

CAS Machine Intelligence: Deep Learning

Week 4

Beate Sick

Institut für Datenanalyse und Prozessdesign

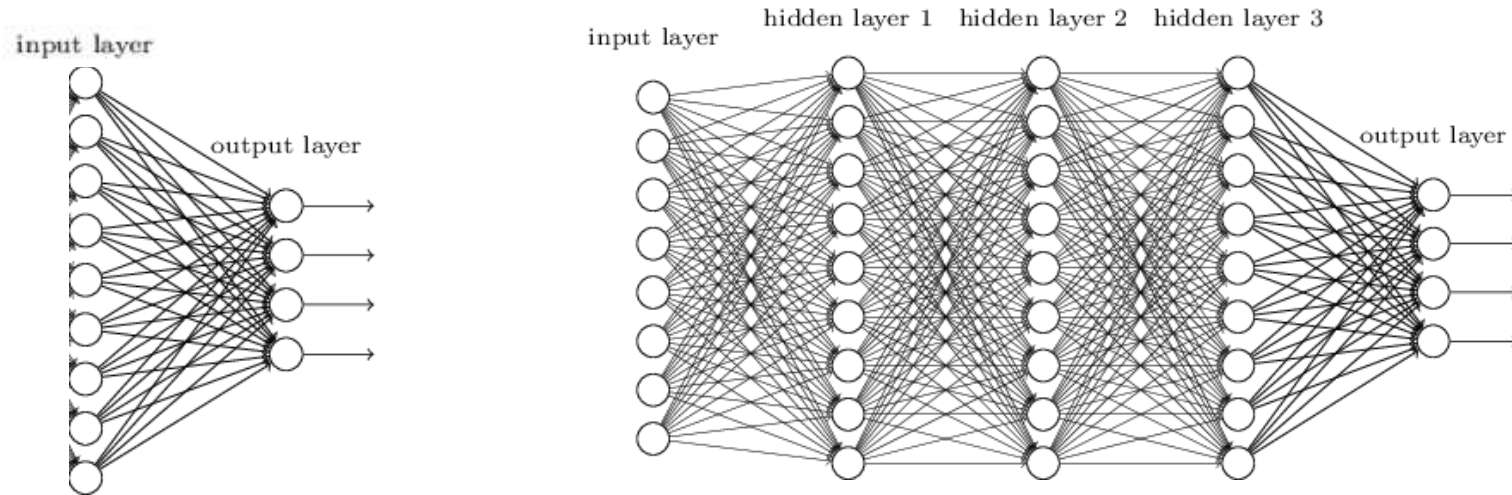
Zürcher Hochschule für Angewandte Wissenschaften

Beate.Sick@zhaw.ch

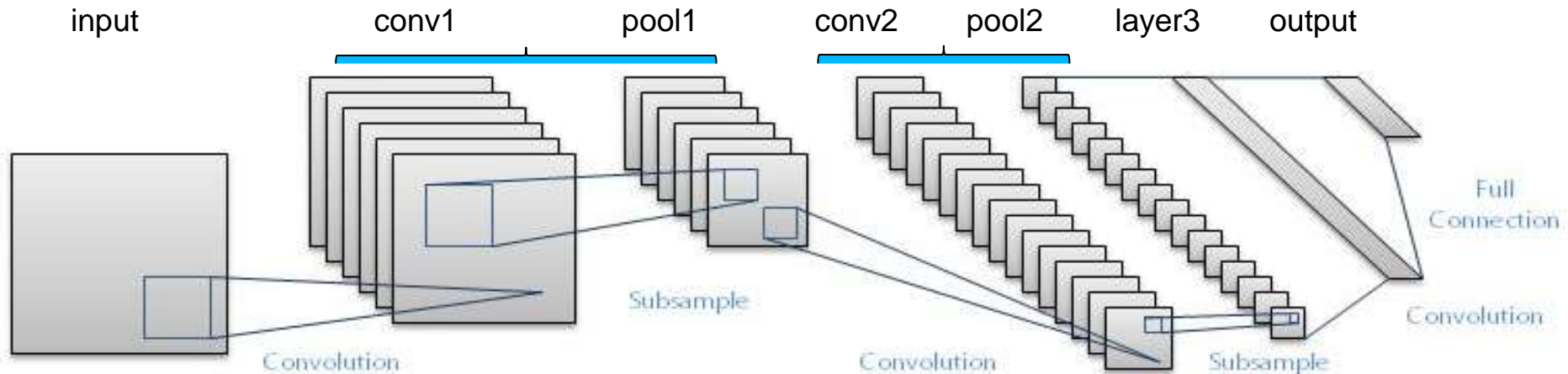
- Recap fully connected NN
- Keras intro (recap)
- Homework: Questions and results
- Motivation of convolutional neural networks (CNNs)
- What is convolution?
- How is convolution performed over several channels/stack of images?
- How does a classical CNN look like?
- Do a CNN yourself

We will go from fully connected NNs to CNNs

Fully connected Neural Networks (fcNN) without and with hidden layers:



Convolutional Neural Network:



At the end of the day

Develop a DL model to solve this task:

For a given image from the internet, decide which out of 8 celebrities is on the image.

Example images:

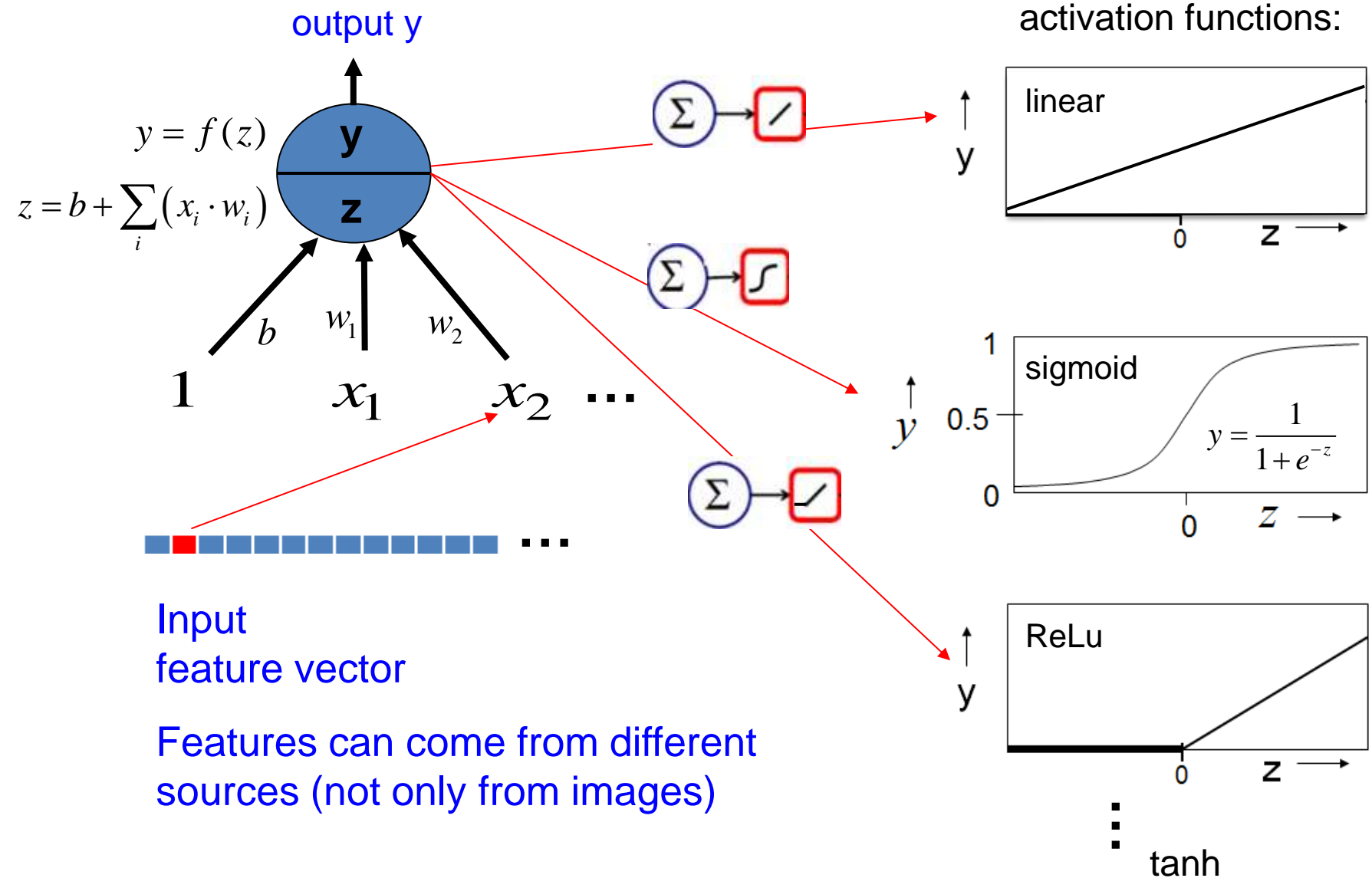
Label: Steve Jobs (entrepreneur)



Label: Emma Stone (actress)



