# Neural Networks

## Semester 1

**Practical works for GRIAT
RCSE master program**

**Teacher: Makhmutova Alisa Zufarovna
AZMakhmutova@kai.ru
@DarkAliceSophie**

# Requirements
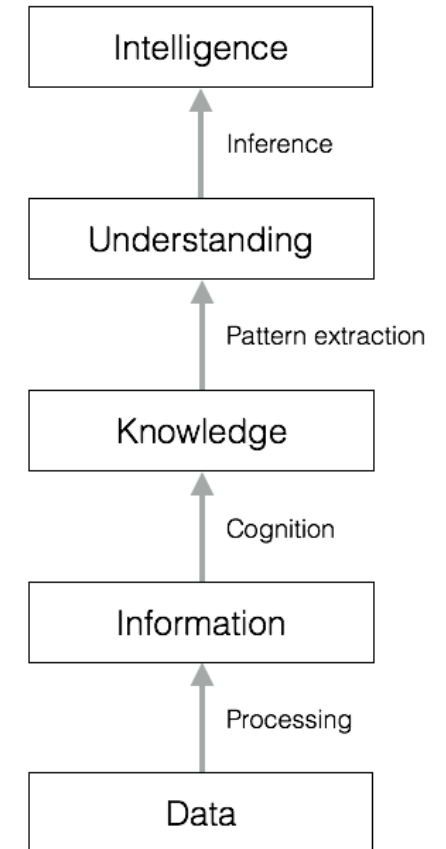
- Basic Python (2, 3) knowledge: syntax, operations, packages etc. (SoloLearn application, courses on edX, tutorials)

- Python IDE (https://www.jetbrains.com/pycharm/)

- Basic matrix calculation

# Motivation for Artificial Intelligence (AI)

- Big amount of data for processing

- Multiple simultaneous sources of data

- Unorganized, mixed data

- Evolving of data (knowledge update)

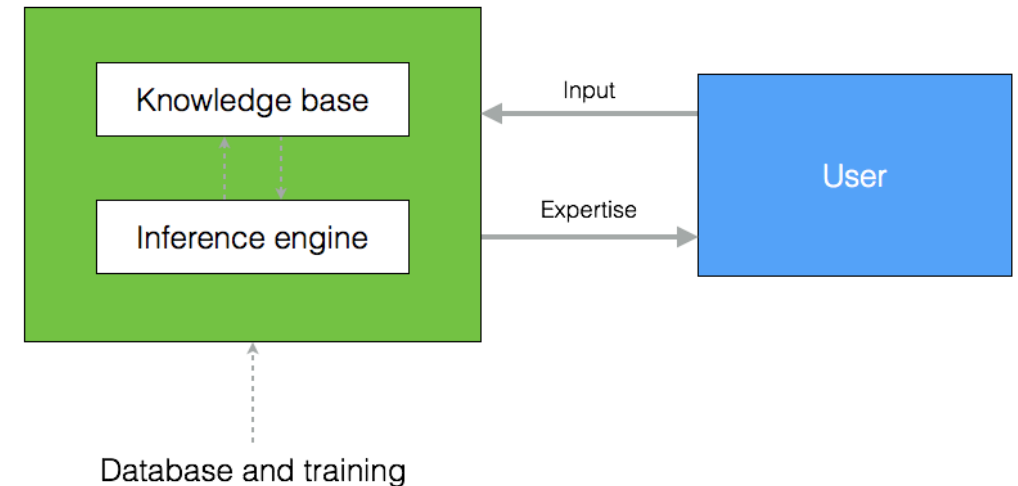- Fast, "human alike" response in real-time

AI techniques can partially cover described problems.

# Which tasks Artificial Intelligence (AI) can solve?

- Computer Vision (CV)

- Natural Language Processing (NLP)

- Speech Recognition

- Expert Systems

- Robotics



AI itself include several branches: search and optimization (e.g. mathematical optimization and evolutionary computation), classifiers and statistical learning methods (machine learning), artificial neural networks

# Types of learning

**1) Knowledge acquisition from expert**

**2) Knowledge acquisition from data:**

- *Supervised learning* – the system is supplied with a set of training examples consisting of inputs and corresponding outputs, and is required to discover the relation or mapping between them.

- *Unsupervised learning* – the system is supplied with a set of training examples consisting only of inputs. It is required to discover what appropriate outputs should be.

Everyday we face with task of recognition of handwritten symbols

Idea – create neuron network, that will recognize handwritten digits

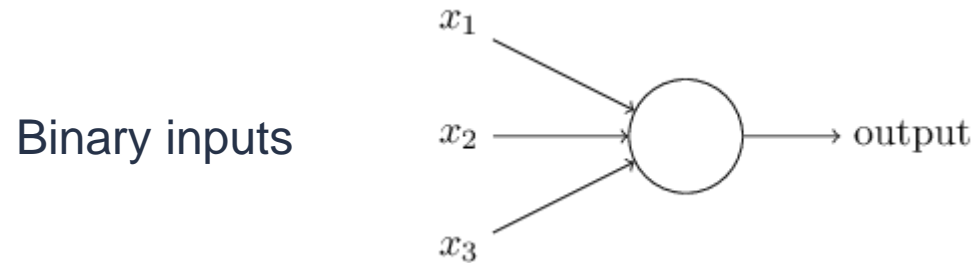**Training examples** – large number of handwritten digits

**Perceptrons** were developed in the **1950s** by the scientist **Frank Rosenblatt.**

Binary inputs

$$x_1$$
$$x_2 \longrightarrow \bigcirc \longrightarrow \text{output}$$
$$x_3$$

Idea of Rosenblatt was in weights, $w_1, w_2, \ldots$, real numbers expressing the importance of the respective inputs to the output

$$\text{output} = \begin{cases} 0 & \text{if } \sum_j w_j x_j \leq \text{ threshold} \\ 1 & \text{if } \sum_j w_j x_j > \text{ threshold} \end{cases}$$
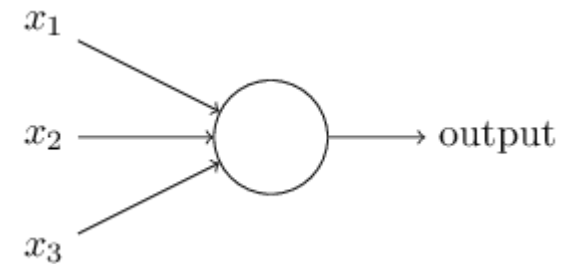
7

You're thinking about going to the cinema. You need to make a decision. You might make your decision by weighing up three factors:

1) Is there a convenient time for the exact movie?

