

Deep Learning Lab

Paulo Rauber

paulo@idsia.ch

Imanol Schlag

imanol@idsia.ch

December 21, 2018



- 1 Overview
- 2 Practical preliminaries
- 3 Introduction to TensorFlow
- 4 Fundamental models
 - Linear regression
 - Feedforward neural networks
 - Convolutional neural networks
 - Recurrent neural networks
 - Long short-term memory networks
- 5 Selected models
 - Supervised learning: Highway/Residual layer, Seq2seq, DNC
 - Unsupervised learning: PixelRNN, GAN, VAE
 - Reinforcement learning: RPG, A3C, TRPO, SNES
- 6 References

Overview

- Sessions: Fridays, 13:30 - 15:15, Sl. 006:
 - September: 21, 28
 - October: 5, 12, 19, 26
 - November: 2, 9, 16, 23, 30
 - December: 7, 14, 21
- Format: lectures and practical sessions
- Grading: four practical assignments

- 1 Overview
- 2 Practical preliminaries**
- 3 Introduction to TensorFlow
- 4 Fundamental models
 - Linear regression
 - Feedforward neural networks
 - Convolutional neural networks
 - Recurrent neural networks
 - Long short-term memory networks
- 5 Selected models
 - Supervised learning: Highway/Residual layer, Seq2seq, DNC
 - Unsupervised learning: PixelRNN, GAN, VAE
 - Reinforcement learning: RPG, A3C, TRPO, SNES
- 6 References

Python

- High-level multi-paradigm programming language
- Additional reading:
 - Dive into Python 3 [Pilgrim, 2011]
 - PEP 8 – Style Guide for Python Code [Rossum et al., 2001]
 - NumPy docstring guide [Num, 2017b]



Virtualenv

- virtualenv: a tool to create isolated python environments

```
virtualenv --system-site-packages -p python3 test_env #creates environment  
echo $PATH #original directories to search for executable files  
source test_env/bin/activate #activates environment  
echo $PATH #starts with /path/to/test_env/bin, containing python3 and pip3  
pip3 install numpy #installs numpy for the current environment  
python3 #this interpreter should be able to import numpy  
deactivate  
python3 #default interpreter (unaffected)
```

NumPy

- NumPy: scientific computing library for Python
 - powerful multidimensional arrays
 - efficient numerical computations
 - sophisticated broadcasting
- Additional reading:
 - NumPy Quickstart tutorial [Num, 2017c]
 - NumPy Broadcasting [Num, 2017a]



Slurm

- Slurm: workload manager for the ICS cluster [ICS, 2017]
 - Connect to hpc.ics.usi.ch using SSH
 - Use squeue to view jobs in the queue
 - Use scancel to send signals to jobs in the queue
 - Use sbatch to run a script that schedules a job. Example: ¹

```
#!/bin/bash -l
#
#SBATCH --job-name="abc"
#SBATCH --output=abc.%j.out
#SBATCH --error=abc.%j.err
#SBATCH --partition=tflow
#SBATCH --time=00:15:00
#SBATCH --exclusive

module load python/3.5.6-DL-Lab
srun python3 expression.py
```

¹Important: never run experiments directly on hpc.ics.usi.ch.

- 1 Overview
- 2 Practical preliminaries
- 3 Introduction to TensorFlow**
- 4 Fundamental models
 - Linear regression
 - Feedforward neural networks
 - Convolutional neural networks
 - Recurrent neural networks
 - Long short-term memory networks
- 5 Selected models
 - Supervised learning: Highway/Residual layer, Seq2seq, DNC
 - Unsupervised learning: PixelRNN, GAN, VAE
 - Reinforcement learning: RPG, A3C, TRPO, SNES
- 6 References

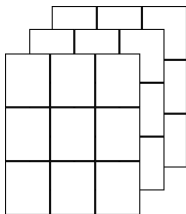
TensorFlow

- TensorFlow: library for numerical computations using data flow graphs
 - Scalable and multi-platform: from mobile devices to clusters
 - Enables transparent GPU usage
 - Widely used by researchers and practitioners
 - Python API



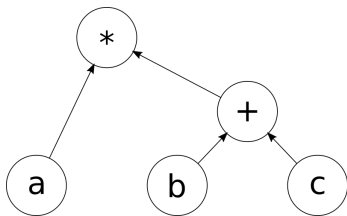
Tensor

- Tensor: for our purposes, synonymous with multidimensional array
- Examples of tensors:
 - Rank-0 tensor: real number
 - Rank-1 tensor: array of real numbers (real vector)
 - Rank-2 tensor: array of real vectors (real matrix)
 - Rank-3 tensor: array of real matrices



Computational graph

- Computational graph consists of nodes and edges:
 - Node: operation that receives zero or more tensors and produces zero or more tensors
 - Edge: connection between node outputs and node inputs



- Note: a constant is represented by a node that receives zero inputs and outputs the desired tensor.

Example: $\mathbf{a} \odot (\mathbf{b} + \mathbf{c})$

```
1 import tensorflow as tf
2
3
4 def main():
5     # Including constants in the default graph (nodes)
6     a = tf.constant([2, 3, 5], dtype=tf.float32)
7     b = tf.constant([1, 1, 3], dtype=tf.float32)
8     c = tf.constant([1, 2, 2], dtype=tf.float32)
9
10    # Including operations in the default graph (nodes)
11    b_plus_c = tf.add(b, c)
12    result = tf.multiply(a, b_plus_c)
13
14    # Using operator overloading, we could accomplish the same by writing
15    # result = a * (b + c)
16
17    # Creating a TensorFlow session
18    session = tf.Session()
19
20    # Using the session to obtain the output for node `result`
21    output = session.run(result) # np.array([4., 9., 25.])
22
23    print(output)
24
25    session.close()
26
27 if __name__ == "__main__":
28     main()
```

Session

- Session: responsible for managing resources to evaluate nodes
- Device management is possible (but often unnecessary):

```
1 # ...
2 with tf.device('/gpu:0'):
3     # Including constants in the default graph (nodes)
4     a = tf.constant([2, 3, 5], dtype=tf.float32)
5     b = tf.constant([1, 1, 3], dtype=tf.float32)
6     c = tf.constant([1, 2, 2], dtype=tf.float32)
7
8     # Including operations in the default graph (nodes)
9     b_plus_c = tf.add(b, c)
10    result = tf.multiply(a, b_plus_c)
11 # ...
```

- Additional reading:
 - TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems [Abadi et al., 2015]
 - TensorFlow: Using GPUs [Tf1, 2017c]

Variables

- Variable instance:
 - Represents state by a tensor
 - Usable as operand
- Variable instance adds several nodes to the graph:
 - Node that outputs initial state
 - Node that changes variable state as a side effect (assign node)
 - Node that outputs the current state (variable node)

Example: variables

```
1 def main():
2     a = tf.Variable([1.0, 1.0, 1.0], dtype=tf.float32) # Variable
3     b = tf.constant([1.0, 2.0, 3.0], dtype=tf.float32)
4
5     c = a * b
6
7     # Operation that assigns initial values to all variables (in our case, `a`)
8     initialize = tf.global_variables_initializer()
9
10    # Operation that assigns 2*a to `a`
11    assign_double = tf.assign(a, 2 * a)
12
13    session = tf.Session()
14
15    # Obtains `initialize` output. Side effect: initializing `a`
16    session.run(initialize)
17    print(session.run(c)) # np.array([1.0, 2.0, 3.0])
18
19    # Obtains `assign_double` output. Side effect: doubling `a`
20    session.run(assign_double)
21    print(session.run(c)) # np.array([2.0, 4.0, 6.0])
22    session.run(assign_double)
23    print(session.run(c)) # np.array([4.0, 8.0, 12.0])
24
25    session.close()
26
27    session = tf.Session()
28    session.run(initialize)
29    print(session.run(c)) # np.array([1.0, 2.0, 3.0])
30    session.close()
```


Placeholders

- Placeholder: unknown tensor during graph creation
 - Usable as operand
 - Must be provided through the feed mechanism

```
1 def main():
2     a = tf.constant([1.0, 2.0, 3.0], dtype=tf.float32)
3     b = tf.placeholder(dtype=tf.float32) # Placeholder, shape omitted
4
5     c = a * b
6
7     session = tf.Session()
8
9     print(session.run(c, feed_dict={b: 2.0})) # np.array([2.0, 4.0, 6.0])
10    print(session.run(c, feed_dict={b: [1.0, 2.0, 3.0]})) # np.array([1.0, 4.0, 9.0])
11
12    session.close()
```

Gradients

- `tf.gradient`: outputs partial derivative of a scalar with respect to each element of a tensor ².
- Example:

$$y = \sum_{i=1}^3 x_i^2 \implies \frac{\partial y}{\partial x_j} = 2x_j.$$

```
1 def main():
2     x = tf.Variable([1.0, 2.0, 3.0])
3     y = tf.reduce_sum(tf.square(x))
4
5     grad = tf.gradients(y, x)[0] # Gradient of y wrt `x`
6
7     initializer = tf.global_variables_initializer()
8
9     session = tf.Session()
10    session.run(initializer)
11    print(session.run(grad)) # np.array([2.0, 4.0, 6.0])
12    session.close()
```

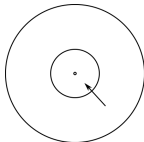
²The function is much more general. See the documentation for details.