Recap: computations performed in a neuron



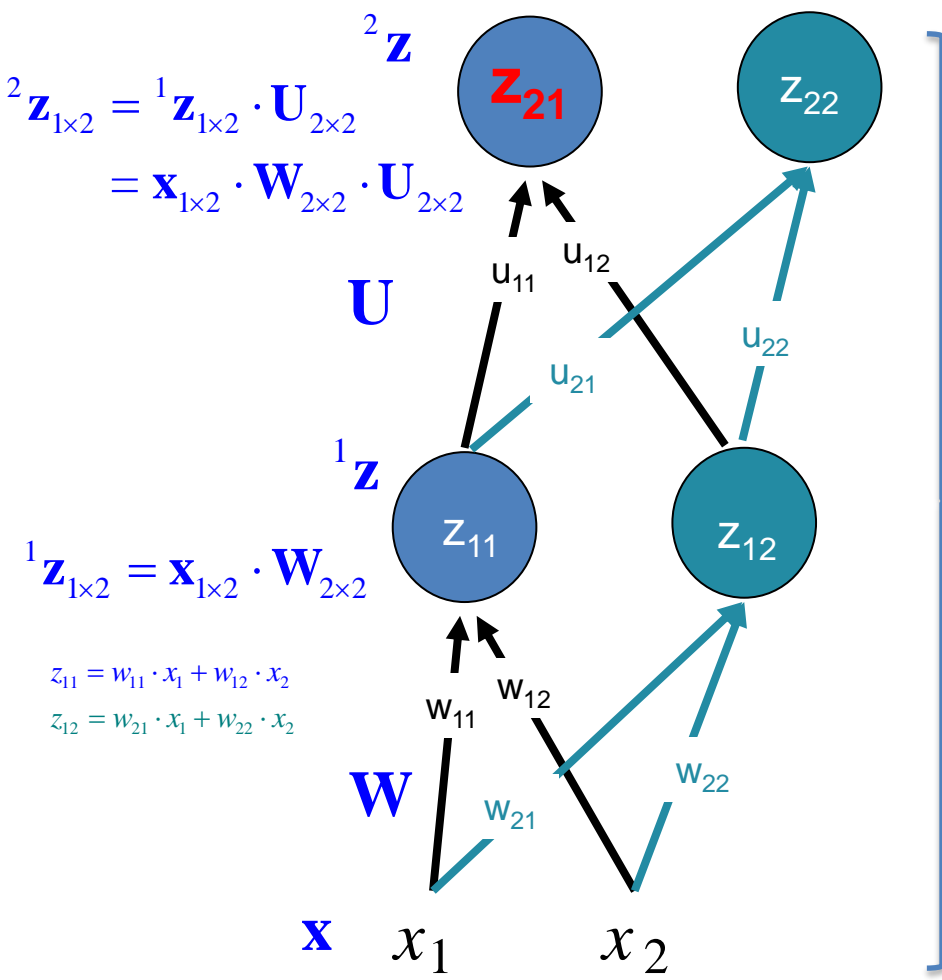


2 linear layers can be replaced by 1 linear layer

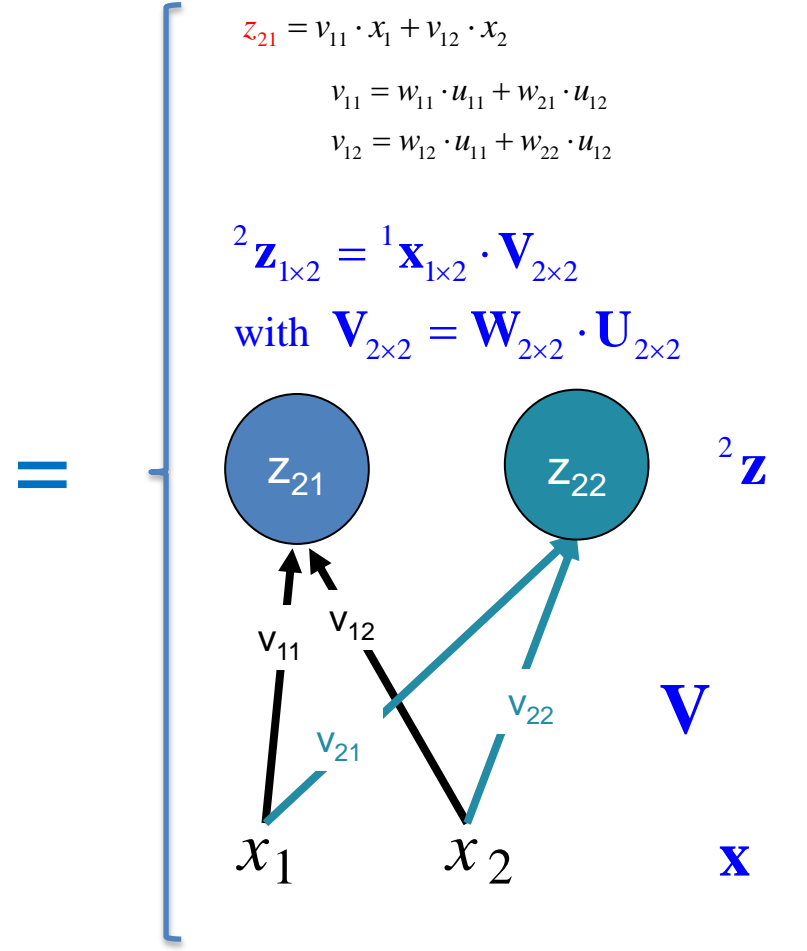
-> we can not go deep with linear layers!

$$z_{21} = z_{11} \cdot u_{11} + z_{12} \cdot u_{12} = (w_{11} \cdot x_1 + w_{12} \cdot x_2) \cdot u_{11} + (w_{21} \cdot x_1 + w_{22} \cdot x_2) \cdot u_{12}$$

$$= x_1 \cdot (w_{11} \cdot u_{11} + w_{21} \cdot u_{12}) + x_2 \cdot (w_{12} \cdot u_{11} + w_{22} \cdot u_{12})$$

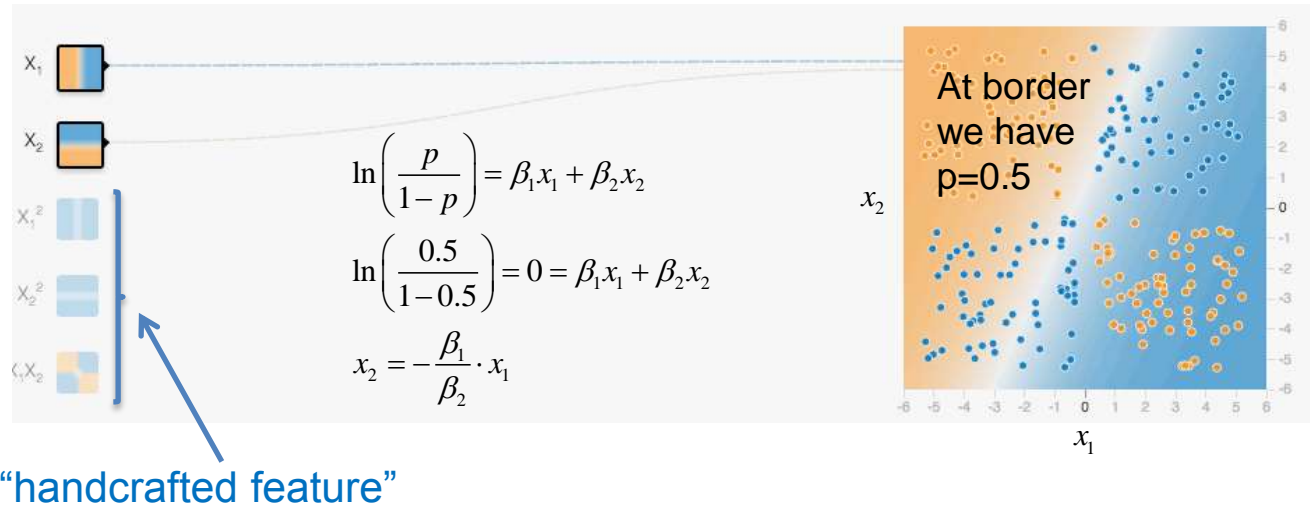
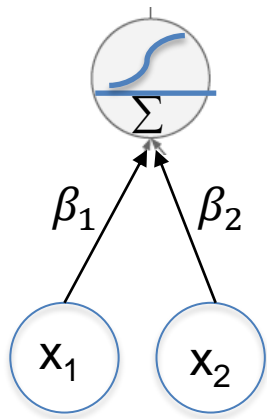


The activations of two stacked linear layers are linear functions of the inputs x_i that can also be achieved with one linear layer.

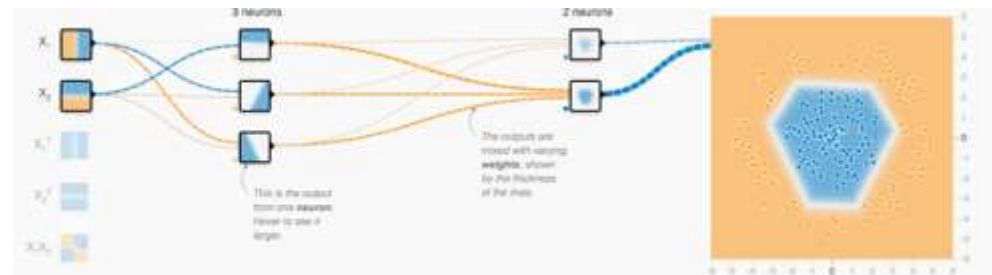
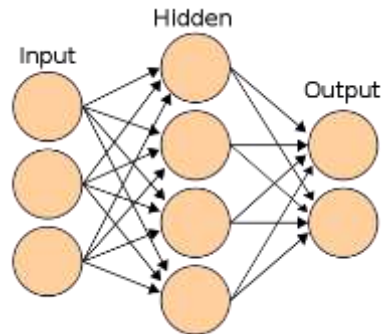


NN can only go deep with non-linear activation functions

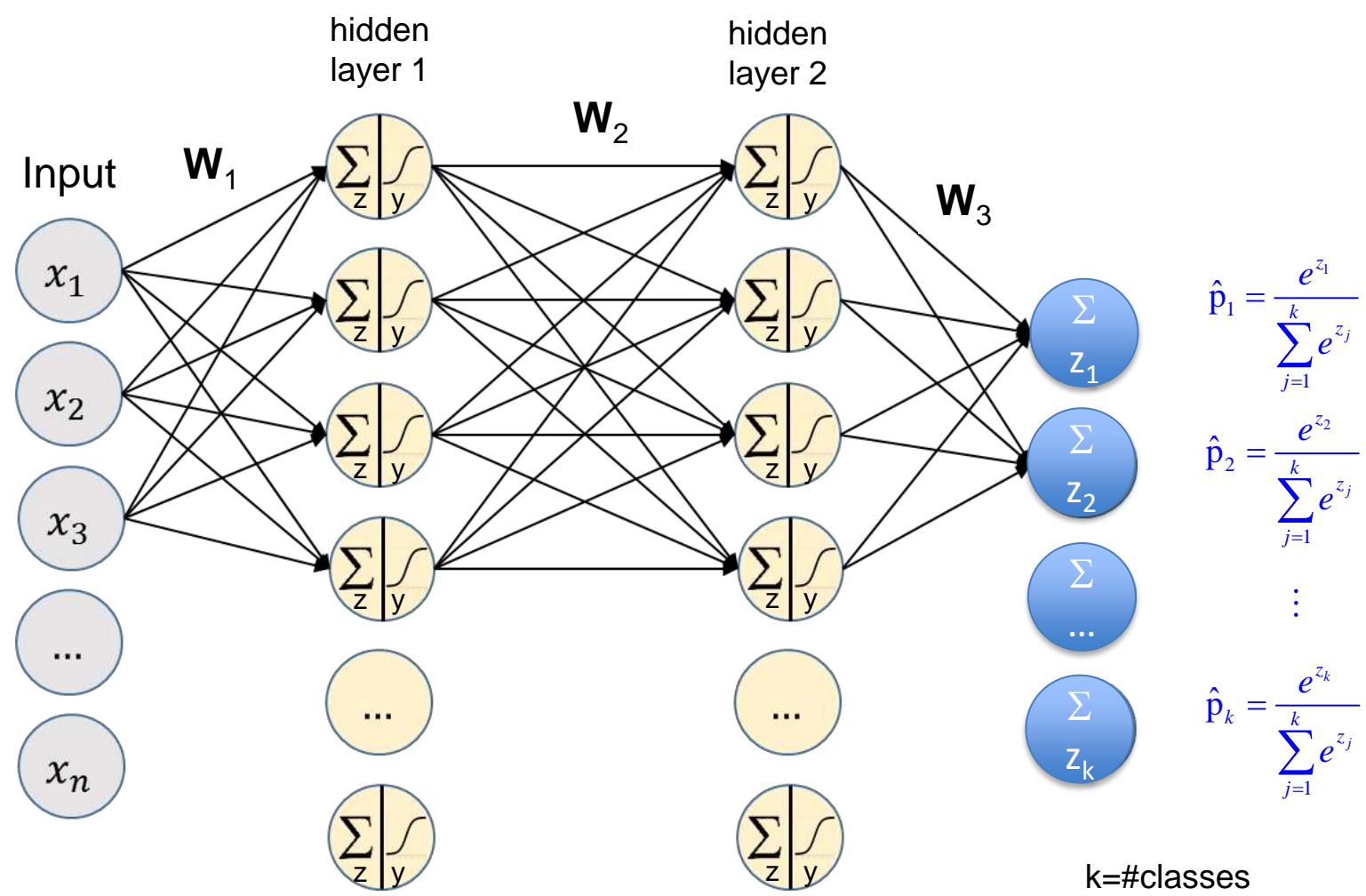
NN **without** hidden layer and sigmoid activation function yields **linear separation** curve.