2) Is there a good price?

3) Do you have a company to go?



Ley's represent these three factors by corresponding binary variables $x_1$, $x_2$, and $x_3$. For instance, we'd have $x_1=1$ if there is a good time, and $x_1=0$ if there is no time. For $x_2$ and $x_3$ we will do the same.

You want to go to see this movie so badly, so you don't care about company or price so much. You can use perceptrons to model this kind of decision-making. One way to do this is to choose a weight $w_1=6$ for the time, and $w_2=2$ and $w_3=2$ for the other conditions.

# Complex network of perceptron

$$\text{output} = \begin{cases} 0 & \text{if } \sum_j w_j x_j \leq \text{threshold} \\ 1 & \text{if } \sum_j w_j x_j > \text{threshold} \end{cases}$$

$w \cdot x \equiv \sum_j w_j x_j$ , where $w$ and $x$ are vectors whose components are the weights and inputs

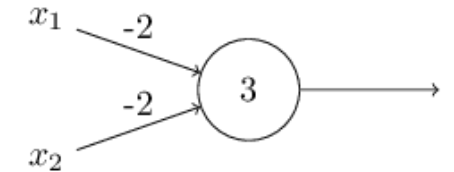$$\sum_j w_j x_j > threshold$$

$$bias, b \equiv -threshold$$

$$output = \begin{cases} 0 & if \ w^T \cdot x + b \leq 0 \\ 1 & if \ w^T \cdot x + b > 0 \end{cases}$$

## Small example



Perceptrons can be used is to compute the elementary logical functions we usually think of as underlying computation, functions such as AND, OR, and NAND. For example, suppose we have a perceptron with two inputs, each with weight −2, and an overall bias of 3.
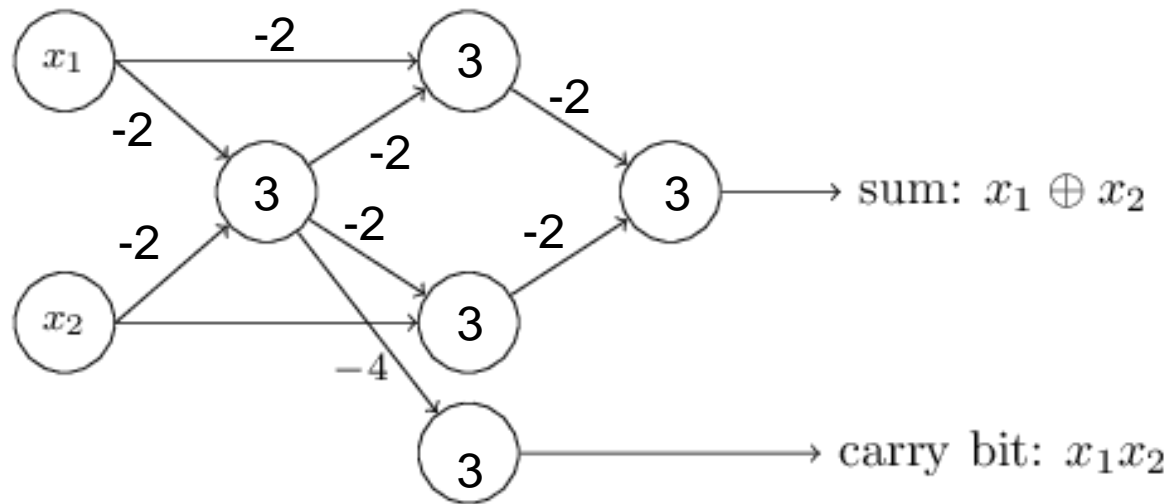
$$W = \begin{pmatrix} -2 \\ -2 \end{pmatrix}$$

Input 00 produces output 1, since $(-2)*0+(-2)*0+3=3$ is positive. Similar calculations show that the inputs 01 and 10 produce output 1. But the input 11 produces output 0, since $(-2)*1+(-2)*1+3=-1$ is negative.

$$W = \begin{pmatrix} w_{11} & w_{12} & \dots & w_{1m} \\ w_{21} & w_{22} & \dots & w_{2m} \\ \dots & \dots & \dots & \dots \\ w_{n1} & w_{n2} & \dots & w_{nn} \end{pmatrix}$$
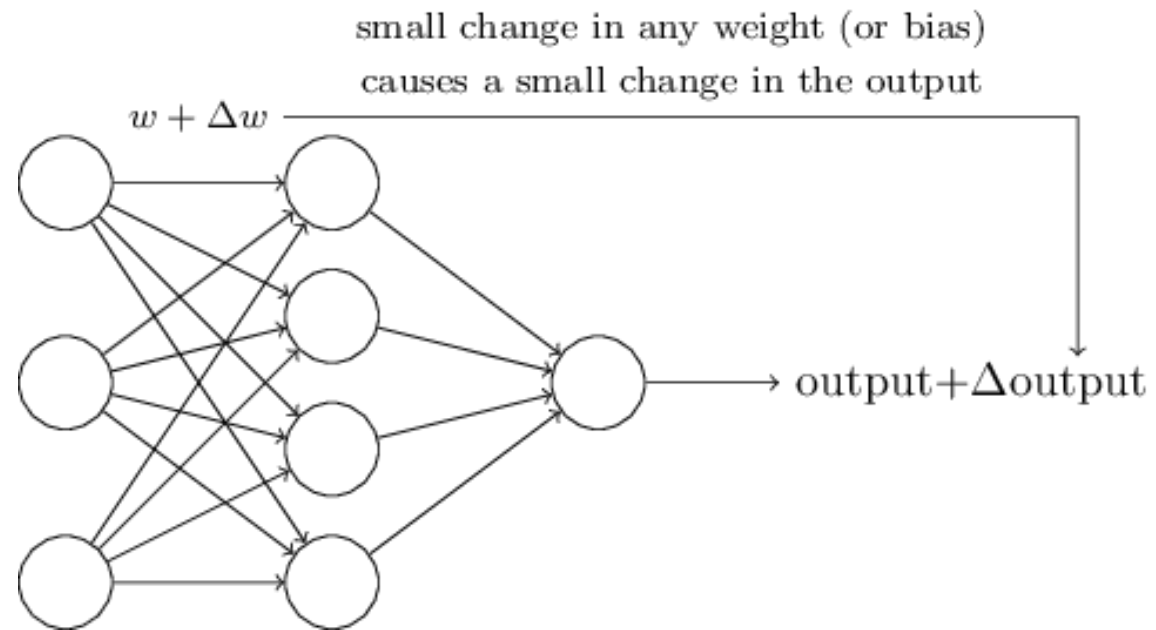
We can use networks of perceptrons to compute any logical function at all. For example, we can use NAND gates to build a circuit which adds two bits, $x_1$ and $x_2$. This requires computing the bitwise sum, $x_1 \oplus x_2$. Let's calculate outputs of the network.