Gradient descent

- Consider the task of minimizing $f : \mathbb{R}^D \rightarrow \mathbb{R}$.
- Gradient descent starts at an arbitrary estimate $\mathbf{x}_0 \in \mathbb{R}^D$ and iteratively updates this estimate using

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t),$$

where η_t is the learning rate at iteration t .

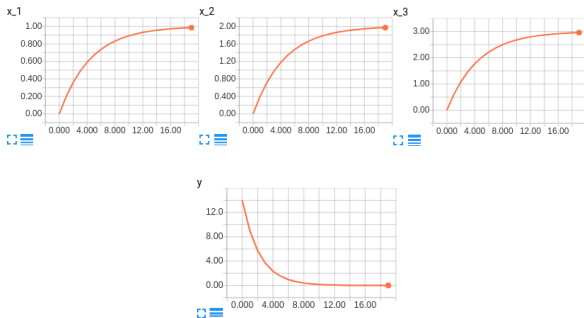


Example: minimizing $\sum_{i=1}^3 (x_i - i)^2$ wrt \mathbf{x}

```
1 def main():
2     n_iterations = 20
3
4     learning_rate = tf.constant(1e-1, dtype=tf.float32)
5
6     # Goal: finding x such that y is minimum
7     x = tf.Variable([0.0, 0.0, 0.0]) # Initial guess
8     y = tf.reduce_sum(tf.square(x - tf.constant([1.0, 2.0, 3.0])))
9
10    grad = tf.gradients(y, x)[0]
11
12    update = tf.assign(x, x - learning_rate * grad) # Gradient descent update
13
14    initializer = tf.global_variables_initializer()
15
16    session = tf.Session()
17    session.run(initializer)
18
19    for _ in range(n_iterations):
20        session.run(update)
21        print(session.run(x)) # State of `x` at this iteration
22
23    session.close()
```

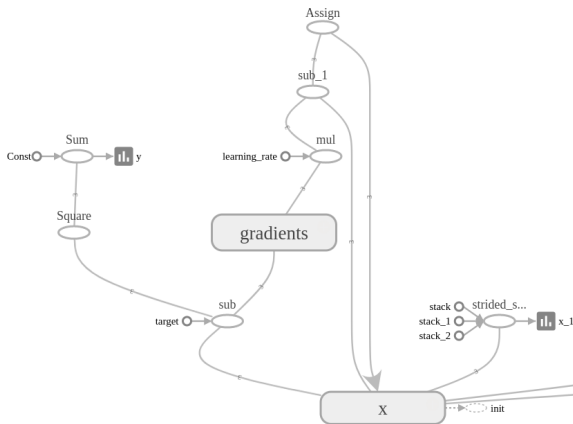
TensorBoard

- TensorBoard: visualizing summary data



TensorBoard

- TensorBoard: visualizing computational graph



TensorBoard

```
1 def main():
2     directory = '/tmp/gradient_descent' # Directory for data storage
3     os.makedirs(directory)
4
5     n_iterations = 20
6
7     # Naming constants/variables to facilitate inspection
8     learning_rate = tf.constant(1e-1, dtype=tf.float32, name='learning_rate')
9     x = tf.Variable([0.0, 0.0, 0.0], name='x')
10    target = tf.constant([1.0, 2.0, 3.0], name='target')
11    y = tf.reduce_sum(tf.square(x - target))
12
13    grad = tf.gradients(y, x)[0]
14
15    update = tf.assign(x, x - learning_rate * grad)
16
17    tf.summary.scalar('y', y) # Includes summary attached to `y`
18    tf.summary.scalar('x_1', x[0]) # Includes summary attached to `x[0]`
19    tf.summary.scalar('x_2', x[1]) # Includes summary attached to `x[1]`
20    tf.summary.scalar('x_3', x[2]) # Includes summary attached to `x[2]`
21
22    # Merges all summaries into single a operation
23    summaries = tf.summary.merge_all()
24
25    initializer = tf.global_variables_initializer()
26
27    # next slide ...
```

TensorBoard

```
1  # ... previous slide
2  session = tf.Session()
3
4  # Creating object that writes graph structure and summaries to disk
5  writer = tf.summary.FileWriter(directory, session.graph)
6
7  session.run(initializer)
8
9  for t in range(n_iterations):
10     # Updates `x` and obtains the summaries for iteration t
11     s, _ = session.run([summaries, update])
12
13     # Stores the summaries for iteration t
14     writer.add_summary(s, t)
15
16  print(session.run(x))
17
18  writer.close()
19  session.close()
20
21  # Run tensorboard --logdir="/tmp/gradient_descent" --port 6006
22  # Access http://localhost:6006 and see scalars/graphs
```


TensorFlow

- Additional reading
 - TensorFlow: Develop [Tf1, 2017a]
 - Get Started
 - Programmer's Guide
 - Tutorials
 - Performance
 - TensorFlow for deep learning research [Tf1, 2017b]
 - TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems [Abadi et al., 2015]

- 1 Overview
- 2 Practical preliminaries
- 3 Introduction to TensorFlow
- 4 Fundamental models**
 - Linear regression**
 - Feedforward neural networks
 - Convolutional neural networks
 - Recurrent neural networks
 - Long short-term memory networks
- 5 Selected models
 - Supervised learning: Highway/Residual layer, Seq2seq, DNC
 - Unsupervised learning: PixelRNN, GAN, VAE
 - Reinforcement learning: RPG, A3C, TRPO, SNES
- 6 References

Linear regression: model

- Consider an iid dataset $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$, where $\mathbf{x}_i \in \mathbb{R}^D$ and $y_i \in \mathbb{R}$.
- Regression: predicting target y given new observation \mathbf{x}
- Simple model:

$$y = \mathbf{w}\mathbf{x} = \sum_{j=1}^D w_j x_j$$

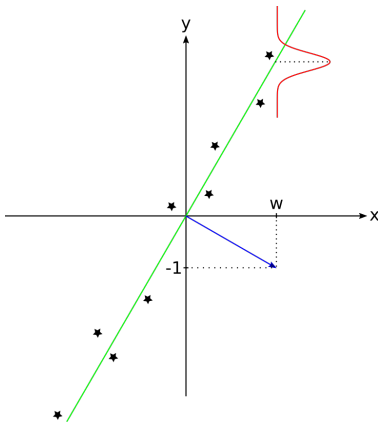
- Linear regression (without a bias term):

$$p(y \mid \mathbf{x}, \mathbf{w}) = \mathcal{N}(y \mid \mathbf{w}\mathbf{x}, \sigma^2),$$
$$\mathbb{E}[Y \mid \mathbf{x}, \mathbf{w}] = \mathbf{w}\mathbf{x}$$

Linear regression: geometry for $D = 1$

- The solutions to $w x - y = 0$ constitute a hyperplane

$$\{(x, y) \mid (w, -1) \cdot (x, y) = 0\}$$



Linear regression: likelihood

- Assuming constant σ^2 , the conditional likelihood is given by

$$p(\mathcal{D} \mid \mathbf{w}) = \prod_{i=1}^N \mathcal{N}(y_i \mid \mathbf{w}\mathbf{x}_i, \sigma^2)$$

- The log-likelihood is given by

$$\log p(\mathcal{D} \mid \mathbf{w}) = -\frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}\mathbf{x}_i)^2$$

- Maximizing the likelihood wrt \mathbf{w} corresponds to minimizing

$$J = \frac{1}{N} \sum_{i=1}^N (y_i - \mathbf{w}\mathbf{x}_i)^2$$

Linear regression: extensions

- If \mathbf{w} maximizes the likelihood, we may predict $y = \mathbf{w}\mathbf{x}$ given \mathbf{x}
 - Alternative: maximum a posteriori estimate (requires a prior)
 - Bayesian alternative: using a posterior predictive distribution
- Using a feature map $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^{D'}$:

$$p(y \mid \mathbf{x}, \mathbf{w}) = \mathcal{N}(y \mid \mathbf{w}\phi(\mathbf{x}), \sigma^2)$$