NN with ≥ 1 hidden layer and sigmoid activation function yields **arbitrary separation** curve.



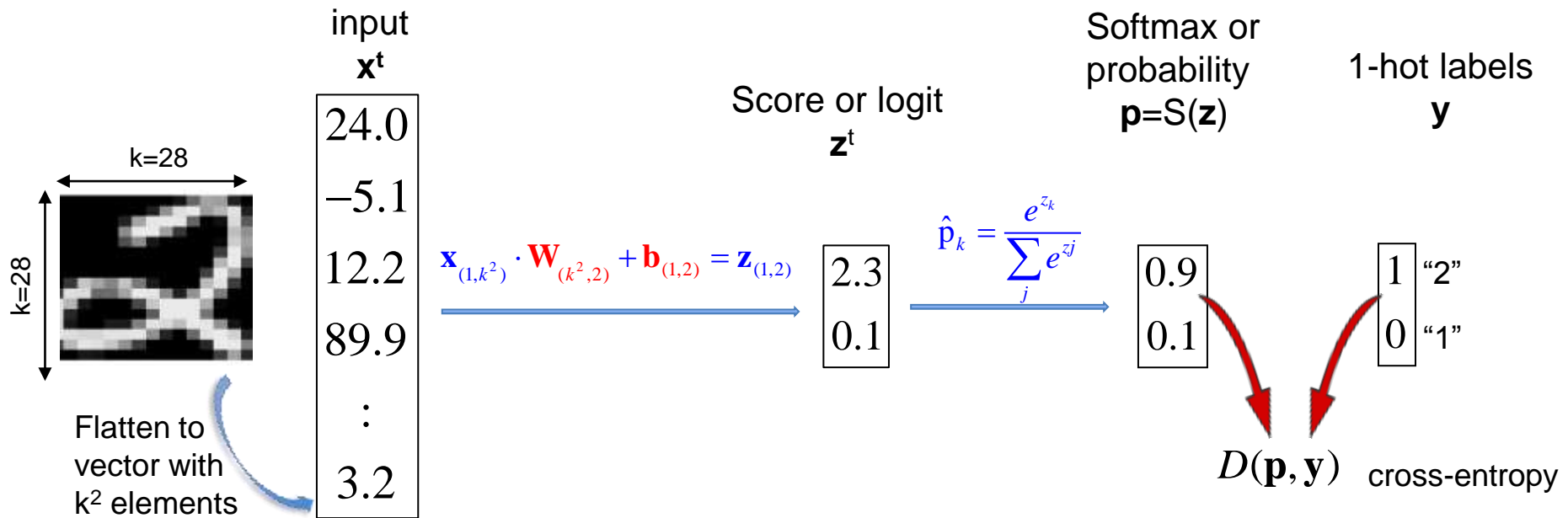
A fully conneted neural networks with 2 hidden layers



$$z = b + \sum_i (x_i \cdot w_i) \quad y = f(z)$$

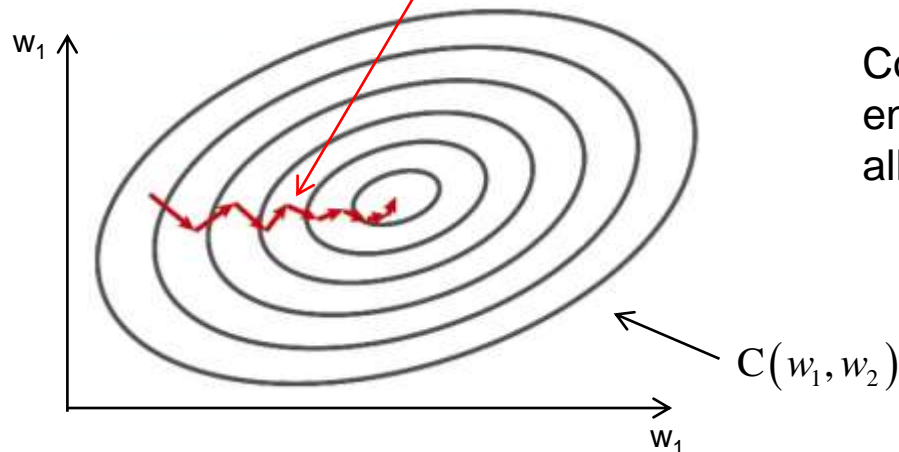
Remark: weight values in the weight matrices that are learned during training.

What is going on in our 1 layer fully connected NN?



Take step in direction of descent gradient:
(the gradient is oriented orthogonal to contour lines)

$$w_i^{(t)} = w_i^{(t-1)} - \varepsilon^{(t)} \left. \frac{\partial C(\mathbf{w})}{\partial w_i} \right|_{w_i = w_i^{(t-1)}}$$



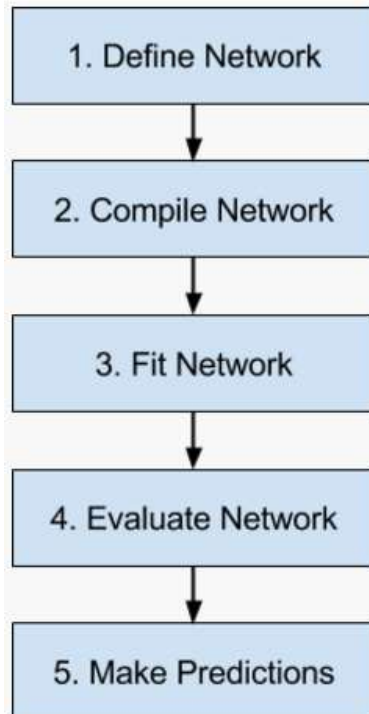
$$- \sum_{k=1}^2 y_k \cdot \log(p_k)$$

Cost C or Loss = cross-entropy averaged over all images in mini-batch

$$C = \frac{1}{N} \sum_i D(\mathbf{p}_i, \mathbf{y}_i)$$

Keras: A high level API with best practice defaults

Keras workflow



- Keras is a high-level NN API for TensorFlow and Theano.
- Is now shipped with TF (or can be imported as python library)
- Can be mixed with TF code
- Offers many predefined layers
- Each layer has a default “best practice choices of parameters”
- Allows for easy and fast prototyping (define only key parameters)
- Supports fcNN, CNNs, RNNs ...
- Supports arbitrary connectivity schemes and NN architectures

Keras: import keras and the required layers

Keras provides two ways to define a model:

- 1) Sequential API: simple, good for linear stack of layers
- 2) Functional API: flexible, required for complicated architectures

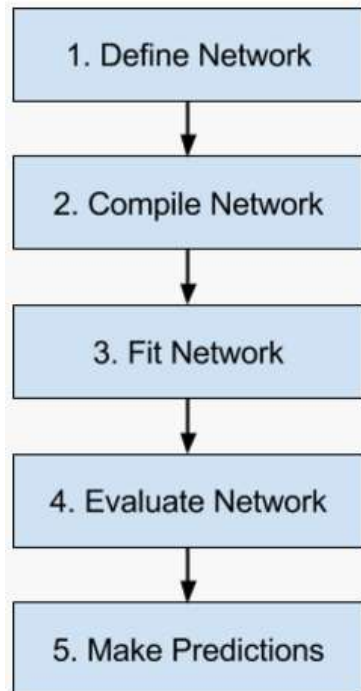
```
import keras
from keras.models import Sequential
from keras.layers import Dense, Activation, Dropout, BatchNormalization
```

For documentation see:

<https://keras.io/>

Remark: Keras can be used as API for theano and tf and which differ in the shape of expected tensors – i.e. #channels is in tf at last position and in theano not.

Keras: A high level API with best practice defaults



```
### Model with default initialization
model = Sequential()

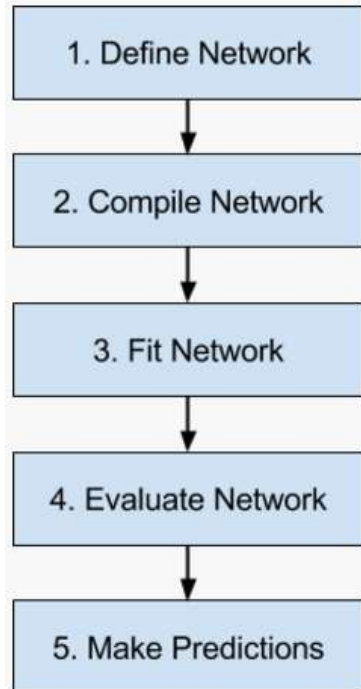
model.add(Dense(500, batch_input_shape=(None, 784)))
model.add(Activation('sigmoid'))

model.add(Dense(50))
model.add(Activation('sigmoid'))

model.add(Dense(10, activation='softmax'))
model.compile(loss='categorical_crossentropy',
              optimizer='adadelta',
              metrics=['accuracy'])
```

Number of neurons in (first)
hidden dense layers
(will be input to next layer)

Keras: A high level API with best practice defaults



```
model.compile(loss='categorical_crossentropy',  
              optimizer='adadelta',  
              metrics=['accuracy'])
```

loss function that will be minimized



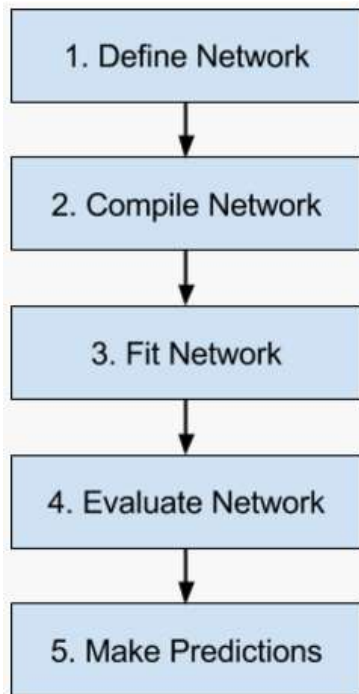
easiest optimizer is SGD
(stochastic gradient descent)