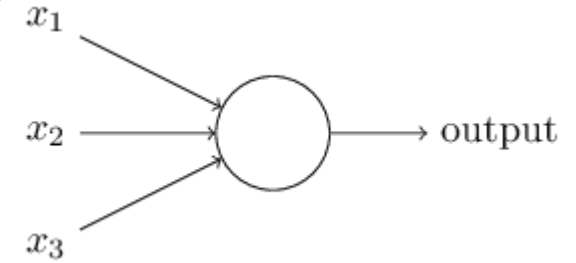| $a$ | $b$ | $a \oplus b$ |
|---|---|---|
| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 1 | 0 |

small change in any weight (or bias)
causes a small change in the output

$w + \Delta w$

output+$\Delta$output

# SIGMOID NEURONS

Sigmoid neurons are similar to perceptrons, but modified so that small changes in their weights and bias cause only a small change in their output.
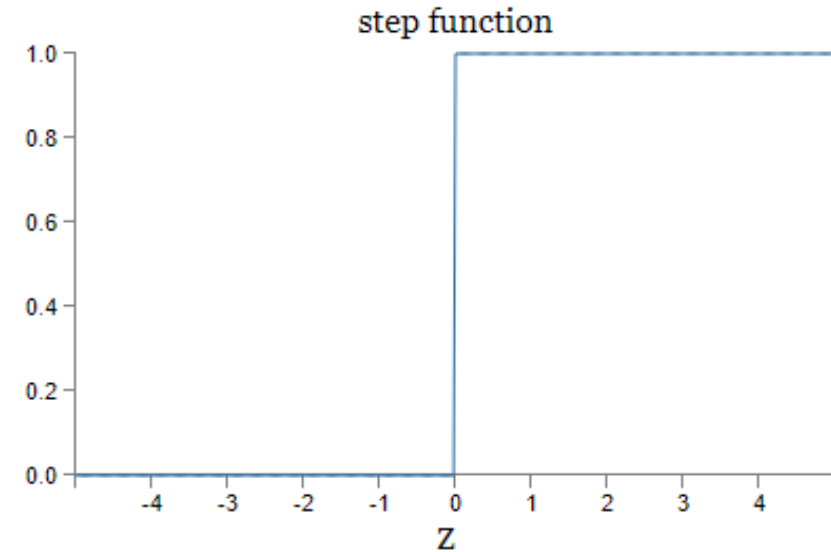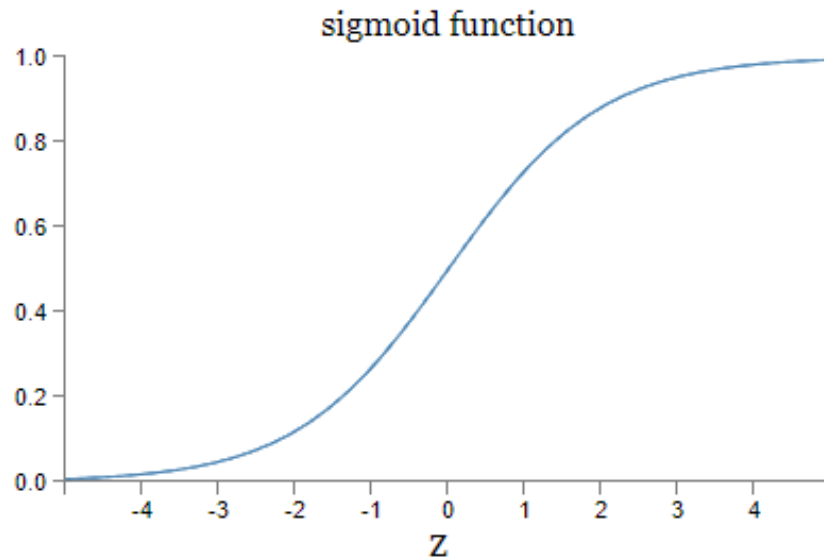
1) Just like a perceptron, the sigmoid neuron has inputs, $x_1, x_2, \ldots$. But instead of being just 0 or 1, these inputs can also take on any values between 0 and 1. So, for instance, 0.638 is a valid input for a sigmoid neuron.

2) The sigmoid neuron has weights for each input, $w_1, w_2, \ldots$, and an overall *bias, b*. But the output is not 0 or 1. Instead, it's *σ(z), where z≡w·x+b and σ* is called the ***sigmoid function****