- Bias-including feature map: $\phi(\mathbf{x}) = (\mathbf{x}, 1)$
 - $\mathbf{w}\phi(\mathbf{x}) = \mathbf{w}_{1:D}\mathbf{x} + w_{D+1}$
- Polynomial feature map ($D = 1$): $\phi(x) = (1, x^1, \dots, x^{D'-1})$
 - $\mathbf{w}\phi(x) = \sum_{j=1}^{D'} w_j x^{j-1}$

Linear regression: additional reading

- Pattern Recognition and Machine Learning (Chapter 3) [Bishop, 2006]
- Machine Learning: a Probabilistic Perspective (Chapter 7) [Murphy, 2012]
- Notes on Machine Learning (Section 7) [Rauber, 2016]

Linear regression: example

```
1 def create_dataset(sample_size, n_dimensions, sigma, seed=None):
2     """Create linear regression dataset (without bias term)"""
3     random_state = np.random.RandomState(seed)
4
5     # True weight vector: np.array([1, 2, ..., n_dimensions])
6     w = np.arange(1, n_dimensions + 1)
7     # Randomly generating observations
8     X = random_state.uniform(-1, 1, (sample_size, n_dimensions))
9     # Computing noisy targets
10    y = np.dot(X, w) + random_state.normal(0.0, sigma, sample_size)
11
12    return X, y
13
14
15 def main():
16     sample_size_train = 100
17     sample_size_val = 100
18
19     n_dimensions = 10
20     sigma = 0.1
21
22     n_iterations = 20
23     learning_rate = 0.5
24
25     # Placeholder for the data matrix, where each observation is a row
26     X = tf.placeholder(tf.float32, shape=(None, n_dimensions))
27     # Placeholder for the targets
28     y = tf.placeholder(tf.float32, shape=(None,))
29
30     # next slide ...
```


Linear regression: example

```
1 # ... previous slide
2 # Variable for the model parameters
3 w = tf.Variable(tf.zeros((n_dimensions, 1)), trainable=True)
4
5 # Loss function
6 prediction = tf.reshape(tf.matmul(X, w), (-1,))
7 loss = tf.reduce_mean(tf.square(y - prediction))
8
9 optimizer = tf.train.GradientDescentOptimizer(learning_rate)
10 train = optimizer.minimize(loss) # Gradient descent update operation
11
12 initializer = tf.global_variables_initializer()
13
14 X_train, y_train = create_dataset(sample_size_train, n_dimensions, sigma)
15
16 session = tf.Session()
17 session.run(initializer)
18
19 for t in range(1, n_iterations + 1):
20     l, _ = session.run([loss, train], feed_dict={X: X_train, y: y_train})
21     print('Iteration {0}. Loss: {1}'.format(t, l))
22
23 X_val, y_val = create_dataset(sample_size_val, n_dimensions, sigma)
24 l = session.run(loss, feed_dict={X: X_val, y: y_val})
25 print('Validation loss: {0}'.format(l))
26
27 print(session.run(w).reshape(-1))
28
29 session.close()
```

- 1 Overview
- 2 Practical preliminaries
- 3 Introduction to TensorFlow
- 4 Fundamental models**
 - Linear regression
 - Feedforward neural networks**
 - Convolutional neural networks
 - Recurrent neural networks
 - Long short-term memory networks
- 5 Selected models
 - Supervised learning: Highway/Residual layer, Seq2seq, DNC
 - Unsupervised learning: PixelRNN, GAN, VAE
 - Reinforcement learning: RPG, A3C, TRPO, SNES
- 6 References

Classification task

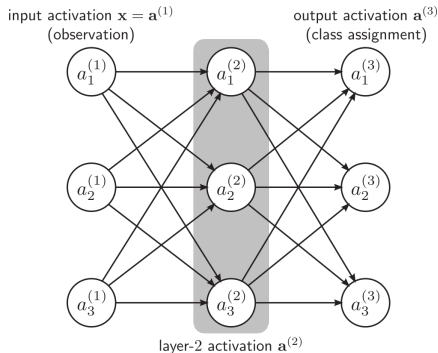
- Consider an iid dataset $\mathcal{D} = (\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)$, where $\mathbf{x}_i \in \mathbb{R}^D$, and $\mathbf{y}_i \in \{0, 1\}^C$
- Given a pair $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$, we assume $y_j = 1$ if and only if observation \mathbf{x} belongs to class j
- Each observation belongs to a single class
- Classification: predicting class assignment \mathbf{y} given new observation \mathbf{x}



$\mapsto (2, 3, 5, 7, 11) \mapsto \text{"dog"}$

Feedforward neural network (MLP)

- Let L be the number of layers in the network



- Let $N^{(l)}$ be the number of neurons in layer l
- Input neurons, hidden neurons, output neurons

Feedforward neural network (MLP)

- Weighted input to neuron j in layer $l > 1$:

$$z_j^{(l)} = b_j^{(l)} + \sum_{k=1}^{N^{(l-1)}} w_{j,k}^{(l)} a_k^{(l-1)},$$

- Activation of neuron j in layer $1 < l < L$:

$$a_j^{(l)} = \sigma(z_j^{(l)}),$$

where σ is a differentiable function, such as $\sigma(z) = \frac{1}{1+e^{-z}}$

Feedforward neural network (MLP)

- Alternatively, the output of each layer $1 < l < L$ can be written as

$$\mathbf{a}^{(l)} = \sigma(\mathbf{W}^{(l)}\mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}),$$

where the activation function is applied element-wise

- The (softmax) activation of output neuron j is given by

$$a_j^{(L)} = \frac{e^{z_j^{(L)}}}{\sum_{k=1}^C e^{z_k^{(L)}}.$$

- The output given $\mathbf{a}^{(1)} = \mathbf{x}$ is simply $\mathbf{a}^{(L)}$

Feedforward neural network (MLP)

- Let θ represent an assignment to weights and biases
- Maximizing the likelihood $p(\mathcal{D} \mid \theta)$ corresponds to minimizing

$$J = -\frac{1}{N} \log p(\mathcal{D} \mid \theta) = -\frac{1}{N} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \sum_{k=1}^C y_k \log a_k^{(L)}$$

with respect to θ

- The gradient $\nabla J(\theta)$ can be computed using backpropagation
- Minimization can be attempted by (stochastic) gradient descent or related techniques [Ruder, 2016]

Feedforward neural network (MLP)

- Additional reading:
 - Pattern Recognition and Machine Learning (Chapter 5) [Bishop, 2006]
 - Machine Learning: a Probabilistic Perspective (Section 16.5) [Murphy, 2012]
 - Neural networks and deep learning (Chapter 1) [Nielsen, 2015]
 - Notes on neural networks (Section 2) [Rauber, 2015]
 - Notes on machine learning (Section 17) [Rauber, 2016]

Example: MNIST classification

```
1 from tensorflow.examples.tutorials.mnist import input_data # tensorflow 1.8
2 def main():
3     tf.set_random_seed(seed=None)
4
5     # Downloads and loads MNIST dataset
6     mnist = input_data.read_data_sets('/tmp/mnist/', one_hot=True)
7     val_size = mnist.validation.num_examples
8
9     # Training procedure hyperparameters
10    learning_rate = 1e-2
11    batch_size = 64
12    n_epochs = 16
13    verbose_freq = 2000
14
15    # Model hyperparameters
16    n_neurons_1 = 784 # Number of input neurons (28 x 28 x 1)
17    n_neurons_2 = 100 # Number of neurons in the second layer (first hidden)
18    n_neurons_3 = 100 # Number of neurons in the third layer (second hidden)
19    n_neurons_4 = 10 # Number of output neurons (and classes)
20
21    X = tf.placeholder(tf.float32, [None, n_neurons_1])
22    Y = tf.placeholder(tf.float32, [None, n_neurons_4])
23
24    # Model parameters. Important: should not be initialized to zero
25    W2 = tf.Variable(tf.truncated_normal([n_neurons_1, n_neurons_2]))
26    W3 = tf.Variable(tf.truncated_normal([n_neurons_2, n_neurons_3]))
27    W4 = tf.Variable(tf.truncated_normal([n_neurons_3, n_neurons_4]))
28    b2 = tf.Variable(tf.zeros(n_neurons_2))
29    b3 = tf.Variable(tf.zeros(n_neurons_3))
30    b4 = tf.Variable(tf.zeros(n_neurons_4))
```