Which metrics besides 'loss' do
we want to collect during training



Keras: A high level API with best practice defaults



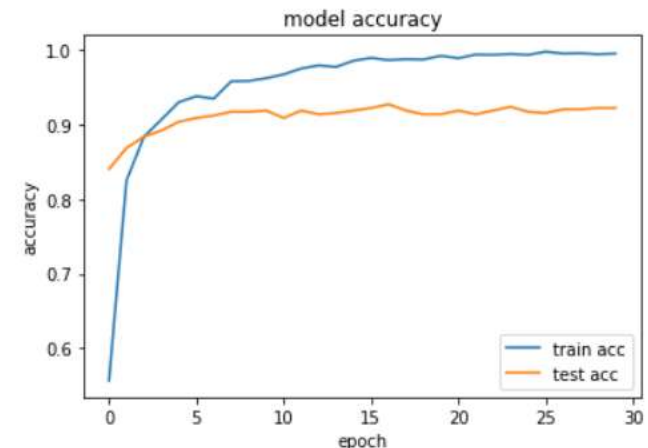
```
# define information required for tensorboard
tensorboard = keras.callbacks.TensorBoard(
    log_dir='tensorboard/mnist_small/' + name + '/',
    write_graph=True,
    histogram_freq=1)

# train the model, memorize training history
history = model.fit(X[0:2400],
                    convertToOneHot(y[0:2400],10),
                    epochs=30,
                    batch_size=128,
                    callbacks=[tensorboard],
                    validation_data=[X[2400:3000],
                                    convertToOneHot(y[2400:3000],10)])
```

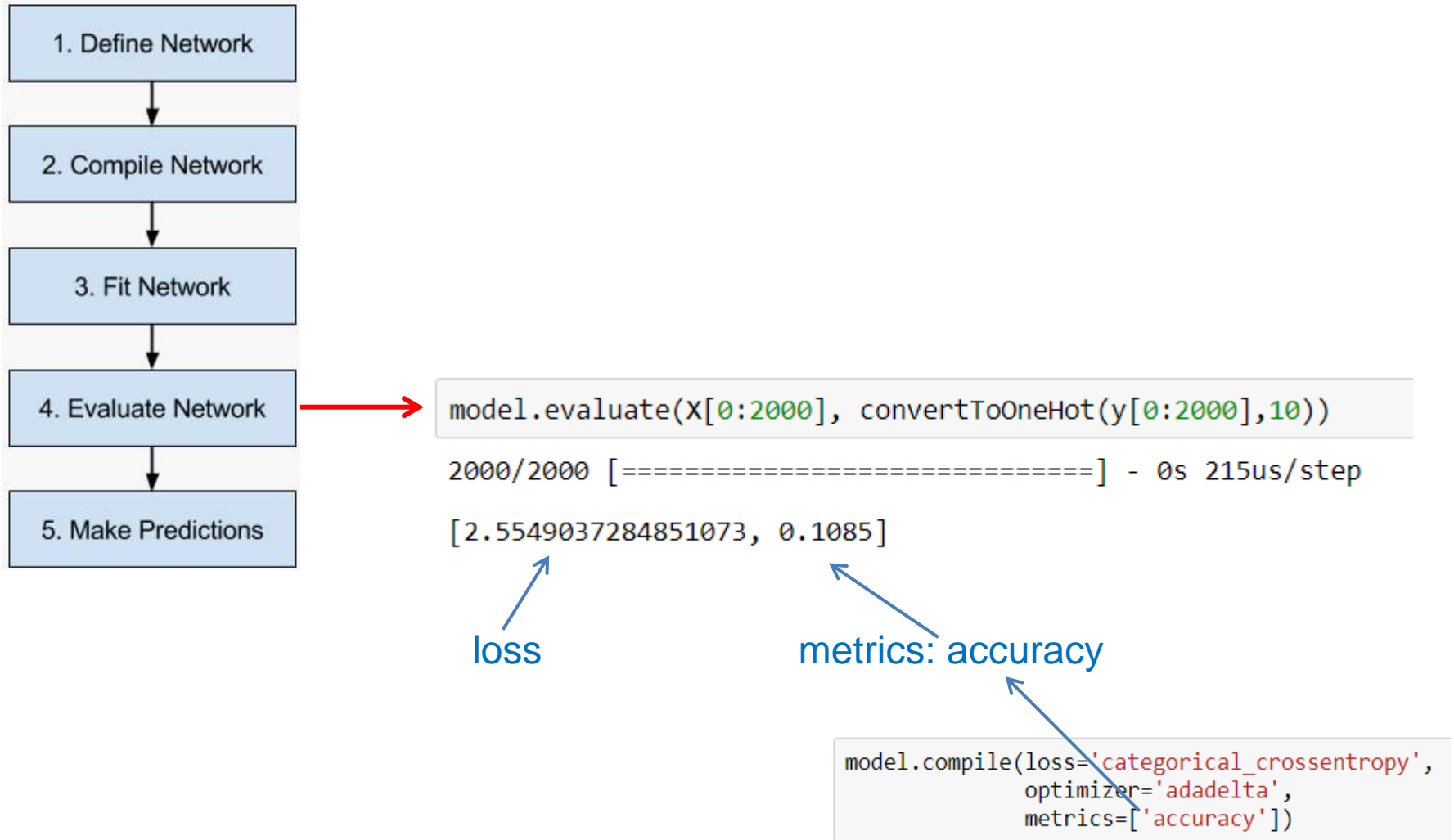
input data (train)
output/label (train)
How and how often provide training data

In history we memorize development of loss and metrics achieved in successive training steps

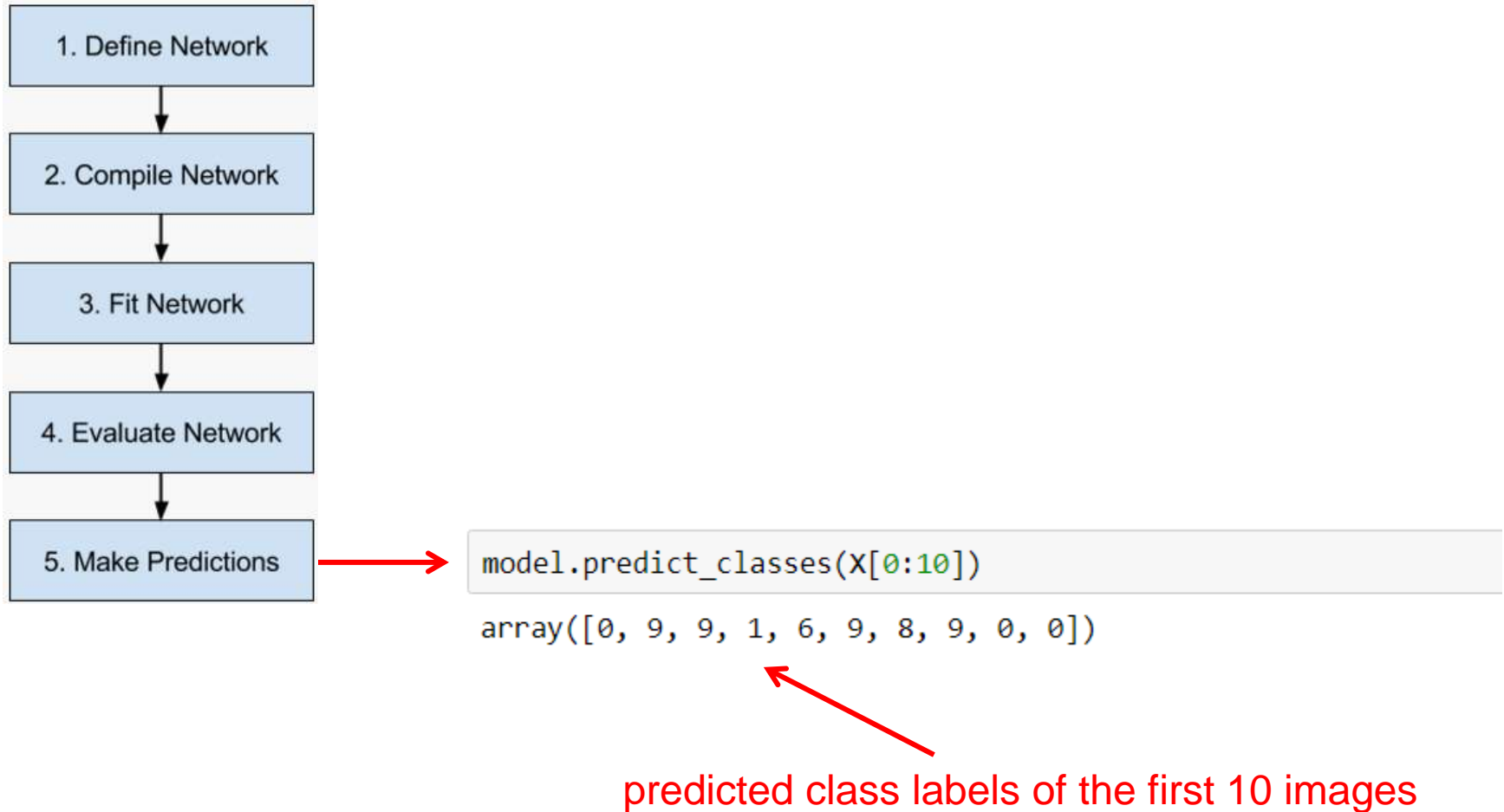
```
plt.plot(history.history['acc'])
plt.plot(history.history['val_acc'])
plt.title('model accuracy')
plt.ylabel('accuracy')
plt.xlabel('epoch')
plt.legend(['train acc', 'test acc'], loc='lower right')
plt.show()
```



Keras: A high level API with best practice defaults



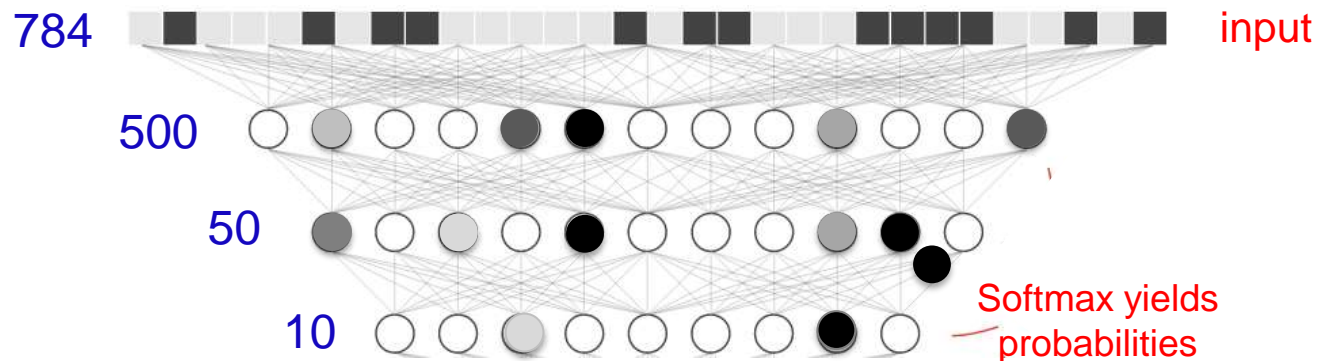
Keras: A high level API with best practice defaults