$$\sigma(z) \equiv \frac{1}{1 + e^{-z}}$$

$$\longrightarrow$$

$$\sigma(z) \equiv \frac{1}{1 + \exp\left(-\sum_j w_j x_j - b\right)}$$

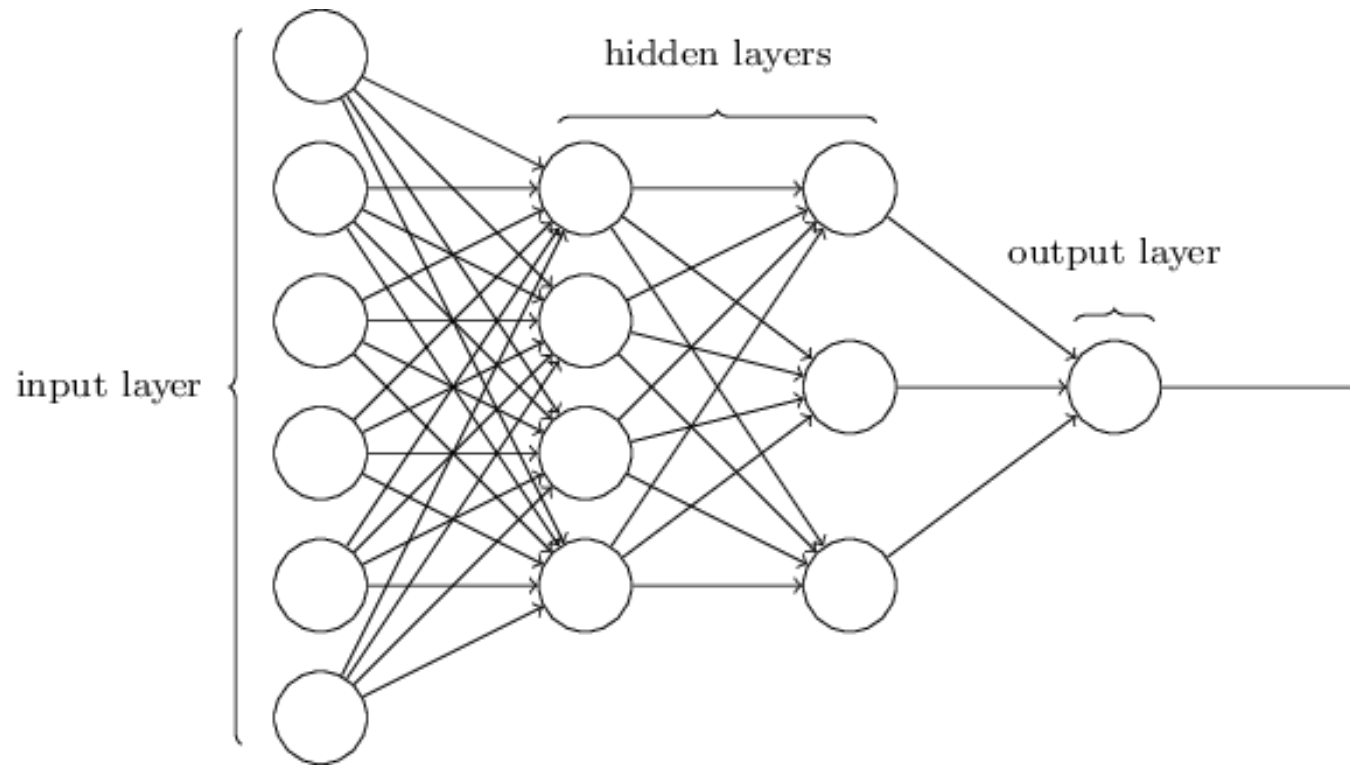*sometimes called the logistic function, and this class of neurons called logistic neurons

sigmoid function



step function

$$\Delta output \approx \sum_{j} \frac{\partial\ output}{\partial\ w_j} \Delta w_j + \frac{\partial\ output}{\partial\ b} \Delta b$$

$\Delta output$ is a **linear function** of the changes $\Delta w_j$ and $\Delta b$ in the weights and bias

13

hidden layers

input layer

output layer

multilayer perceptrons or MLPs

$$\sigma(w_1a_1+w_2a_2+w_3a_3+w_4a_4+\dots+w_na_n + bias)$$

Only activate meaningfully
when weighted sum > **threshold**

$x$ is training input - 28×28=784-dimensional vector, where each entry in the vector represents the grey value for a single pixel in the image

$y = y(x)$ is a desired output - 10-dimensional vector ( $y(x)$=(0,0,0,0,0,0,1,0,0,0)$^T$ for digit 6)

For determination of how good we are in approximating $y(x)$ we will define cost function*:

$$C(w, b) = \frac{1}{2n} \sum_x \|y(x) - a\|^2$$

where $w$ denotes the collection of all weights in the network, $b$ all the biases, $n$ is the total number of training inputs, $a$ is the vector of outputs from the network when $x$ is input, and the sum is over all training inputs, x. In the other words – the smaller $C(w, b)$≈0 means that training algorithm has done a good job

*There are 3 different notions: loss function, cost function and objective function.
Loss function is usually a function defined on a data point, prediction and label, and measures the penalty.
Cost function is usually more general. It might be a sum of loss functions over your training set plus some model complexity penalty (regularization). Objective function is the most general term for any function that you optimize during training.

|  | Neural network function | Cost function |
|---|---|---|
| Input | 784 numbers (pixels) | 11 935 |
| Output | 10 numbers | 1 number (the cost or the error) |
| Parameters | 11 935 | Many training examples |

$$(784 \times 15 + 15 \times 10) + (15 + 10) = 11\ 935$$
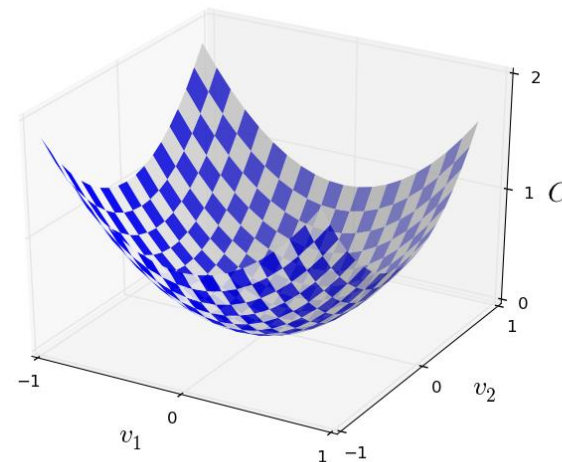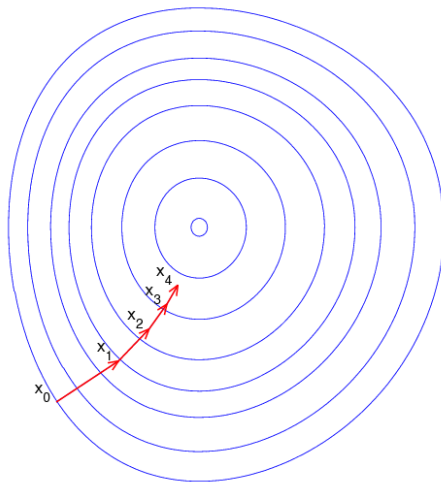weights          biases