Example: MNIST classification

```
1  # Model definition
2  # The rectified linear activation function is given by  $a = \max(z, 0)$ 
3  A2 = tf.nn.relu(tf.matmul(X, W2) + b2)
4  A3 = tf.nn.relu(tf.matmul(A2, W3) + b3)
5  Z4 = tf.matmul(A3, W4) + b4
6
7  # Loss definition
8  # Important: this function expects weighted inputs, not activations
9  loss = tf.nn.softmax_cross_entropy_with_logits(labels=Y, logits=Z4)
10 loss = tf.reduce_mean(loss)
11
12 hits = tf.equal(tf.argmax(Z4, axis=1), tf.argmax(Y, axis=1))
13 accuracy = tf.reduce_mean(tf.cast(hits, tf.float32))
14
15 # Using Adam instead of gradient descent
16 optimizer = tf.train.AdamOptimizer(learning_rate)
17 train = optimizer.minimize(loss)
18
19 # Allows saving model to disc
20 saver = tf.train.Saver()
```

Example: MNIST classification

```
1 session = tf.Session()
2 session.run(tf.global_variables_initializer())
3
4 # Using mini-batches instead of entire dataset
5 n_batches = n_epochs * (mnist.train.num_examples // batch_size) # roughly
6 for t in range(n_batches):
7     X_batch, Y_batch = mnist.train.next_batch(batch_size)
8     session.run(train, {X: X_batch, Y: Y_batch})
9
10    # Computes validation loss every `verbose_freq` batches
11    if verbose_freq > 0 and t % verbose_freq == 0:
12        X_batch, Y_batch = mnist.validation.next_batch(val_size)
13        l = session.run(loss, {X: X_batch, Y: Y_batch})
14
15        print('Batch: {0}. Validation loss: {1}'.format(t, l))
16
17 saver.save(session, '/tmp/mnist.ckpt')
18 session.close()
19
20 # Loading model from file
21 session = tf.Session()
22 saver.restore(session, '/tmp/mnist.ckpt')
23
24 # In a proper experiment, test set results are computed only once, and
25 # absolutely never considered during the choice of hyperparameters
26 X_batch, Y_batch = mnist.test.next_batch(val_size)
27 acc = session.run(accuracy, {X: X_batch, Y: Y_batch})
28 print('Test accuracy: {0}'.format(acc))
29
30 session.close()
```

- 1 Overview
- 2 Practical preliminaries
- 3 Introduction to TensorFlow
- 4 Fundamental models**
 - Linear regression
 - Feedforward neural networks
 - Convolutional neural networks**
 - Recurrent neural networks
 - Long short-term memory networks
- 5 Selected models
 - Supervised learning: Highway/Residual layer, Seq2seq, DNC
 - Unsupervised learning: PixelRNN, GAN, VAE
 - Reinforcement learning: RPG, A3C, TRPO, SNES
- 6 References

Convolutional neural network: overview

- Convolutional neural network (CNN):
 - Parameterized function
 - Parameters may be adapted to minimize a cost function using gradient descent
 - Suitable for image classification: explores the spatial relationships between pixels
 - Three important types of layers: convolutional layers, max-pooling layers, and fully connected layers

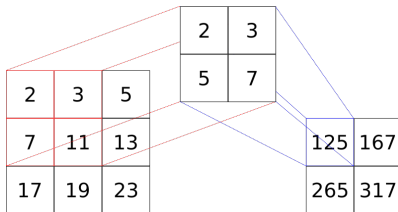
Convolutional neural network: notation

- Image: a function $\mathbf{f} : \mathbb{Z}^2 \rightarrow \mathbb{R}^c$
 - $\mathbf{a} \in \mathbb{Z}^2$ is a pixel
 - $\mathbf{f}(\mathbf{a})$ is the value of pixel \mathbf{a}
 - If $\mathbf{f}(\mathbf{a}) = (f_1(\mathbf{a}), \dots, f_c(\mathbf{a}))$, then f_i is channel i
 - Window $W \subset \mathbb{Z}^2$ is a finite set $W = [s_1, S_1] \times [s_2, S_2]$ that corresponds to a rectangle in the image domain
 - If the domain Z of an image \mathbf{f} is a window, it is possible to flatten \mathbf{f} into a vector $\mathbf{x} \in \mathbb{R}^{c|Z|}$
- Consider an iid dataset $\mathcal{D} = (\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)$, such that $\mathbf{x}_i \in \mathbb{R}^D$ and $\mathbf{y}_i \in \{0, 1\}^C$. Each vector \mathbf{x}_i corresponds to a distinct image $\mathbb{Z}^2 \mapsto \mathbb{R}^c$, and all images are defined on the same window Z , such that $D = c|Z|$ ³

³Note that we denote the number of colors by c and the number of classes by C .

Convolutional layer

- A neuron in a convolutional layer is not necessarily connected to the activations of all neurons in the previous layer, but only to the activations in a particular $w \times h$ window W
- A neuron in a convolutional layer is replicated through parameter sharing for all windows of size $w \times h$ in the domain Z whose centers are offset by pre-defined steps (strides)



Convolutional layer

- Receives an input image \mathbf{f} and outputs an image \mathbf{o}
- Each artificial neuron h in a convolutional layer l receives as input the values in a window $W = [s_1, S_1] \times [s_2, S_2] \subset Z$ of size $w \times h$, where Z is the domain of \mathbf{f} . The weighted input $z_h^{(l)}$ of that neuron is given by

$$z_h^{(l)} = b_h^{(l)} + \sum_{i=1}^c \sum_{j=s_1}^{S_1} \sum_{k=s_2}^{S_2} w_{h,i,j,k}^{(l)} a_{i,j,k}^{(l-1)},$$

where $a_{i,j,k}^{(l-1)} = f_i(j, k)$ is the value of pixel (j, k) in channel i of the input image \mathbf{f}

- Activation function is typically rectified linear: $a_h^{(l)} = \max(0, z_h^{(l)})$

Convolutional layer

- An output image $\mathbf{o} : \mathbb{Z}^2 \rightarrow \mathbb{R}^n$ is obtained by replicating n neurons over the whole domain of the input image
- The activations corresponding to a neuron replicated in this way correspond to the values in a single channel of the output image \mathbf{o} (appropriately arranged in \mathbb{Z}^2)
- The total number of free parameters in a convolutional layer is only $n(cwh + 1)$.

Convolutional layer

- If the parameters in a convolutional layer were not shared by replicated neurons, the number of parameters would be $mn(cwh + 1)$, where m is the number of windows of size $w \times h$ that fit into \mathbf{f} (for the given strides)
- A convolutional layer is fully specified by the size of the filters (window size), the number of filters (number of channels in the output image), horizontal and vertical strides (which are usually 1)

Max-pooling layer

- Goal: achieving similar results to using comparatively larger convolutional filters in the next layers with less parameters
- Receives an input image $\mathbf{f} : \mathbb{Z}^2 \rightarrow \mathbb{R}^c$ and outputs an image $\mathbf{o} : \mathbb{Z}^2 \rightarrow \mathbb{R}^c$
- Reduces the size of the window domain Z of \mathbf{f} by an operation that acts independently on each channel

$$o_i(j, k) = \max_{\mathbf{a} \in W_{j,k}} f_i(\mathbf{a}),$$

where $i \in \{1, \dots, c\}$, $(j, k) \in \mathbb{Z}^2$, Z is the window domain of \mathbf{f} , and $W_{j,k} \subseteq Z$ is the input window corresponding to output pixel (j, k) .

- A max-pooling layer is fully specified by the size of a pooling window and vertical/horizontal strides

Fully connected layer

- Receives a vector (or flattened image) and outputs a vector
- Analogous to a layer in a multilayer perceptron
- Typically only followed by other fully connected layers
- In a classification task, the output layer is typically fully connected with C neurons

Convolutional neural network

- Additional reading:
 - Pattern Recognition and Machine Learning (Chapter 5) [Bishop, 2006]
 - Machine Learning: a Probabilistic Perspective (Section 16.5) [Murphy, 2012]
 - Neural networks and deep learning (Chapter 6) [Nielsen, 2015]
 - Convolutional Neural Networks for Visual Recognition [Li and Karpathy, 2015]
 - Notes on neural networks (Section 5) [Rauber, 2015]
 - Notes on machine learning (Section 17) [Rauber, 2016]