Keras: gives nice summary of model architecture

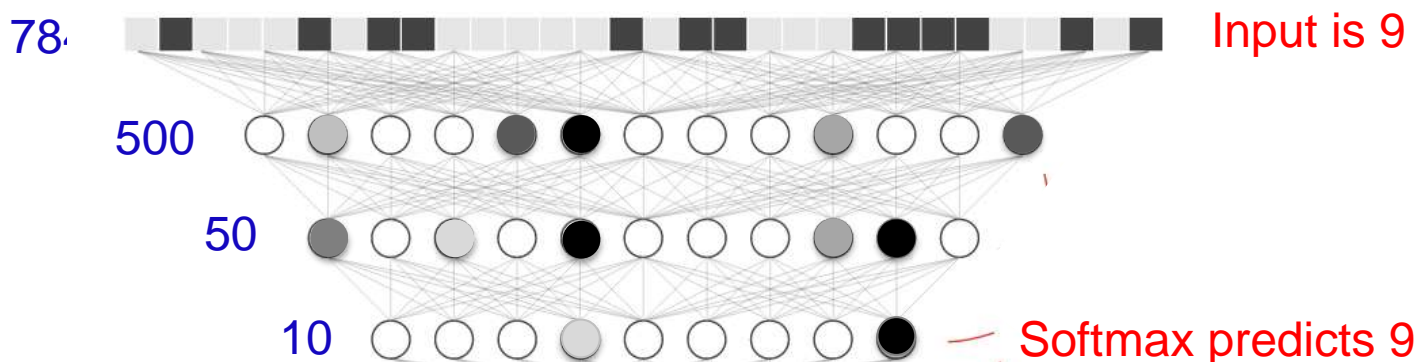
```
model.summary()
```

Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 500)	392500
activation_1 (Activation)	(None, 500)	0
dense_2 (Dense)	(None, 50)	25050
activation_2 (Activation)	(None, 50)	0
dense_3 (Dense)	(None, 10)	510
Total params: 418,060		
Trainable params: 418,060		
Non-trainable params: 0		

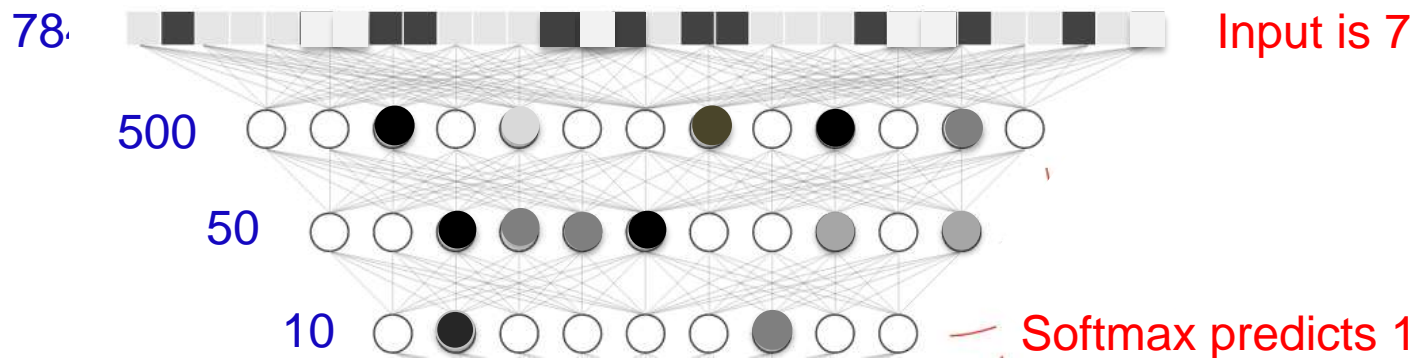


Cheated: we should have 10 (not 9) nodes ;-)

A trained NN has fixed weights and input specific activation patterns in each layer



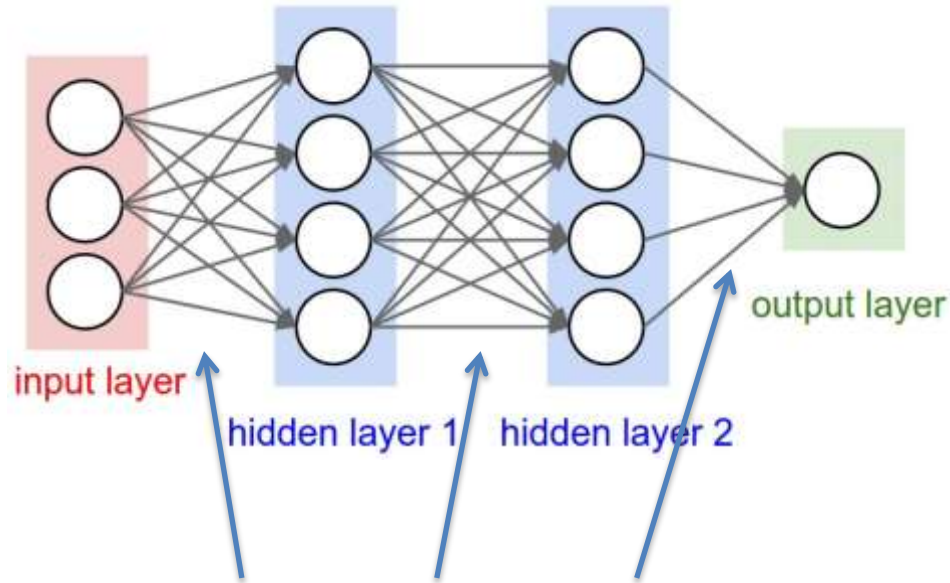
Cheated: we should have 10 (not 9) nodes ;-)



Cheated: we should have 10 (not 9) nodes ;-)

Question in home work:

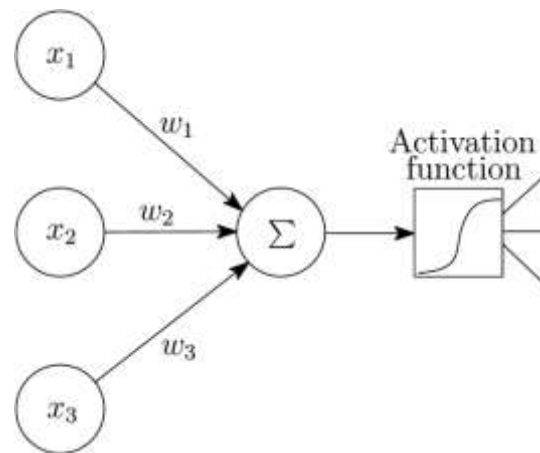
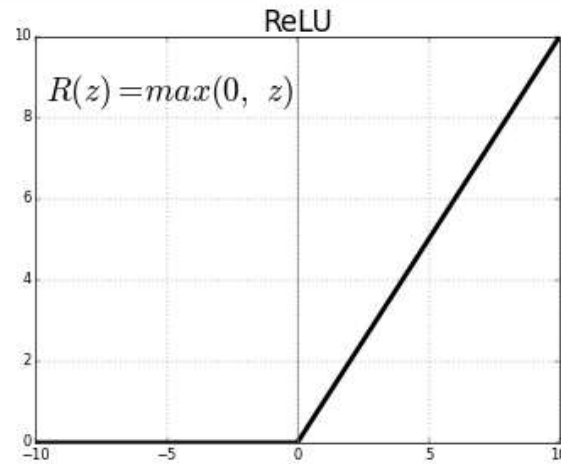
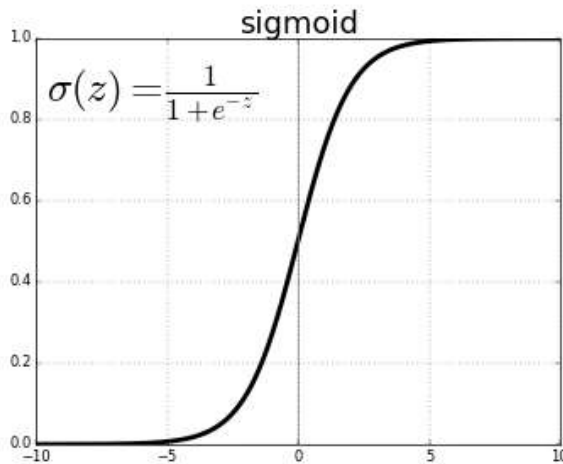
How to initialize the weights? Does it matter?



How to initialize the weights?

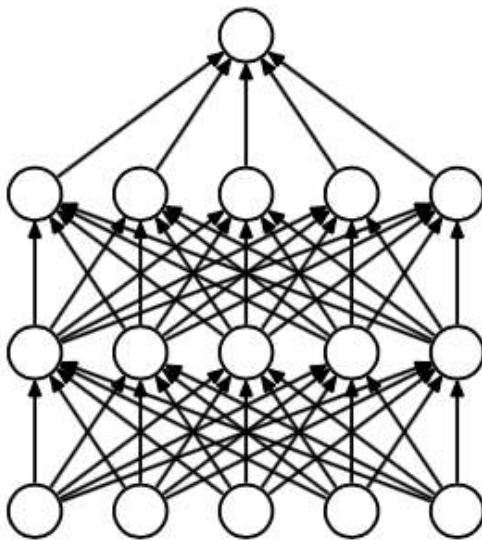
Question in home work:

Which activation function should we use? Does it matter?

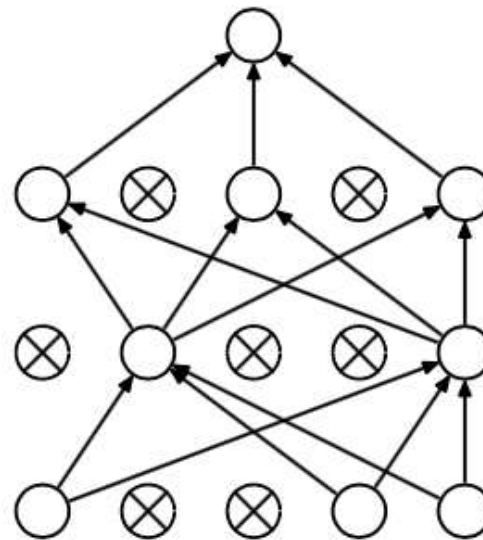


Question in home work: Dropout – does it help?

- After each weight update, we **randomly “delete” a certain percentage of neurons**, which will not be updated in the next step – then repeat.
- In each training step we optimize a slightly different NN model.