So our goal in training a neural network is to find weights and biases, which minimize the quadratic cost function

$$C(w, b) = \frac{1}{2n} \sum_{x} \|y(x) - a\|^2$$

Let's replace $w, b$ with some variables $v = v_1, v_2, \ldots$ => $C(v)$. To minimize $C(v)$ it helps to imagine $C$ as a function of just two variables, which is $v_1$ and $v_2$:

In "moving a ball" for a small amount $\Delta v_1$ in the $v_1$ direction, and a small amount $\Delta v_2$ in the $v_2$ direction. C changes the follows:
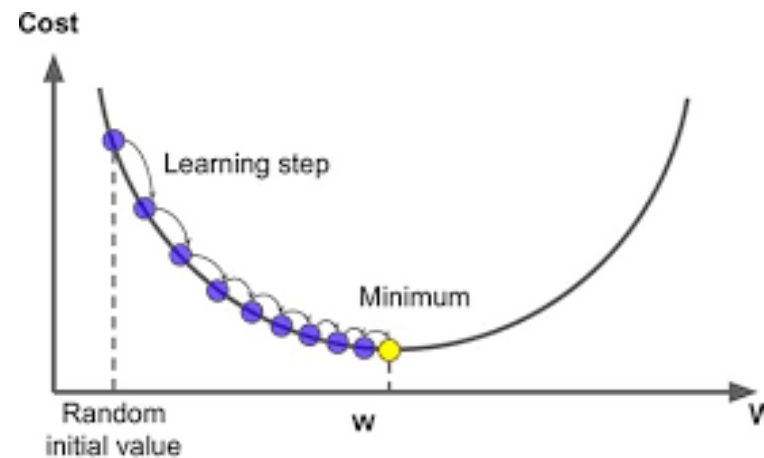
$$\Delta C \approx \frac{\partial C}{\partial v_1} \Delta v_1 + \frac{\partial C}{\partial v_2} \Delta v_2$$

We need to find the way of choosing $\Delta v_1$ and $\Delta v_2$ to make $\Delta C$ negative
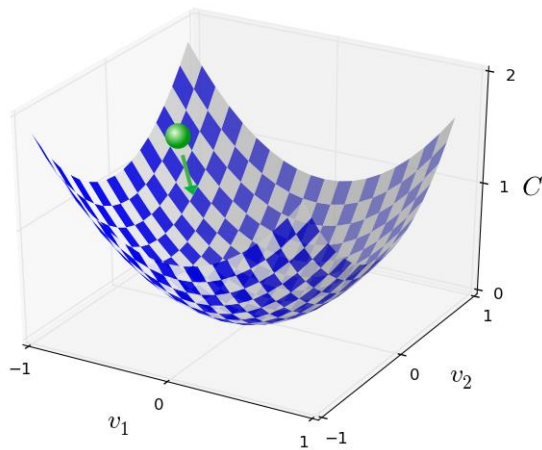
We can re-write this in vector form:

$$\Delta C \approx \frac{\partial C}{\partial v_1}\Delta v_1 + \frac{\partial C}{\partial v_2}\Delta v_2 \rightarrow \Delta C \approx \nabla C \cdot \Delta v$$

$$\Delta v = (\Delta v_1, \Delta v_2)^T \qquad \text{- vector of changes}$$

$$\nabla C = \left(\frac{\partial C}{\partial v_1}, \frac{\partial C}{\partial v_2}\right)^T \qquad \text{- gradient vector}$$

We need to find the way of choosing $\Delta v_1$ and $\Delta v_2$ to make $\Delta C$ negative.

$$\Delta v = -\eta \cdot \nabla C \rightarrow \Delta C \approx -\eta \cdot ||\nabla C||^2$$

$\eta$ – learning rate

$v' = v - \eta \cdot \nabla C$ – update rule

22

$$\vec{\mathbf{W}} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_{13,000} \\ w_{13,001} \\ w_{13,002} \end{bmatrix}$$

$$w_k \rightarrow w_k' = w_k - \eta \frac{\partial C}{\partial w_k}$$

$$b_l \rightarrow b_l' = b_l - \eta \frac{\partial C}{\partial b_l}.$$

$$w_k \rightarrow w_k' = w_k - \frac{\eta}{m} \sum_j \frac{\partial C_{X_j}}{\partial w_k}$$

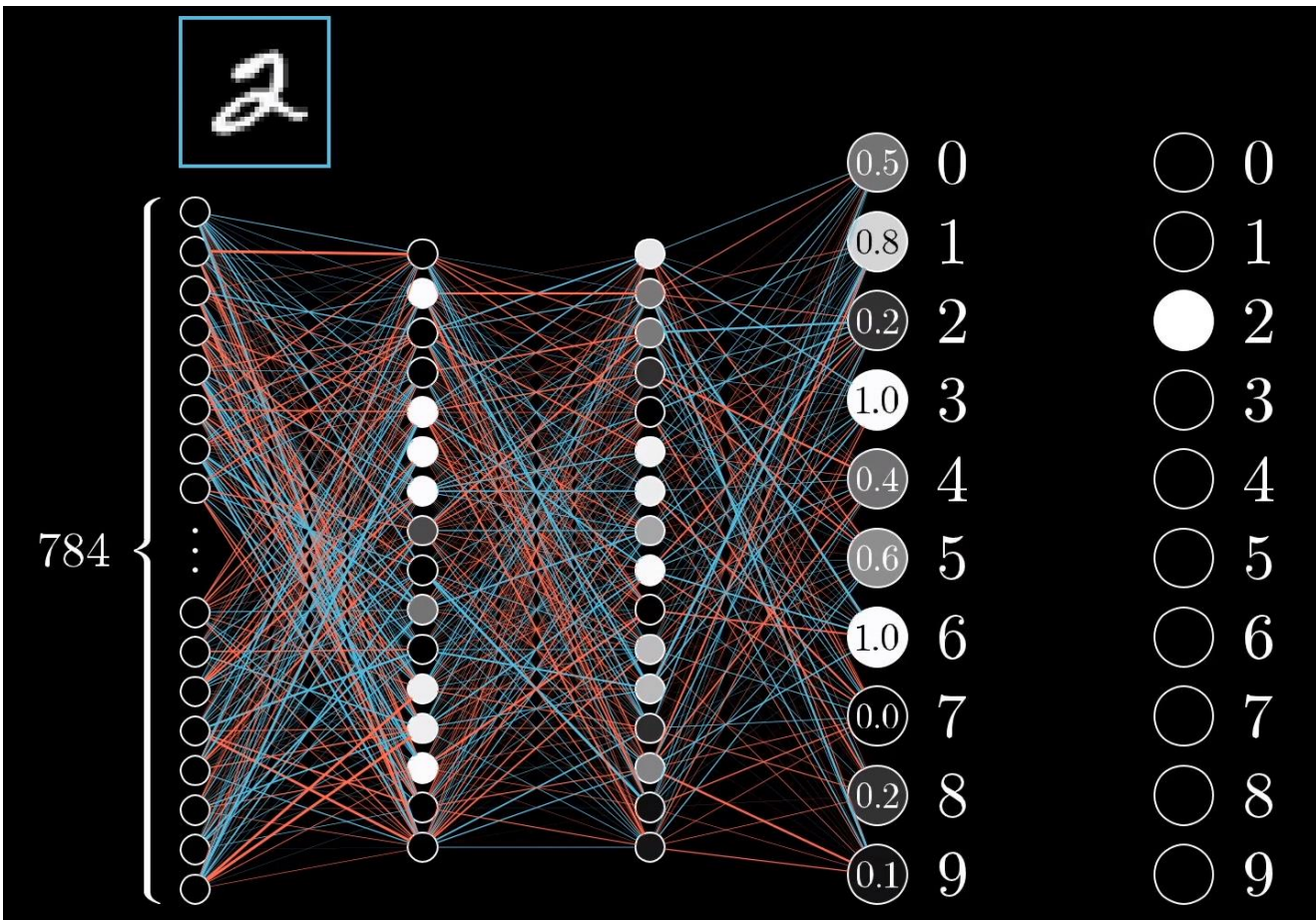$$b_l \rightarrow b_l' = b_l - \frac{\eta}{m} \sum_j \frac{\partial C_{X_j}}{\partial b_l},$$