Example: MNIST classification

```
1  # The placeholder `X` is the same as in the previous example
2  X_img = tf.reshape(X, [-1, 28, 28, 1])  # ? x 28 x 28 x 1
3
4  W_conv1 = tf.Variable(tf.truncated_normal([5, 5, 1, 32], stddev=0.1))  # 32 filters
5  b_conv1 = tf.Variable(tf.zeros(shape=(32,)))
6  A_conv1 = tf.nn.relu(tf.nn.conv2d(X_img, W_conv1, strides=[1, 1, 1, 1],
7                                padding='SAME') + b_conv1)  # ? x 28 x 28 x 32
8
9  A_pool1 = tf.nn.max_pool(A_conv1, ksize=[1, 2, 2, 1], strides=[1, 2, 2, 1],
10                        padding='SAME')  # ? x 14 x 14 x 32
11
12  W_conv2 = tf.Variable(tf.truncated_normal([5, 5, 32, 64], stddev=0.1))  # 64 filters
13  b_conv2 = tf.Variable(tf.zeros(shape=(64,)))
14  A_conv2 = tf.nn.relu(tf.nn.conv2d(A_pool1, W_conv2, strides=[1, 1, 1, 1],
15                                padding='SAME') + b_conv2)  # ? x 14 x 14 x 64
16
17  A_pool2 = tf.nn.max_pool(A_conv2, ksize=[1, 2, 2, 1], strides=[1, 2, 2, 1],
18                        padding='SAME')  # -1 x 7 x 7 x 64
19  A_pool2_flat = tf.reshape(A_pool2, [-1, 7 * 7 * 64])  # ? x 3136
20
21  W_fc1 = tf.Variable(tf.truncated_normal([7 * 7 * 64, 1024], stddev=0.1))
22  b_fc1 = tf.Variable(tf.zeros(shape=(1024,)))
23
24  A_fc1 = tf.nn.relu(tf.matmul(A_pool2_flat, W_fc1) + b_fc1)  # ? x 1024
25
26  W_fc2 = tf.Variable(tf.truncated_normal([1024, 10], stddev=0.1))
27  b_fc2 = tf.Variable(tf.zeros(shape=(10,)))
28
29  Z = tf.matmul(A_fc1, W_fc2) + b_fc2  # ? x 10
```

- 1 Overview
- 2 Practical preliminaries
- 3 Introduction to TensorFlow
- 4 Fundamental models**
 - Linear regression
 - Feedforward neural networks
 - Convolutional neural networks
 - Recurrent neural networks**
 - Long short-term memory networks
- 5 Selected models
 - Supervised learning: Highway/Residual layer, Seq2seq, DNC
 - Unsupervised learning: PixelRNN, GAN, VAE
 - Reinforcement learning: RPG, A3C, TRPO, SNES
- 6 References

Recurrent neural network: overview

- Recurrent neural network (RNN):
 - Parameterized function
 - Parameters may be adapted to minimize a cost function using gradient descent
 - Suitable for receiving a sequence of vectors and producing a sequence of vectors

Recurrent neural network: notation

- Let A^+ denote the set of non-empty sequences of elements from the set A , and let $|X|$ denote the length of a sequence $X \in A^+$
- Let $X[t]$ denote the t -th element of sequence X
- Consider the dataset:

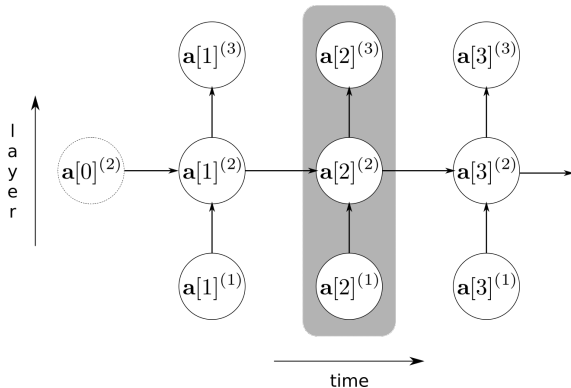
$$\mathcal{D} = \{(X_i, Y_i) \mid i \in \{1, \dots, N\}, X_i \in (\mathbb{R}^D)^+, Y_i \in (\mathbb{R}^C)^+\},$$

and let $|X| = |Y|$ for every $(X, Y) \in \mathcal{D}$

- In words, the dataset \mathcal{D} is composed of pairs (X, Y) of sequences of the same length. Each element of the two sequences is a real vector, but $X[t]$ and $Y[t]$ do not necessarily have the same dimension
- Sequence element classification: finding a function f that is able to generalize from \mathcal{D}

Recurrent neural network

- A recurrent neural network summarizes a sequence of vectors $X[1 : t - 1]$ into an activation vector
- This summary is combined with the input $X[t]$ to produce the output and the summary for the next timestep



Recurrent neural network

- We consider recurrent neural networks with a single recurrent layer and a softmax output layer
- In that case, the weighted input to neuron j in the recurrent layer at time step t is given by

$$z[t]_j^{(2)} = b_j^{(2)} + \sum_{k=1}^{N^{(1)}} w_{j,k}^{(2)} a[t]_k^{(1)} + \sum_{k=1}^{N^{(2)}} \omega_{j,k}^{(2)} a[t-1]_k^{(2)},$$

where $\mathbf{a}[0]^{(l)}$ is usually zero (or learnable)

- The corresponding activation is given by $a[t]_j^{(2)} = \sigma(z[t]_j^{(2)})$
- The output activation $\mathbf{a}[t]^{(3)}$ is computed from $\mathbf{a}[t]^{(2)}$ as usual

Recurrent neural network

- The output of the recurrent neural network on input $X = \mathbf{a}[1]^{(1)}, \dots, \mathbf{a}[T]^{(1)}$ is the sequence $\mathbf{a}[1]^{(3)}, \dots, \mathbf{a}[T]^{(3)}$
- Intuitively, the sequence X is presented to the network element by element
- The network behaves similarly to a single hidden layer feedforward neural network, except for the fact that the output activation $\mathbf{a}[t]^{(2)}$ of the hidden layer at time t possibly affects the weighted input $\mathbf{z}[t+1]^{(2)}$ of the hidden layer at time $t+1$
- An ideal recurrent neural network would be capable of representing a sequence $X[1:t]$ by its hidden layer activation $\mathbf{a}[t]^{(2)}$ to allow correct classification of $X[t+1]$
- Parameters are shared across time

Recurrent neural network

- Consider a sequence element classification cost function J given by

$$J = -\frac{1}{NT} \sum_{(X,Y) \in \mathcal{D}} \sum_{t=1}^T \sum_{j=1}^C Y[t]_j \log a[t]_j^{(3)}$$

- Let θ represent an assignment to weights and biases
- The gradient $\nabla J(\theta)$ can be computed using backpropagation through time
- Minimization can be attempted by (stochastic) gradient descent or related techniques [Ruder, 2016]

Recurrent neural network

- Additional reading:
 - Supervised sequence labelling with recurrent neural networks (Sec. 3.2) [Graves, 2012]
 - Notes on Neural networks (Sec. 6) [Rauber, 2015]
 - The Unreasonable Effectiveness of Recurrent Neural Networks [Karpathy, 2015]
 - Understanding LSTM Networks [Olah, 2015]
 - Recurrent Neural Networks in Tensorflow I [R2R, 2016]

Example: N-back using RNNs

```
1 import numpy as np
2 import tensorflow as tf
3
4
5 def nback(n, k, length, random_state):
6     """Creates n-back task given n, number of digits k, and sequence length.
7
8     Given a sequence of integers `xi`, the sequence `yi` has yi[t] = 1 if and
9     only if xi[t] == xi[t - n].
10    """
11     xi = random_state.randint(k, size=length) # Input sequence
12     yi = np.zeros(length, dtype=int) # Target sequence
13
14     for t in range(n, length):
15         yi[t] = (xi[t - n] == xi[t])
16
17     return xi, yi
```

Example: N-back using RNNs

```
1 def nback_dataset(n_sequences, mean_length, std_length, n, k, random_state):
2     """Creates dataset composed of n-back tasks."""
3     X, Y, lengths = [], [], []
4
5     for _ in range(n_sequences):
6         # Choosing length for current task
7         length = random_state.normal(loc=mean_length, scale=std_length)
8         length = int(max(n + 1, length))
9
10        # Creating task
11        xi, yi = nback(n, k, length, random_state)
12
13        # Storing task
14        X.append(xi)
15        Y.append(yi)
16        lengths.append(length)
17
18        # Creating padded arrays for the tasks
19        max_len = max(lengths)
20        Xarr = np.zeros((n_sequences, max_len), dtype=np.int64)
21        Yarr = np.zeros((n_sequences, max_len), dtype=np.int64)
22
23        for i in range(n_sequences):
24            Xarr[i, 0: lengths[i]] = X[i]
25            Yarr[i, 0: lengths[i]] = Y[i]
26
27    return Xarr, Yarr, lengths
```


Example: N-back using RNNs

```
1 def main():
2     seed = 0
3
4     # Task parameters
5     n = 3 # n-back
6     k = 4 # Input dimension
7     mean_length = 20 # Mean sequence length
8     std_length = 5 # Sequence length standard deviation
9     n_sequences = 512 # Number of training/validation sequences
10
11     # Creating datasets
12     random_state = np.random.RandomState(seed=seed)
13     X_train, Y_train, lengths_train = nback_dataset(n_sequences, mean_length,
14                                                    std_length, n, k,
15                                                    random_state)
16
17     X_val, Y_val, lengths_val = nback_dataset(n_sequences, mean_length,
18                                              std_length, n, k, random_state)
19
20     # Model parameters
21     hidden_units = 64 # Number of recurrent units
22
23     # Training procedure parameters
24     learning_rate = 1e-2
25     n_epochs = 256
26
27     # Model definition
28     X_int = tf.placeholder(shape=[None, None], dtype=tf.int64)
29     Y_int = tf.placeholder(shape=[None, None], dtype=tf.int64)
30     lengths = tf.placeholder(shape=[None], dtype=tf.int64)
```