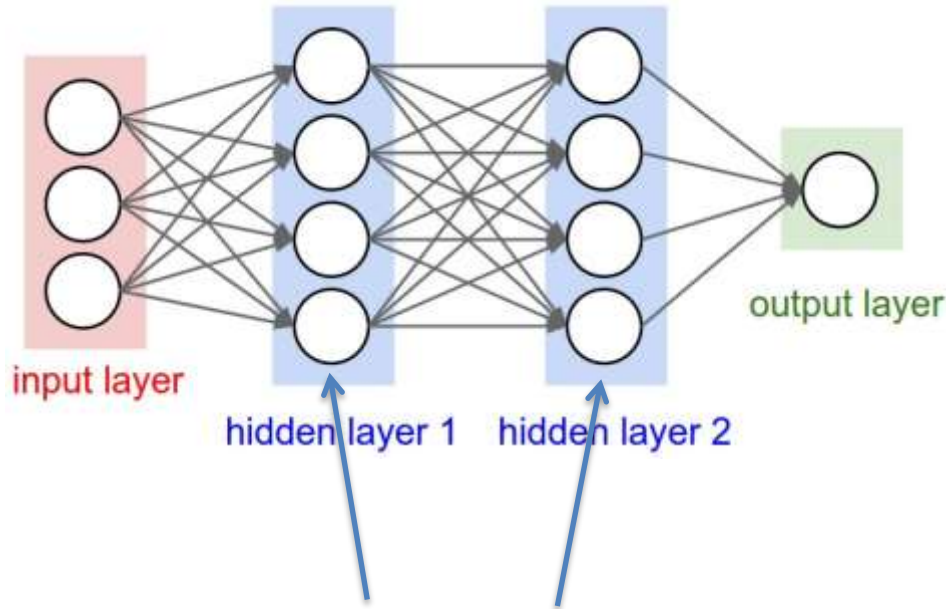


(a) Standard Neural Net

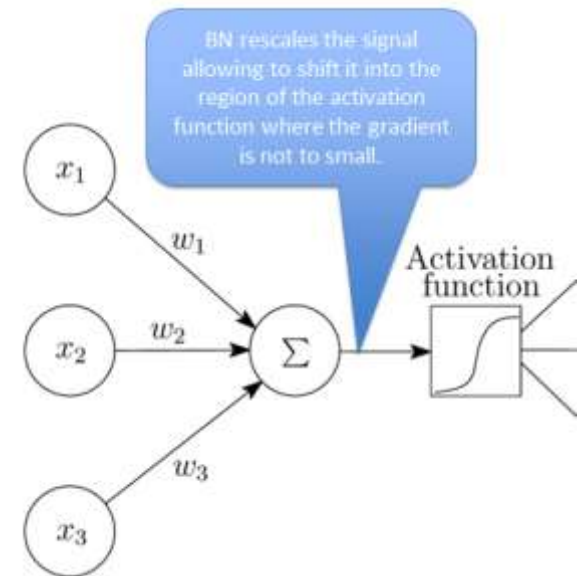


(b) After applying dropout.

Question in home work: Should we allow the NN to normalize the data between layers (batch_norm)? Does it matter?



Should we allow the NN to normalize intermediate data (activations), so that they have mean=0 and sd=1?



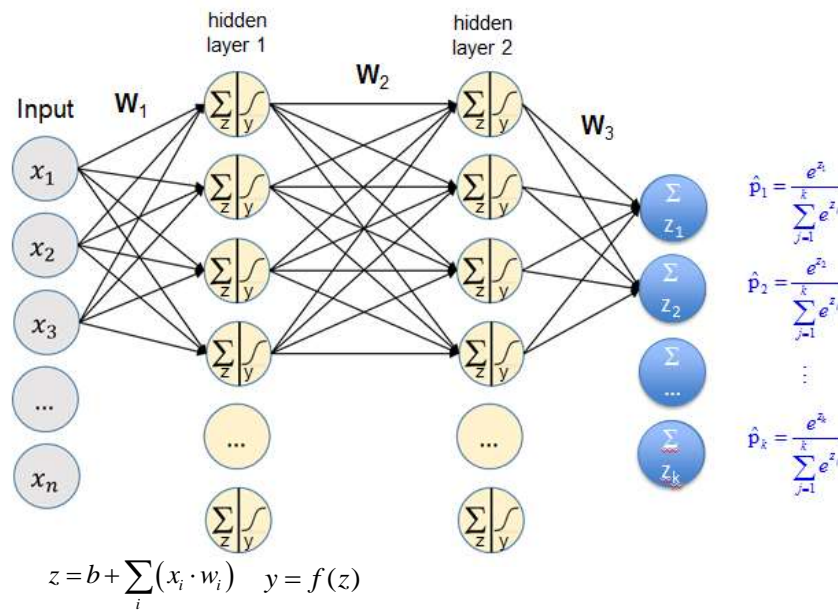
Home work - main result

With this small training set of 4000 images and 100 epochs training we get the **best test accuracy of ~92% when working with random initialization** (reducing weights if number of input data increases), ReLu, dropout and BN (here BN does not improve things – in many applications it does!).



Remark: we will later discuss why BatchNorm and Dropout help.

Bad idea: initializing all weights with the same value



$$w_i^{(t)} = w_i^{(t-1)} - \epsilon^{(t)} \left. \frac{\partial C(\mathbf{w})}{\partial w_i} \right|_{w_i = w_i^{(t-1)}}$$

chain rule:

$$\left. \frac{\partial C}{\partial w_{2n}} \right|_{\text{current NN values}} = \left. \frac{\partial C}{\partial p_k} \right|_{cv} \cdot \left. \frac{\partial \hat{p}_k}{\partial z_{2m}} \right|_{cs} \cdot \left. \frac{\partial z_{2m}}{\partial y_{2m}} \right|_{cv} \cdot \left. \frac{\partial y_{2m}}{\partial z_{1n}} \right|_{cv} \cdot \left. \frac{\partial z_{1n}}{\partial w_{2n}} \right|_{cv}$$

Forward pass: initialize all weights with the same value

⇒ all units get the same values $y_j = f(z_j) = f(b + \sum_i x_i w_i)$

⇒ ... all outputs are the same.. (Initializing all weights=0 will give all units the value 0!)

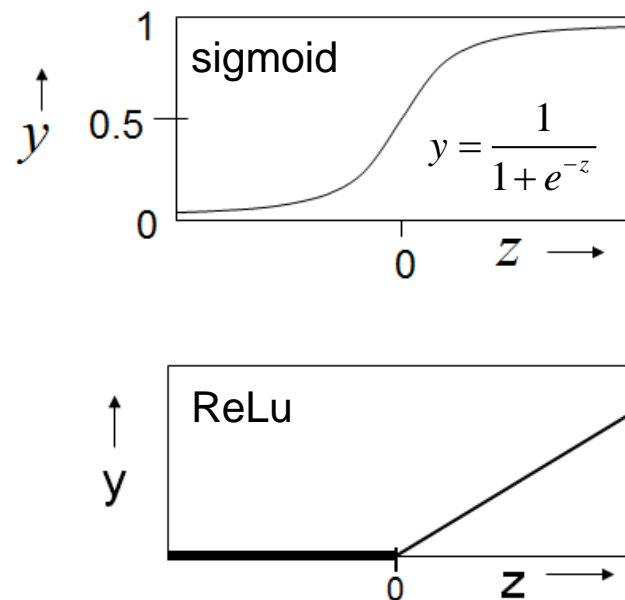
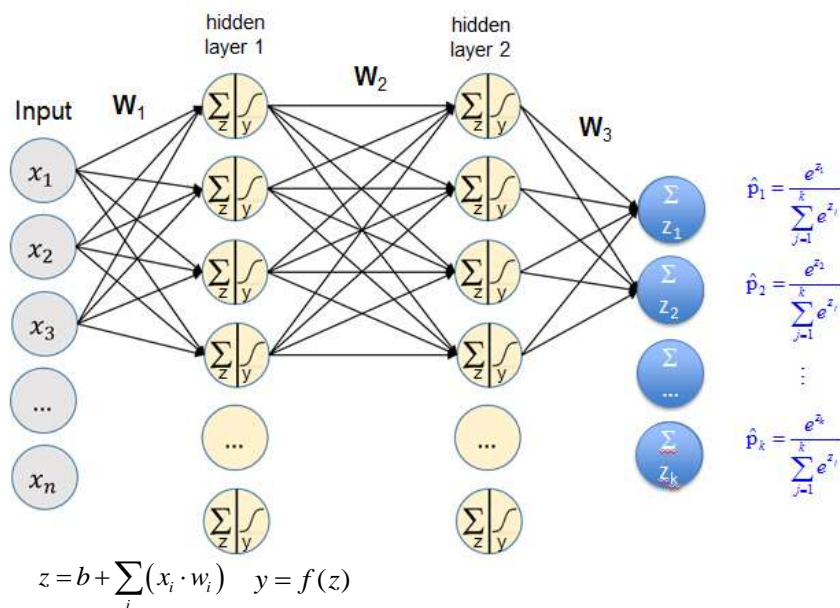
Backward pass: all weights and units have same values & all functions same

⇒ all gradients are the same

⇒ all weights get the same update and get again the same value!

⇒ **no learning**

Bad idea: initializing with high values



Initialize weights with partly large values.

⇒ large absolute z -values

⇒ flat parts of activation function

⇒ According chain rule we multiply with $\frac{\partial y}{\partial z} \approx 0$

⇒ gradient is zero we cannot update the weights ⇒ **no learning**

↓

$$w_i^{(t)} = w_i^{(t-1)} - \epsilon^{(t)} \left. \frac{\partial C(\mathbf{w})}{\partial w_i} \right|_{w_i = w_i^{(t-1)}}$$

What is the default initializer in Keras?

Dense <https://keras.io/layers/core/>

Also in conv layer 'glorot_uniform' is used as default initializer.

[source]

```
keras.layers.Dense(units, activation=None, use_bias=True, kernel_initializer='glorot_uniform', bias_i
```

glorot_uniform

Recommended 2010 by Glorot & Bengio

<http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>

```
glorot_uniform(seed=None)
```

Glorot uniform initializer, also called Xavier uniform initializer.

guarantees random small numbers

It draws samples from a uniform distribution within $[-limit, limit]$ where `limit` is

`$\sqrt{6 / (fan_in + fan_out)}$` where `fan_in` is the number of input units in the weight tensor and `fan_out` is the number of output units in the weight tensor.