$$-\nabla C(\vec{\mathbf{W}}) = \begin{bmatrix} 0.31 \\ 0.03 \\ -1.25 \\ \vdots \\ 0.78 \\ -0.37 \\ 0.16 \end{bmatrix}$$

$w_0$ should increase somewhat

$w_1$ should increase a little

$w_2$ should decrease a lot

$w_{13,000}$ should increase a lot

$w_{13,001}$ should decrease somewhat

$w_{13,002}$ should increase a little

After first training we pick out another randomly chosen mini-batch and train with those. And so on, until we've exhausted the training inputs, which is said to complete an **epoch** of training. At that point we start over with a new training epoch.
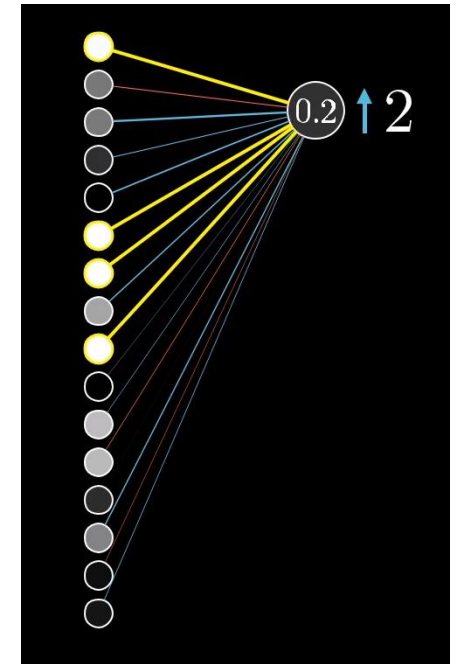
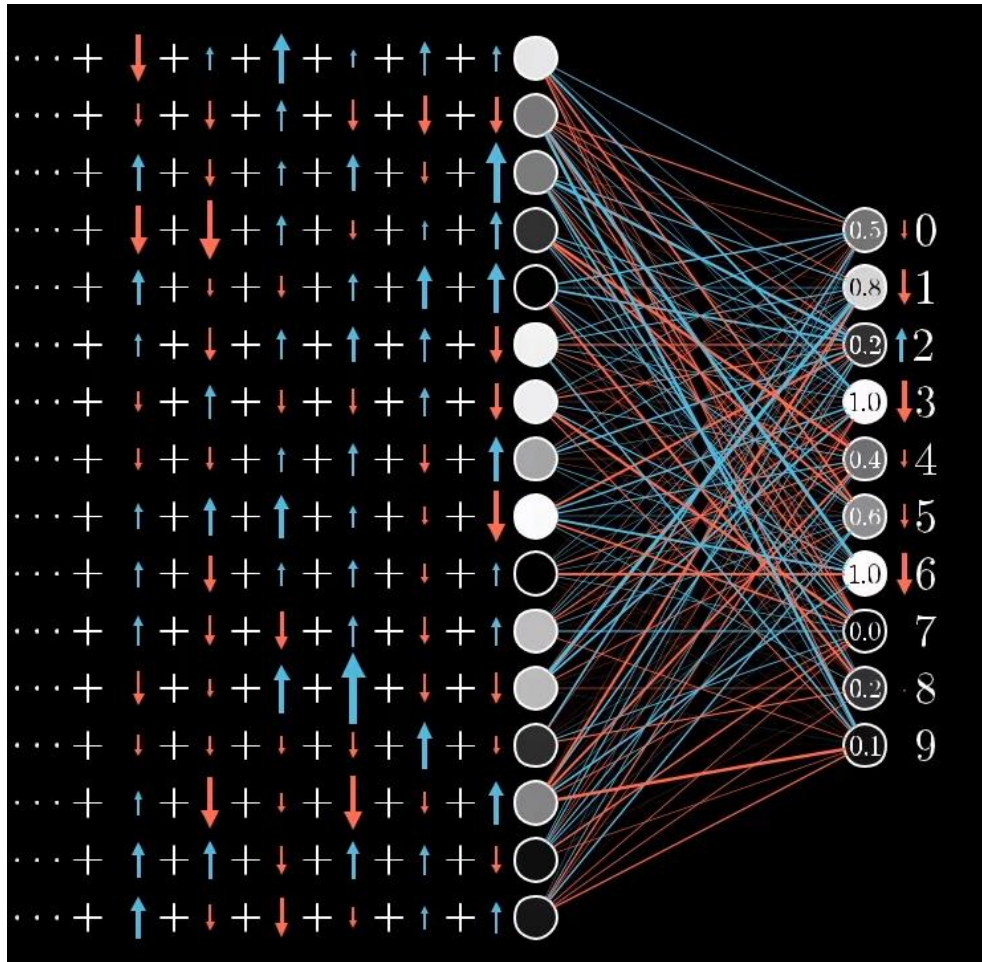$$0.2 = \sigma(w_1 a_1 + w_2 a_2 + w_3 a_3 + w_4 a_4 + \ldots + w_n a_n + bias)$$