Example: N-back using RNNs

```
1 batch_size = tf.shape(X_int)[0]
2 max_len = tf.shape(X_int)[1]
3
4 # One-hot encoding X_int
5 X = tf.one_hot(X_int, depth=k) # shape: (batch_size, max_len, k)
6 # One-hot encoding Y_int
7 Y = tf.one_hot(Y_int, depth=2) # shape: (batch_size, max_len, 2)
8
9 cell = tf.contrib.rnn.BasicRNNCell(num_units=hidden_units)
10 init_state = cell.zero_state(batch_size, dtype=tf.float32)
11
12 # rnn_outputs shape: (batch_size, max_len, hidden_units)
13 rnn_outputs, \
14     final_state = tf.nn.dynamic_rnn(cell, X, sequence_length=lengths,
15                                     initial_state=init_state)
16
17 # rnn_outputs_flat shape: ((batch_size * max_len), hidden_units)
18 rnn_outputs_flat = tf.reshape(rnn_outputs, [-1, hidden_units])
19
20 # Weights and biases for the output layer
21 Wout = tf.Variable(tf.truncated_normal(shape=(hidden_units, 2),
22                                         stddev=0.1))
23
24 bout = tf.Variable(tf.zeros(shape=[2]))
25
26 # Z shape: ((batch_size * max_len), 2)
27 Z = tf.matmul(rnn_outputs_flat, Wout) + bout
28
29 Y_flat = tf.reshape(Y, [-1, 2]) # shape: ((batch_size * max_len), 2)
```

Example: N-back using RNNs

```
1  # Creates a mask to disregard padding
2  mask = tf.sequence_mask(lengths, dtype=tf.float32)
3  mask = tf.reshape(mask, [-1]) # shape: (batch_size * max_len)
4
5  # Network prediction
6  pred = tf.argmax(Z, axis=1) * tf.cast(mask, dtype=tf.int64)
7  pred = tf.reshape(pred, [-1, max_len]) # shape: (batch_size, max_len)
8
9  hits = tf.reduce_sum(tf.cast(tf.equal(pred, Y_int), tf.float32))
10 hits = hits - tf.reduce_sum(1 - mask) # Disregards padding
11
12 # Accuracy: correct predictions divided by total predictions
13 accuracy = hits/tf.reduce_sum(mask)
14
15 # Loss definition (masking to disregard padding)
16 loss = tf.nn.softmax_cross_entropy_with_logits(labels=Y_flat, logits=Z)
17 loss = tf.reduce_sum(loss*mask)/tf.reduce_sum(mask)
18
19 optimizer = tf.train.AdamOptimizer(learning_rate)
20 train = optimizer.minimize(loss)
```

Example: N-back using RNNs

```
1 session = tf.Session()
2 session.run(tf.global_variables_initializer())
3
4 for e in range(1, n_epochs + 1):
5     feed = {X_int: X_train, Y_int: Y_train, lengths: lengths_train}
6     l, _ = session.run([loss, train], feed)
7     print('Epoch: {0}. Loss: {1}'.format(e, l))
8
9 feed = {X_int: X_val, Y_int: Y_val, lengths: lengths_val}
10 accuracy_ = session.run(accuracy, feed)
11 print('Validation accuracy: {0}'.format(accuracy_))
12
13 # Shows first task and corresponding prediction
14 xi = X_val[0, 0: lengths_val[0]]
15 yi = Y_val[0, 0: lengths_val[0]]
16 print('Sequence:')
17 print(xi)
18 print('Ground truth:')
19 print(yi)
20 print('Prediction:')
21 print(session.run(pred, {X_int: [xi], lengths: [len(xi)]})[0])
22
23 session.close()
```

- 1 Overview
- 2 Practical preliminaries
- 3 Introduction to TensorFlow
- 4 Fundamental models**
 - Linear regression
 - Feedforward neural networks
 - Convolutional neural networks
 - Recurrent neural networks
 - Long short-term memory networks**
- 5 Selected models
 - Supervised learning: Highway/Residual layer, Seq2seq, DNC
 - Unsupervised learning: PixelRNN, GAN, VAE
 - Reinforcement learning: RPG, A3C, TRPO, SNES
- 6 References

Long short-term memory network: overview

- Long short-term memory network (LSTM):
 - Parameterized function
 - Parameters may be adapted to minimize a cost function using gradient descent
 - Suitable for receiving a sequence of vectors and producing a sequence of vectors
 - Mitigates the vanishing gradients problem
 - Better than simple recurrent neural networks at learning dependencies between input and target vectors that manifest after many time steps

Long short-term memory network: overview

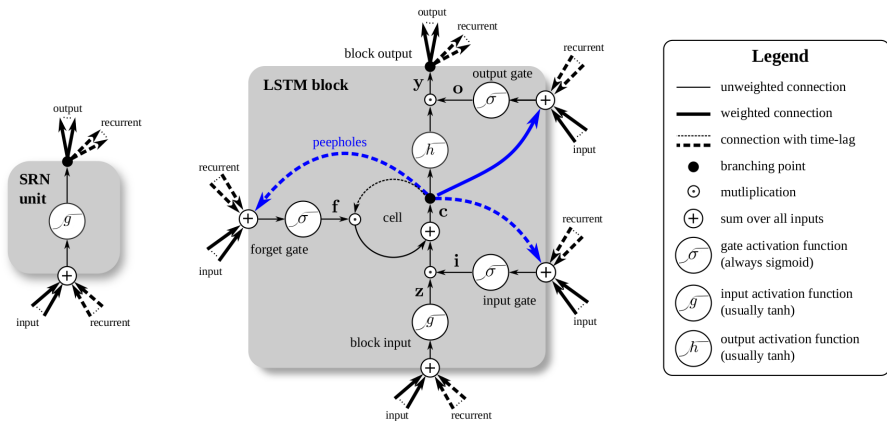


Image from [Greff et al., 2016]

Long short-term memory network

- We consider long short-term memory networks with a single hidden layer for the task of sequence element classification
- The input activation for the network at time t for $(X, Y) \in \mathcal{D}$ is defined as $X[t] = \mathbf{a}[t]^{(1)}$
- Similarly to a neuron in the hidden layer of a simple recurrent neural network, memory block j also receives the vectors $\mathbf{a}[t]^{(1)}$ and $\mathbf{a}[t-1]^{(2)}$ at time step t , and outputs a scalar $a[t]_j^{(2)}$
- However, the computations performed in a memory block are considerably more involved than those in a simple recurrent artificial neuron

Long short-term memory network

- A memory block is composed of four modules: cell, input gate I , forget gate F , and output gate O
- The weighted input $z[t]_j^{(2)}$ to the cell in memory block j is defined as

$$z[t]_j^{(2)} = b_j^{(2)} + \sum_{k=1}^{N^{(1)}} w_{j,k}^{(2)} a[t]_k^{(1)} + \sum_{k=1}^{N^{(2)}} \omega_{j,k}^{(2)} a[t-1]_k^{(2)},$$

where $\mathbf{a}[0]^{(2)}$ may be zero

- This is analogous to the weighted input for neuron j in the hidden layer of a simple recurrent network

Long short-term memory network

- The activation $s[t]_j^{(2)}$ of the cell in memory block j is defined as

$$s[t]_j^{(2)} = a[t]_{F,j}^{(2)} s[t-1]_j^{(2)} + a[t]_{I,j}^{(2)} g(z[t]_j^{(2)}),$$

where $s[0]^{(2)}$ may be zero, and g is a differentiable activation function

- The terms $a[t]_{F,j}^{(2)}$ and $a[t]_{I,j}^{(2)}$ correspond to the activations of the forget and input gates, respectively, and will be defined shortly
- Because each of these two scalars is usually between 0 and 1, they control how much the previous activation of the cell and the current weighted input to the cell affect its current activation