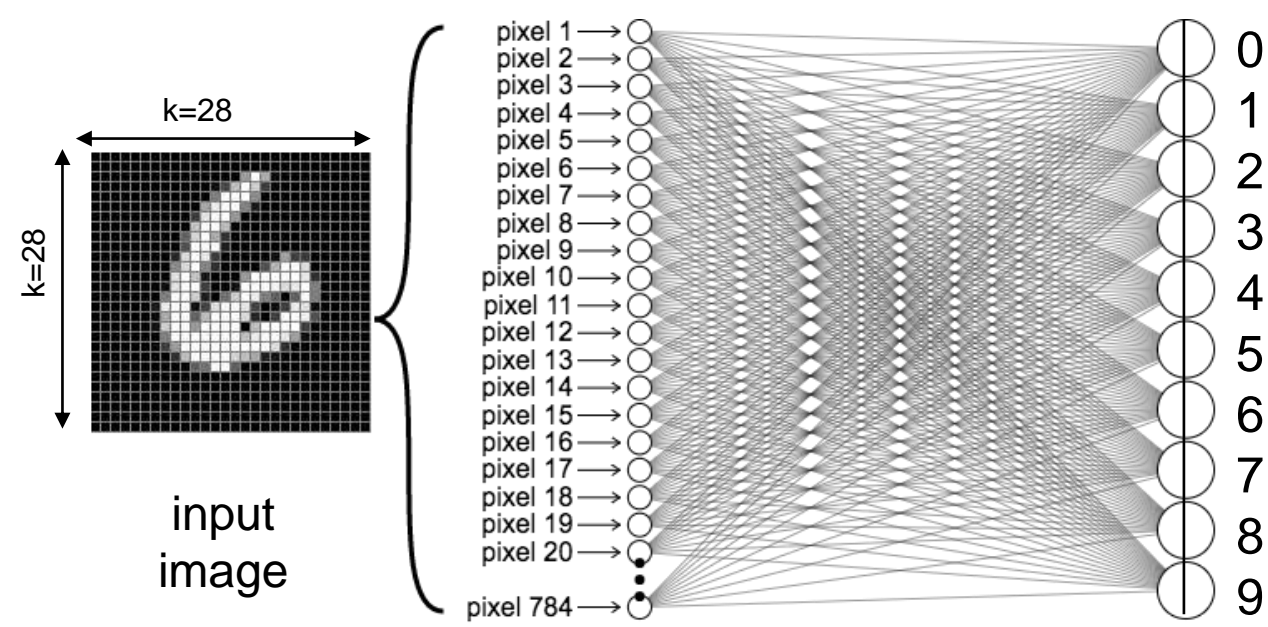
he_normal Great initializer for CNN with ReLu – see He's 2015 paper <https://arxiv.org/abs/1502.01852>

```
he_normal(seed=None)
```

He normal initializer.

It draws samples from a truncated normal distribution centered on 0 with `$stddev = \sqrt{2 / fan_in}$` where `fan_in` is the number of input units in the weight tensor.

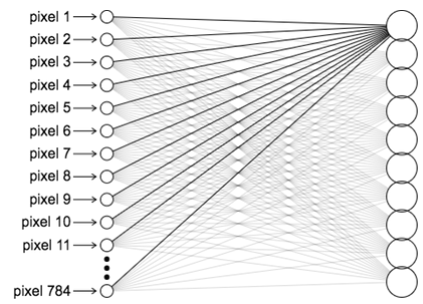
MNIST classification with a 1 layer fully connected NN



$$z_j = \sum_i (x_i \cdot w_{ij})$$

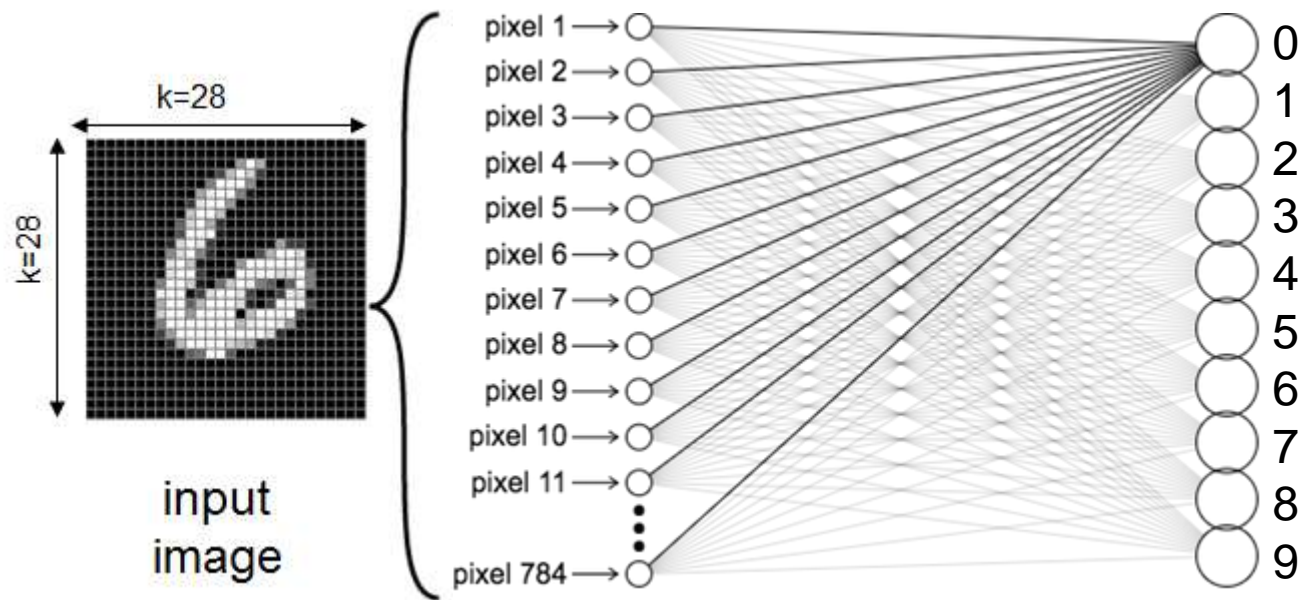
$$p_j = \frac{e^{z_j}}{\sum_{k=0}^9 e^{z_k}}$$

Each output node z_j gets 784 inputs connected via 784 weights which (usually the weights are different for each node) can also be arranged as **28x28 weight image**



W_1	W_2	W_3	\dots	W_{28}
W_{29}	W_{30}	W_{31}		W_{56}
W_{57}	W_{58}	W_{59}		W_{84}
\vdots			\ddots	
W_{757}	W_{758}	W_{759}		W_{784}

MNIST classification with a 1 layer fc NN

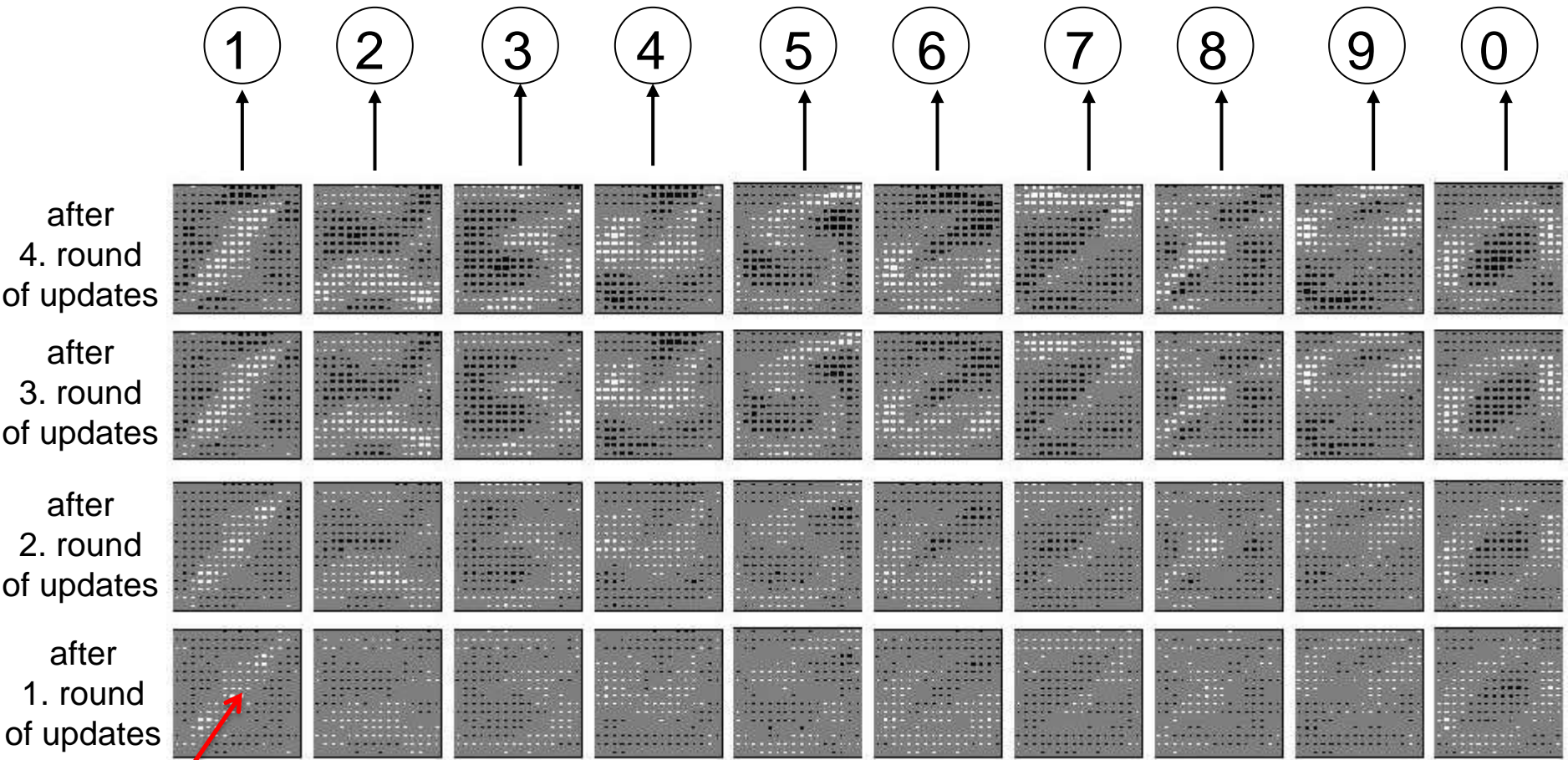


$$p_j = \frac{e^{z_j}}{\sum_{k=0}^9 e^{z_k}}$$

$$z_j = \sum_i (x_i \cdot w_{ij})$$

- For the correct digit we want to get a high probability p_j
 \Rightarrow we need a large z_j
- z_j is large if at large signal pixels x the w are also large
and at small signal pixels x the w are also small
- For maximal p_j the j^{th} weight image looks like the input image!

Displaying learned weights at different steps of training via backpropagation

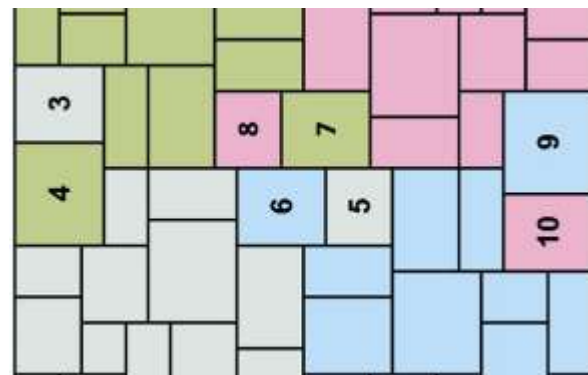
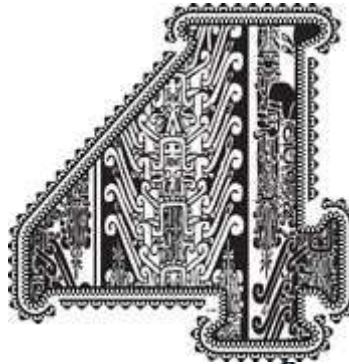


A 1 layer fully connected NN learns rigid template for each class!