- Increase b

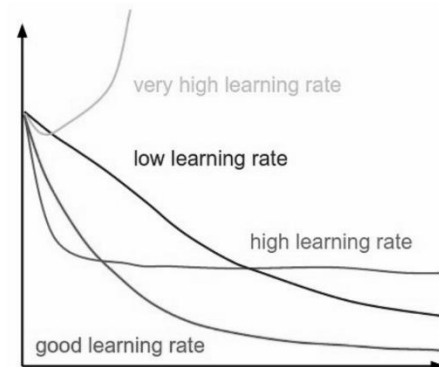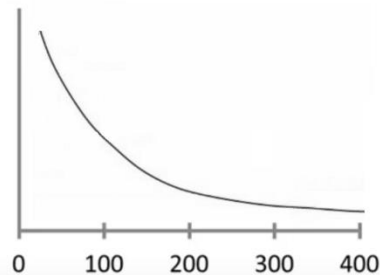- Increase $w_i$ in proportion to $a_i$

- Change $a_i$ in proportion to $w_i$

$- \eta \cdot \nabla C(W)$ calculates over all training data

Types of Gradient Descent:
- **Batch** - calculates the error for each example within the training dataset, but only after all training examples have been evaluated, the model gets updated. This whole process is like a cycle and called a training epoch.
- **Stochastic** - it updates the parameters for each training example, one by one. This can make SGD faster than Batch Gradient Descent
- **Mini Batch** - it's a combination of the concepts of SGD and Batch Gradient Descent. It simply splits the training dataset into small batches and performs an update for each of these batches. Therefore it creates a balance between the robustness of stochastic gradient descent and the efficiency of batch gradient descent.



26

$$f(x, y) = x^2 + xy + (x + y)^2$$



a     б     в

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial a}\frac{\partial a}{\partial x} + \frac{\partial f}{\partial b}\frac{\partial b}{\partial x} + \frac{\partial f}{\partial c}\frac{\partial c}{\partial x} = 2x + y + 2(x + y)$$

$$f(x, y) = x^2 + xy + (x + y)^2$$

## Setting up the environment

- Download and install Anaconda Data Science Libraries 3.6
- https://www.anaconda.com/download/

- Install Anaconda and add it to the Path variable

$$\boxed{\frac{\partial C_0}{\partial a_k^{(L-1)}}} = \underbrace{\sum_{j=0}^{n_L-1} \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial C_0}{\partial a_j^{(L)}}}_{\text{Sum over layer L}}$$

$$z_j^{(L)} = \cdots + w_{jk}^{(L)} a_k^{(L-1)} + \cdots$$

$$a_j^{(L)} = \sigma(z_j^{(L)})$$

$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$

$a_0^{(L)}$

$a_1^{(L)}$

0.94  0.10  0.83  0.80  0.40  0.63  0.98  0.81

## Setting up the environment

■ Go to the Anaconda Prompt (Start->Programs->Anaconda3->Anaconda Prompt)

*conda create -n venv pip python=3.6*

We are creating virtual enviroment to isolate package installation from the system

■Then activate the enviroment by *conda activate venv*

■And then install simple version of tensorflow (for CPU only) by

*conda --upgrade tensorflow*

```
Anaconda Prompt                                                    —    □    ×

(base) C:\Users\AZMakhmutova>conda activate venv

(venv) C:\Users\AZMakhmutova>pip install --upgrade tensorflow
Collecting tensorflow
  Downloading https://files.pythonhosted.org/packages/c1/e1/a5693a158f3867417d7b50bd0514c83304706242fae74463705b8c373777/tensorflow-1.11
.0-cp36-cp36m-win_amd64.whl (46.9MB)
    100% |████████████████████████████████| 46.9MB 319kB/s
Collecting gast>=0.2.0 (from tensorflow)
  Using cached https://files.pythonhosted.org/packages/5c/78/ff794fcae2ce8aa6323e789d1f8b3b7765f601e7702726f430e814822b96/gast-0.2.0.tar
.gz
Collecting absl-py>=0.1.6 (from tensorflow)
  Downloading https://files.pythonhosted.org/packages/16/db/cce5331638138c178dd1d5fb69f3f55eb3787a12efd9177177ae203e847f/absl-py-0.5.0.t
ar.gz (90kB)
    100% |████████████████████████████████| 92kB 2.9MB/s
Collecting keras-preprocessing>=1.0.3 (from tensorflow)
  Downloading https://files.pythonhosted.org/packages/fc/94/74e0fa783d3fc07e41715973435dd051ca89c550881b3454233c39c73e69/Keras_Preproces
sing-1.0.5-py2.py3-none-any.whl
Collecting six>=1.10.0 (from tensorflow)
  Using cached https://files.pythonhosted.org/packages/67/4b/141a581104b1f6397bfa78ac9d43d8ad29a7ca43ea90a2d863fe3056e86a/six-1.11.0-py2
.py3-none-any.whl
Collecting termcolor>=1.1.0 (from tensorflow)
  Using cached https://files.pythonhosted.org/packages/8a/48/a76be51647d0eb9f10e2a4511bf3ffb8cc1e6b14e9e4fab46173aa79f981/termcolor-1.1.
0.tar.gz
Collecting astor>=0.6.0 (from tensorflow)
  Downloading https://files.pythonhosted.org/packages/35/6b/11530768cac581a12952a2aad00e1526b89d242d0b9f59534ef6e6a1752f/astor-0.7.1-py2
.py3-none-any.whl
Collecting tensorboard<1.12.0,>=1.11.0 (from tensorflow)
  Downloading https://files.pythonhosted.org/packages/9b/2f/4d788919b1feef04624d63ed6ea45a49d1d1c834199ec50716edb5d310f4/tensorboard-1.1
1.0-py3-none-any.whl (3.0MB)
    100% |████████████████████████████████| 3.0MB 5.7MB/s
Collecting protobuf>=3.6.0 (from tensorflow)
  Downloading https://files.pythonhosted.org/packages/e8/df/d606d07cff0fc8d22abcc54006c0247002d11a7f2d218eb008d48e76851d/protobuf-3.6.1-
```