Long short-term memory network

- The weighted input $z[t]_{G,j}^{(2)}$ of a gate $G = I, F$ or O in memory block j is defined as

$$z[t]_{G,j}^{(2)} = b_{G,j}^{(2)} + \psi_{G,j}^{(2)} s[t-1]_j^{(2)} + \sum_{k=1}^{N^{(1)}} w_{G,j,k}^{(2)} a[t]_k^{(1)} + \sum_{k=1}^{N^{(2)}} \omega_{G,j,k}^{(2)} a[t-1]_k^{(2)},$$

where $\psi_{G,j}$ is the so-called peephole weight

- The activation $a[t]_{G,j}^{(2)}$ of a gate G in memory block j is defined as $a[t]_{G,j}^{(2)} = f(z[t]_{G,j}^{(2)})$, where f is typically the sigmoid function
- Each gate G in memory block j has its own parameters and behaves similarly to a simple recurrent neuron

Long short-term memory network

- The output activation $a[t]_j^{(2)}$ of memory block j is defined as

$$a[t]_j^{(2)} = a[t]_{O_j}^{(2)} h(s[t]_j^{(2)}),$$

where h is a differentiable activation function

- The activation of the output gate controls how much the current activation of the cell affects the output of the memory block
- A memory block can be interpreted as a parameterized circuit. By training the network, a memory block may learn when to store, output and erase its memory (cell activation), given the current input activation to the network and the previous activation of the memory blocks

Long short-term memory network

- The output activation $\mathbf{a}[t]^{(3)}$ is computed from $\mathbf{a}[t]^{(2)}$ as usual
- The output of the long short-term memory network on input $X = \mathbf{a}[1]^{(1)}, \dots, \mathbf{a}[T]^{(1)}$ is the sequence $\mathbf{a}[1]^{(3)}, \dots, \mathbf{a}[T]^{(3)}$
- An ideal LSTM would be capable of representing a sequence $X[1 : t]$ by the activation of its memory blocks $\mathbf{a}[t]^{(2)}$ and cells $\mathbf{s}[t]^{(2)}$ to allow correct classification of $X[t + 1]$

Long short-term memory network

- Additional reading:
 - Supervised sequence labelling with recurrent neural networks (Chap. 4)[Graves, 2012]
 - LSTM: A search space odyssey [Greff et al., 2016]
 - Notes on Neural networks (Sec. 7) [Rauber, 2015]
 - The Unreasonable Effectiveness of Recurrent Neural Networks [Karpathy, 2015]
 - Understanding LSTM Networks [Olah, 2015]

Example: N-back using LSTMs

```
1 # ...
2 # One-hot encoding X_int
3 X = tf.one_hot(X_int, depth=k) # shape: (batch_size, max_len, k)
4 # One-hot encoding Y_int
5 Y = tf.one_hot(Y_int, depth=2) # shape: (batch_size, max_len, 2)
6
7 # There is a single change from the previous n-back example:
8 # cell = tf.contrib.rnn.BasicRNNCell(num_units=hidden_units)
9 cell = tf.contrib.rnn.LSTMCell(num_units=hidden_units)
10
11 init_state = cell.zero_state(batch_size, dtype=tf.float32)
12
13 # rnn_outputs shape: (batch_size, max_len, hidden_units)
14 rnn_outputs, \
15     final_state = tf.nn.dynamic_rnn(cell, X, sequence_length=lengths,
16                                     initial_state=init_state)
17
18 # rnn_outputs_flat shape: ((batch_size * max_len), hidden_units)
19 rnn_outputs_flat = tf.reshape(rnn_outputs, [-1, hidden_units])
20 # ...
```

- 1 Overview
- 2 Practical preliminaries
- 3 Introduction to TensorFlow
- 4 Fundamental models
 - Linear regression
 - Feedforward neural networks
 - Convolutional neural networks
 - Recurrent neural networks
 - Long short-term memory networks
- 5 Selected models
 - Supervised learning: Highway/Residual layer, Seq2seq, DNC**
 - Unsupervised learning: PixelRNN, GAN, VAE
 - Reinforcement learning: RPG, A3C, TRPO, SNES
- 6 References

Highway/residual layer

- Idea: information should be able to flow across layers unaltered
- Traditional layer:

$$\mathbf{a}^{(l)} = \mathbf{f}(\mathbf{W}^{(l)}\mathbf{a}^{(l-1)} + \mathbf{b}^{(l)})$$

- Residual layer [He et al., 2016]:

$$\mathbf{a}^{(l)} = \mathbf{a}^{(l-1)} + \mathbf{f}(\mathbf{W}^{(l)}\mathbf{a}^{(l-1)} + \mathbf{b}^{(l)})$$

- Highway layer (with coupled gates) [Srivastava et al., 2015]:

$$\mathbf{a}^{(l)} = \mathbf{a}^{(l-1)} \odot \mathbf{g}(\mathbf{a}^{(l-1)}) + \mathbf{f}(\mathbf{W}^{(l)}\mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}) \odot (\mathbf{1} - \mathbf{g}(\mathbf{a}^{(l-1)})),$$

where

$$\mathbf{g}(\mathbf{a}^{(l-1)}) = \sigma(\mathbf{W}^{(l,g)}\mathbf{a}^{(l-1)} + \mathbf{b}^{(l,g)})$$

Sequence to sequence model

- Idea: using an encoding phase followed by a decoding phase to map between sequences of arbitrary lengths
[Cho et al., 2014, Sutskever et al., 2014]

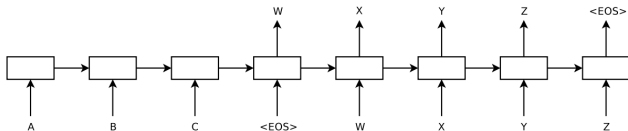


Image from [Sutskever et al., 2014]

- The recurrent networks that perform encoding and decoding are not necessarily the same

Differentiable neural computer

- Idea: a neural network can learn to read and write from a memory matrix using gating mechanisms [Graves et al., 2016]

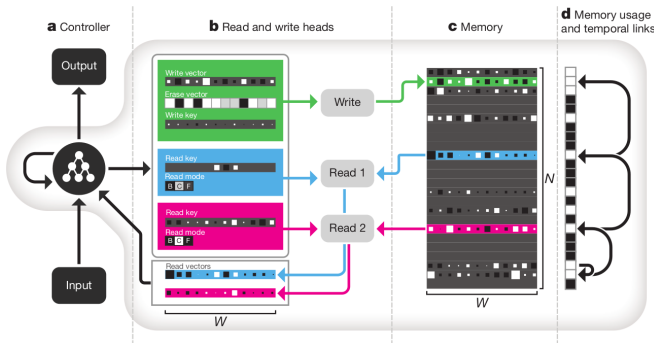


Image from [Graves et al., 2016]

- 1 Overview
- 2 Practical preliminaries
- 3 Introduction to TensorFlow
- 4 Fundamental models
 - Linear regression
 - Feedforward neural networks
 - Convolutional neural networks
 - Recurrent neural networks
 - Long short-term memory networks
- 5 Selected models
 - Supervised learning: Highway/Residual layer, Seq2seq, DNC
 - Unsupervised learning: PixelRNN, GAN, VAE
 - Reinforcement learning: RPG, A3C, TRPO, SNES
- 6 References

PixelRNN

- Idea: using a recurrent neural network trained to predict each pixel given the previous pixels as a probabilistic model
[van den Oord et al., 2016]

$$p(\mathbf{x} \mid \boldsymbol{\theta}) = \prod_{j=1}^d p(x_j \mid x_1, \dots, x_{j-1}, \boldsymbol{\theta})$$

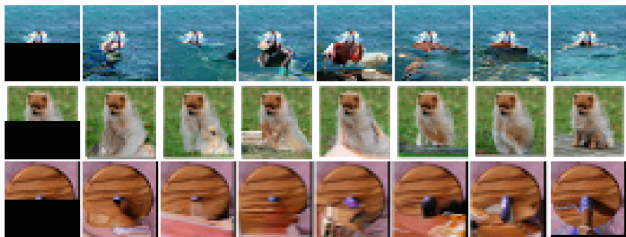


Image from [van den Oord et al., 2016]

Generative adversarial network

- Idea: training a (discriminator) network to discriminate between real and synthetic observations and training another (generator) network to generate synthetic observations from noise that fool the discriminator [Goodfellow et al., 2014]



Image from [Goodfellow et al., 2014]

Variational autoencoder

- Idea: training a model with (easy to sample) hidden variables by maximizing a particular lower bound on the log-likelihood
[Kingma and Welling, 2014, Rezende et al., 2014]

$$\int_{\text{Val}(\mathbf{Z})} p(\mathbf{x} | \mathbf{z}, \theta) p(\mathbf{z} | \theta) d\mathbf{z} = \int_{\text{Val}(\mathbf{Z})} \mathcal{N}(\mathbf{x} | \mathbf{f}(\mathbf{z}, \theta), \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}) d\mathbf{z}$$



Image from [Doersch, 2016]

