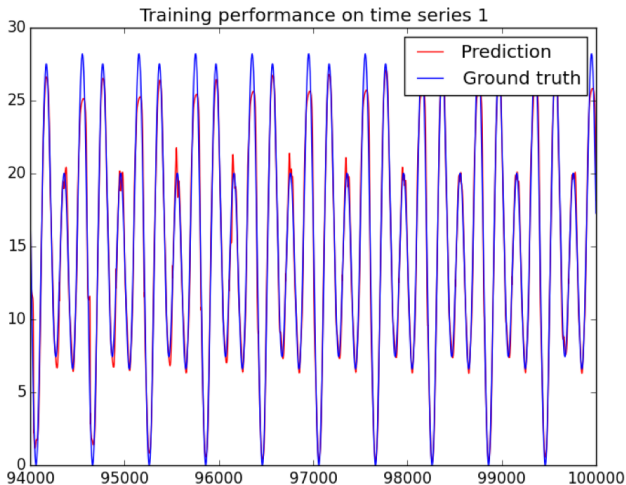
# Example: time series forecasting



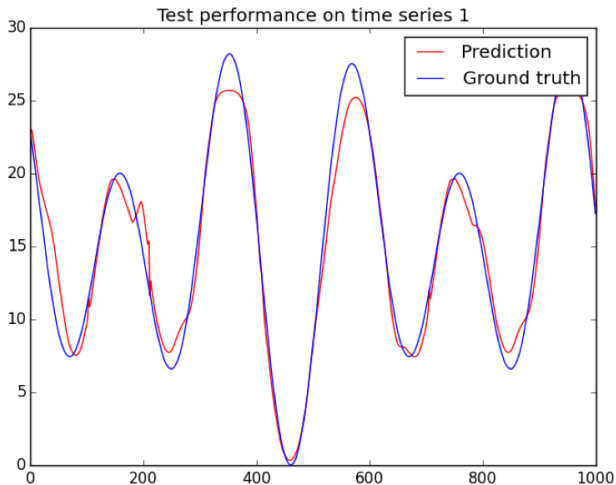
# Example: time series forecasting



# Example: time series forecasting



# Example: time series forecasting





# Example: time series forecasting

```
# Input layer for 6 inputs, batch size 1
input_layer = tf.placeholder(tf.float32,
                             [1, INPUT_DIM * 3])

# Initialization of LSTM layer
lstm_layer = rnn_cell.BasicLSTMCell(INPUT_DIM * 3)
# LSTM state, initialized to 0
lstm_state = tf.Variable(
    tf.zeros([1, lstm_layer.state_size]))
# Connect input layer to LSTM
lstm_output, lstm_state_output1 = lstm_layer(
    input_layer, lstm_state)
# Update of LSTM state
lstm_update = lstm_state.assign(lstm_state_output1)
```

## Example: time series forecasting

```
# Regression output layer
# Weights and biases
output_W = tf.Variable(
    tf.truncated_normal([INPUT_DIM * 3, INPUT_DIM]))
output_b = tf.Variable(tf.zeros([INPUT_DIM]))
output_layer = tf.matmul(lstm_output, output_W)
                + output_b

# Input for correct output (for training)
output_ground_truth = tf.placeholder(
    tf.float32, [1, INPUT_DIM])

# Sum of squared error terms
error = tf.pow(tf.sub(output_layer,
                      output_ground_truth), 2)

# Adam optimizer
optimizer = tf.train.AdamOptimizer(0.0006)
            .minimize(error)
```

- ▶ Add some noise for more realistic synthetic data
- ▶ Real-world load forecasting problem: <http://www.kaggle.com/c/global-energy-forecasting-competition-2012-load-forecasting>
- ▶ Models can be applied to other regression problems or time series classification (e.g. for detection of electricity theft)
- ▶ Usually more features need to be added
- ▶ Model selection in order to tweak hyper parameters (architecture, learning rate, etc.)



## Do we have to be worried?

- ▶ Specialized AIs have made significant progress and started to outperform humans
- ▶ Do we have to be worried about machines taking over?
- ▶ When will we achieve the singularity, the point in time when machines will become more intelligent than humans?
- ▶ Fears are spread by Stephen Hawking and other researchers

## From a researcher who actually works on AI

"There's also a lot of hype, that AI will create evil robots with super-intelligence. That's an unnecessary distraction. [...] Those of us on the frontline shipping code, we're excited by AI, but we don't see a realistic path for our software to become sentient. [...] If we colonize Mars, there could be too many people there, which would be a serious pressing issue. But there's no point working on it right now, and that's why I can't productively work on not turning AI evil." (Andrew Ng)<sup>a</sup>

---

<sup>a</sup>[http://www.theregister.co.uk/2015/03/19/andrew\\_ng\\_baidu\\_ai/](http://www.theregister.co.uk/2015/03/19/andrew_ng_baidu_ai/)

## Some thoughts

- ▶ The fear of an out-of-control AI is exaggerated
- ▶ Fears are mostly spread by people who do not work on AI, such as Stephen Hawking
- ▶ A lot of work needs to be done to work towards an **artificial general intelligence**
- ▶ Working towards simulating the brain may achieve the singularity in the late 21st century<sup>a</sup>
- ▶ In any case, many jobs will disappear in the next decades
- ▶ If computers only do a larger fraction of today's jobs, this will put pressure on salaries

---

<sup>a</sup>M. Shanahan, "The Technological Singularity", MIT Press, 2015.

- ▶ Deep neural networks can learn complex feature hierarchies
- ▶ TensorFlow is a easy-to-use Deep Learning framework
- ▶ Significant speedup of training on GPUs or Spark
- ▶ Interfaces for Python and C++
- ▶ Offers rich functionality and advanced features, such as LSTMs
- ▶ Udacity class and lots of documentation and examples available
- ▶ AI will not turn evil so soon