■After successfull installation deactivate the enviroment by

*conda deactivate*

## Setting up the environment

Install PyCharm - Python IDE from
https://www.jetbrains.com/pycharm/download/#section=windows


As a student you can have free individual licenses (fastest way with your KAI email)
https://www.jetbrains.com/student/


Install PyCharm by default settings

Create a new project with Project Interpreter from virtual environment

# Setting up the environment

Choose Project Interpreter from virtual environment that we just created

# Setting up the environment

■You can always check/change the current interpreter in the project

■*File -> Settings -> Project:Somename ->Project Interpreter*

■Let's check the work of tensorflow. Create simple python file in the project and run following code:

```python
import tensorflow as tf


hello = tf.constant ('Hello, TensorFlow!')

sess = tf.Session()
print(sess.run(hello))

a = tf.constant(10)
b = tf.constant(32)

print (sess.run(a+b))
```

■The output should be 42. If so, you set up the environment!

# Variables in tf :

```python
w = tf.Variable(tf.random_normal([3, 2], mean=0.0, stddev=0.4), name='weights')
b = tf.Variable(tf.zeros([2]), name='biases')



with tf.device('/gpu:0'):
w = tf.Variable(tf.random_normal([3, 2], mean=0.0, stddev=0.4), name='weights')



with tf.device('/job:ps/task:0'):
w = tf.Variable(tf.random_normal([3, 2], mean=0.0, stddev=0.4), name='weights')
```

Initialization of the variables in tf :

```
init = tf.initialize_global_variables()


w2 = tf.Variable(w.initialized_value(), name='w2')



saved = saver.save(sess, 'model.ckpt')

saver.restore('model.ckpt')
```

■tf.placeholder:

```python
import tensorflow as tf

x = tf.placeholder(tf.float32)
y = tf.placeholder(tf.float32)
output = tf.multiply(x, y)
with tf.Session() as sess:
    result = sess.run(output, feed_dict={x: 2, y: 3})

print(result)
```

■broadcasting:

```
m = tf.Variable(tf.random_normal([10, 100], mean=0.0, stddev=0.4), name='matrix')
v = tf.Variable(tf.random_normal([100], mean=0.0, stddev=0.4), name='vector')
result = m + v

print(result)
```

◼variable scopes:

◼tf.get_variable(<name>, <shape>, <initializer>) creates and returns a variable with the specified name;

◼• tf.variable_scope(<scope_name>) manages the namespaces used in tf.get_variable();

```python
def linear_transform(vec, shape):
    with tf.variable_scope('transform'):
        w = tf.get_variable('matrix', shape,
        initializer=tf.random_normal_initializer())
    return tf.matmul(vec, w)


with tf.variable_scope('linear_transformers') as scope:
        result1 = linear_transform(vec1, shape)
        scope.reuse_variables()
        result2 = linear_transform(vec2, shape)
```

# Linear regression training

```python
import numpy as np,tensorflow as tf

n_samples, batch_size, num_steps = 1000, 100, 20000
X_data = np.random.uniform(1, 10, (n_samples, 1))
y_data = 2 * X_data + 1 + np.random.normal(0, 2, (n_samples, 1))
X = tf.placeholder(tf.float32, shape=(batch_size, 1))
y = tf.placeholder(tf.float32, shape=(batch_size, 1))

with tf.variable_scope('linear-regression'):
    k = tf.Variable(tf.random_normal((1, 1), stddev=0.0), name='slope')
    b = tf.Variable(tf.zeros((1,)), name='bias')

y_pred = tf.matmul(X, k) + b
loss = tf.reduce_mean((y - y_pred) ** 2)
optimizer = tf.train.GradientDescentOptimizer(0.001).minimize(loss)

display_step = 50

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    for i in range(num_steps):
        indices = np.random.choice(n_samples, batch_size)
        X_batch, y_batch = X_data[indices], y_data[indices]
        _, loss_val, k_val, b_val = sess.run([ optimizer, loss, k, b ],
            feed_dict = { X : X_batch, y : y_batch })
        if (i+1) % display_step == 0:
            print('Epoch %d: %.8f, k=%.4f, b=%.4f' % (i+1, loss_val, k_val, b_val))
```

$$L = \sum_{i=1}^{N} (\hat{y} - y)^2 \to \min$$

*f = kx + b, k=2, b=1*

1) Create 1000 random points in [0:1]
2) Calculate for each point x the correct answer $y = 2x+1+\epsilon$, $\epsilon$  $N(\epsilon; 0; 2$- variance)
3) Initialization of k and b
4) Calculating the initial loss function
5) Calculating the optimizer, learning rate = 0.001
6) Optimization cycle

MNIST data set - MNIST's name comes from the fact that it is a modified subset of two data sets collected by NIST, the United States' National Institute of Standards and Technology

http://yann.lecun.com/exdb/mnist/

Training data (set) is used for train the NN, but tests are made on test data!

# Handwriting digits recognition

```python
import tensorflow as tf

from tensorflow.examples.tutorials.mnist import input_data
mnist = input_data.read_data_sets("MNIST_data/", one_hot=True)

x = tf.placeholder(tf.float32, [None, 784])

W = tf.Variable(tf.zeros([784, 10]))
b = tf.Variable(tf.zeros([10]))

y = tf.nn.softmax(tf.matmul(x, W) + b)

y_ = tf.placeholder(tf.float32, [None, 10])

cross_entropy = tf.reduce_mean(-tf.reduce_sum(y_ * tf.log(y), reduction_indices=[1]))

train_step = tf.train.GradientDescentOptimizer(0.5).minimize(cross_entropy)

init = tf.global_variables_initializer()
sess = tf.Session()
sess.run(init)
for i in range(1000):
    batch_xs, batch_ys = mnist.train.next_batch(100)
    sess.run(train_step, feed_dict={x: batch_xs, y_: batch_ys})

correct_prediction = tf.equal(tf.argmax(y, 1), tf.argmax(y_, 1))
accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
print("Accuracy: %s" % sess.run(accuracy, feed_dict={x: mnist.test.images, y_: mnist.test.labels}))
```