- 1 Overview
- 2 Practical preliminaries
- 3 Introduction to TensorFlow
- 4 Fundamental models
 - Linear regression
 - Feedforward neural networks
 - Convolutional neural networks
 - Recurrent neural networks
 - Long short-term memory networks
- 5 Selected models
 - Supervised learning: Highway/Residual layer, Seq2seq, DNC
 - Unsupervised learning: PixelRNN, GAN, VAE
 - Reinforcement learning: RPG, A3C, TRPO, SNES
- 6 References

Recurrent policy gradient

- Idea: a recurrent neural network represents a policy by a probability distribution over actions given the history of observations and actions [Wierstra et al., 2009]
- The goal is to maximize the expected return J given by

$$J(\theta) = \mathbb{E} \left[\sum_{t=1}^T R_t \mid \theta \right] = \sum_{\tau} p(\tau \mid \theta) \sum_{t=1}^T r_t,$$

where θ are the policy parameters and τ denotes a trajectory

- It can be shown that $\nabla J(\theta)$ is given by

$$\nabla J(\theta) = \mathbb{E} \left[\sum_{t=1}^{T-1} \nabla_{\theta} \log p(A_t \mid X_{1:t}, A_{1:t-1}, \theta) \sum_{t'=t+1}^T R_{t'} \mid \theta \right]$$

- A Monte Carlo estimate may be used for gradient ascent

Asynchronous advantage actor-critic

- Idea: asynchronously updating policy parameters shared by several threads using policy gradients with a value-function baseline [Mnih et al., 2016]



Image from [DM1, 2016]

Trust region policy optimization

- Idea: approximating a minorization-maximization procedure that would lead to updates that never deteriorate the policy [Schulman et al., 2015]

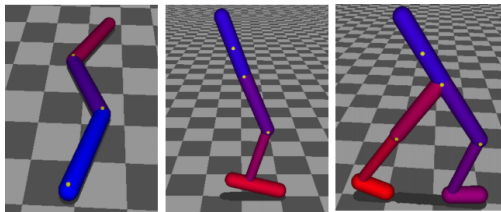


Image from [Schulman et al., 2015]

Separable natural evolution strategies

- We consider a simplified version of separable natural evolution strategies [Wierstra et al., 2014] that was applied to current benchmarks [Salimans et al., 2017]
- Let $J(\theta)$ denote the expected return of following a policy parameterized by θ
- Let $p(\theta \mid \psi) = \mathcal{N}(\theta \mid \psi, \sigma^2 \mathbf{I})$ and consider the task of maximizing the expected return η given by

$$\eta(\psi) = \mathbb{E}[J(\Theta) \mid \psi] = \int_{\text{Val}(\Theta)} p(\theta \mid \psi) J(\theta) d\theta$$

- It can be shown that $\nabla \eta(\psi)$ is given by

$$\nabla \eta(\psi) = \sigma^{-1} \mathbb{E}[J(\psi + \sigma \mathcal{E}) \mathcal{E}], \text{ where } \mathcal{E} \sim \mathcal{N}(\cdot \mid \mathbf{0}, \mathbf{I})$$

- A Monte Carlo estimate may be used for gradient ascent

Deep Learning Lab

Paulo Rauber

paulo@idsia.ch

Imanol Schlag

imanol@idsia.ch

December 21, 2018



References I



(2016).

Asynchronous methods for deep reinforcement learning: Labyrinth.

<https://www.youtube.com/watch?v=nMR5mjCFZCw>.



(2016).

Recurrent neural networks in tensorflow I.

[https://r2rt.com/
recurrent-neural-networks-in-tensorflow-i.html](https://r2rt.com/recurrent-neural-networks-in-tensorflow-i.html).



(2017).

Institute of computational science HPC.

<https://intranet.ics.usi.ch/HPC>.



(2017a).

Numpy broadcasting.

[https://docs.scipy.org/doc/numpy-1.13.0/user/basics.
broadcasting.html](https://docs.scipy.org/doc/numpy-1.13.0/user/basics.broadcasting.html).

References II



(2017b).

Numpy docstring guide.

<https://numpydoc.readthedocs.io/en/latest/format.html#docstring-standard>.



(2017c).

Numpy quickstart tutorial.

<https://docs.scipy.org/doc/numpy/user/quickstart.html>.



(2017a).

TensorFlow: Develop.

https://www.tensorflow.org/get_started/.



(2017b).

TensorFlow for deep learning research.

<http://web.stanford.edu/class/cs20si/syllabus.html>.

References III



(2017c).

TensorFlow: Using gpus.

https://www.tensorflow.org/tutorials/using_gpu.



Abadi, M. et al. (2015).

TensorFlow: Large-scale machine learning on heterogeneous systems.

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



Bishop, C. M. (2006).

Pattern Recognition and Machine Learning.

Springer.

References IV



Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014).

Learning phrase representations using rnn encoder–decoder for statistical machine translation.

In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP).



Doersch, C. (2016).

Tutorial on variational autoencoders.

arXiv preprint arXiv:1606.05908.



Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014).

Generative adversarial nets.

In Advances in neural information processing systems.

References V



Graves, A. (2012).

Supervised sequence labelling with recurrent neural networks.
Springer.



Graves, A., Wayne, G., Reynolds, M., Harley, T., Danihelka, I.,
Grabska-Barwińska, A., Colmenarejo, S. G., Grefenstette, E.,
Ramalho, T., Agapiou, J., et al. (2016).

Hybrid computing using a neural network with dynamic external
memory.
Nature.







Greff, K., Srivastava, R. K., Koutník, J., Steunebrink, B. R., and
Schmidhuber, J. (2016).




LSTM: A search space odyssey.

IEEE transactions on neural networks and learning systems.

References VI

-  He, K., Zhang, X., Ren, S., and Sun, J. (2016).
Deep residual learning for image recognition.
In Proceedings of the IEEE conference on computer vision and pattern recognition.
-  Karpathy, A. (2015).
The unreasonable effectiveness of recurrent neural networks.
<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>.
-  Kingma, D. P. and Welling, M. (2014).
Auto-encoding variational bayes.
In International Conference on Learning Representations.
-  Li, F.-F. and Karpathy, A. (2015).
Convolutional neural networks for visual recognition.
<http://cs231n.github.io/convolutional-networks>.

References VII

-  Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. (2016).
Asynchronous methods for deep reinforcement learning.
In International Conference on Machine Learning.
-  Murphy, K. P. (2012).
Machine learning: a probabilistic perspective.
MIT Press.
-  Nielsen, M. (2015).
Neural networks and deep learning.
Determination Press.
<http://neuralnetworksanddeeplearning.com>.

References VIII



Olah, C. (2015).
Understanding LSTM networks.

[http:
//colah.github.io/posts/2015-08-Understanding-LSTMs/](http://colah.github.io/posts/2015-08-Understanding-LSTMs/).



Pilgrim, M. (2011).
Dive into Python 3.

<http://www.diveintopython3.net/>.



Rauber, P. E. (2015).
Notes on neural networks.





http://paulorauber.com/notes/neural_networks.pdf.







Rauber, P. E. (2016).
Notes on machine learning.

http://paulorauber.com/notes/machine_learning.pdf.



References IX

-  Rezende, D. J., Mohamed, S., and Wierstra, D. (2014).
Stochastic backpropagation and approximate inference in deep generative models.
In International Conference on Machine Learning.
-  Rossum, G. v., Warsaw, B., and Coghlan, N. (2001).
PEP 8 – style guide for Python code.
<https://www.python.org/dev/peps/pep-0008/>.
-  Ruder, S. (2016).
An overview of gradient descent optimization algorithms.
<http://ruder.io/optimizing-gradient-descent/>.
-  Salimans, T., Ho, J., Chen, X., and Sutskever, I. (2017).
Evolution strategies as a scalable alternative to reinforcement learning.
[arXiv preprint arXiv:1703.03864](https://arxiv.org/abs/1703.03864).

References X

-  Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. (2015).
Trust region policy optimization.
In International Conference on Machine Learning.
-  Srivastava, R. K., Greff, K., and Schmidhuber, J. (2015).
Training very deep networks.
In Advances in neural information processing systems.
-  Sutskever, I., Vinyals, O., and Le, Q. V. (2014).
Sequence to sequence learning with neural networks.
In Advances in neural information processing systems.
-  van den Oord, A., Kalchbrenner, N., Espeholt, L., Vinyals, O., Graves, A., et al. (2016).
Conditional image generation with pixelCNN decoders.
In Advances in Neural Information Processing Systems.

References XI

-  Wierstra, D., Förster, A., Peters, J., and Schmidhuber, J. (2009).
Recurrent policy gradients.
Logic Journal of IGPL, 18(5).
-  Wierstra, D., Schaul, T., Glasmachers, T., Sun, Y., Peters, J., and Schmidhuber, J. (2014).
Natural evolution strategies.
Journal of Machine Learning Research.