“weight images”

What is wrong with a rigid template? Find the 4!

1 I II III IV V VI VII VIII IX X
2 3 4 5 6 7 8 9 10
XI XII
LXXX XC XCVIII XCIX C D M
LABB





Some tasks are hard - how can solve them?

The training data consists of many different pictures of Oliver Dür and Albert Einstein

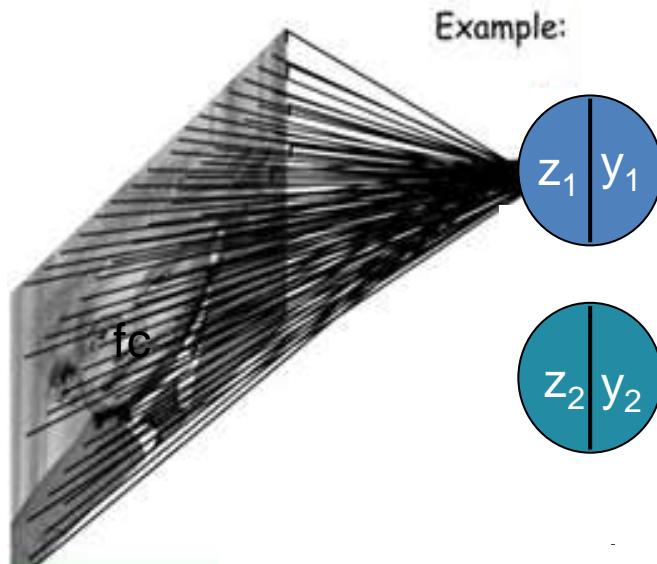


How to distinguish Oliver from Einstein?

A hard task !



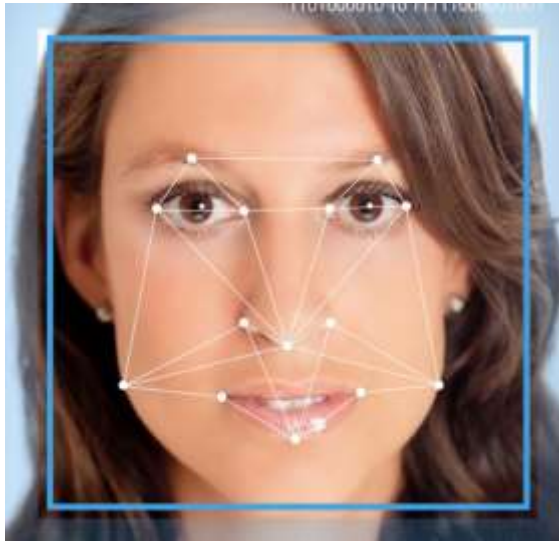
At least with only 2 neurons...



A rigid template is not enough!

We should go beyond fc NN!

How to do computer vision? Traditional and DL approach

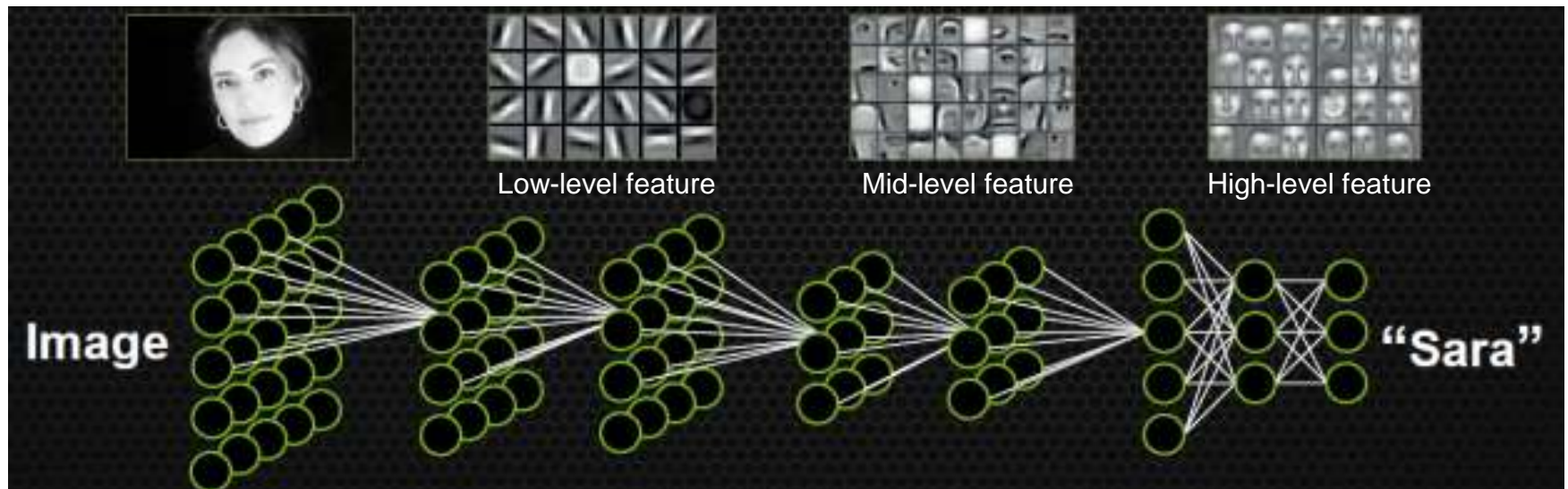


Traditional:

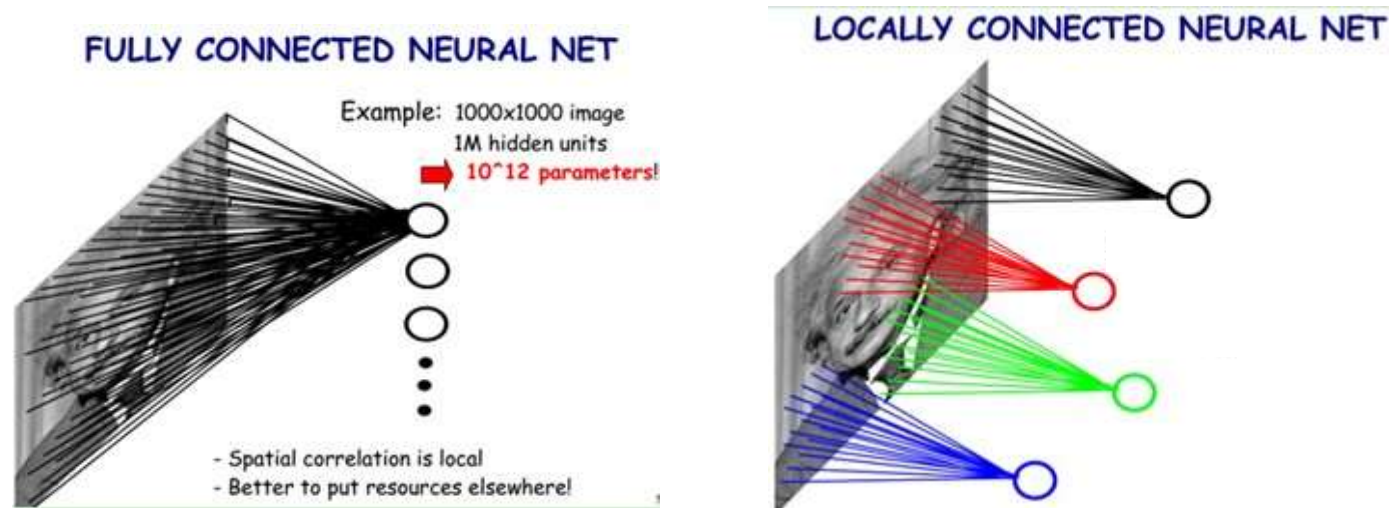
Extract **handcrafted features** & use these features to **train / fit a model** (e.g. logistic regression) and use fitted model to perform classification/prediction.

Deep learning using CNNs:

Based on variants of **deep** neural networks which **learn** during training/fitting appropriate **hierarchical features** e.g. from images. The fitted model (artificial neural net) can be used for classification/prediction.



Convolution extracts local information using few weights



Shared weights:

by using the **same weights for each patch** of the image we need much **less parameters** than in the fully connected NN and get from each patch the same kind of **local feature information** such as the presence of a edge.

Convolutional networks use neighborhood information and replicated local feature extraction

In a locally connected network the calculation rule

$$z = b + \sum_i (x_i \cdot w_i)$$

Corresponds to convolution of a filter with the image and the pattern of weights represent a filter.

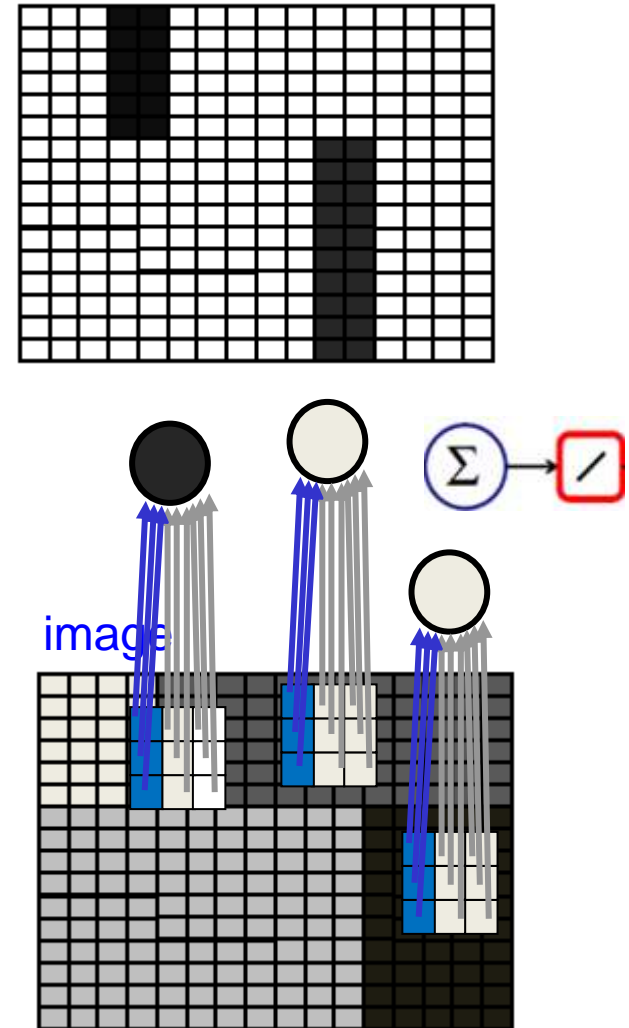
w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9

0.9	0.1	0.1
0.9	0.1	0.1
0.9	0.1	0.1

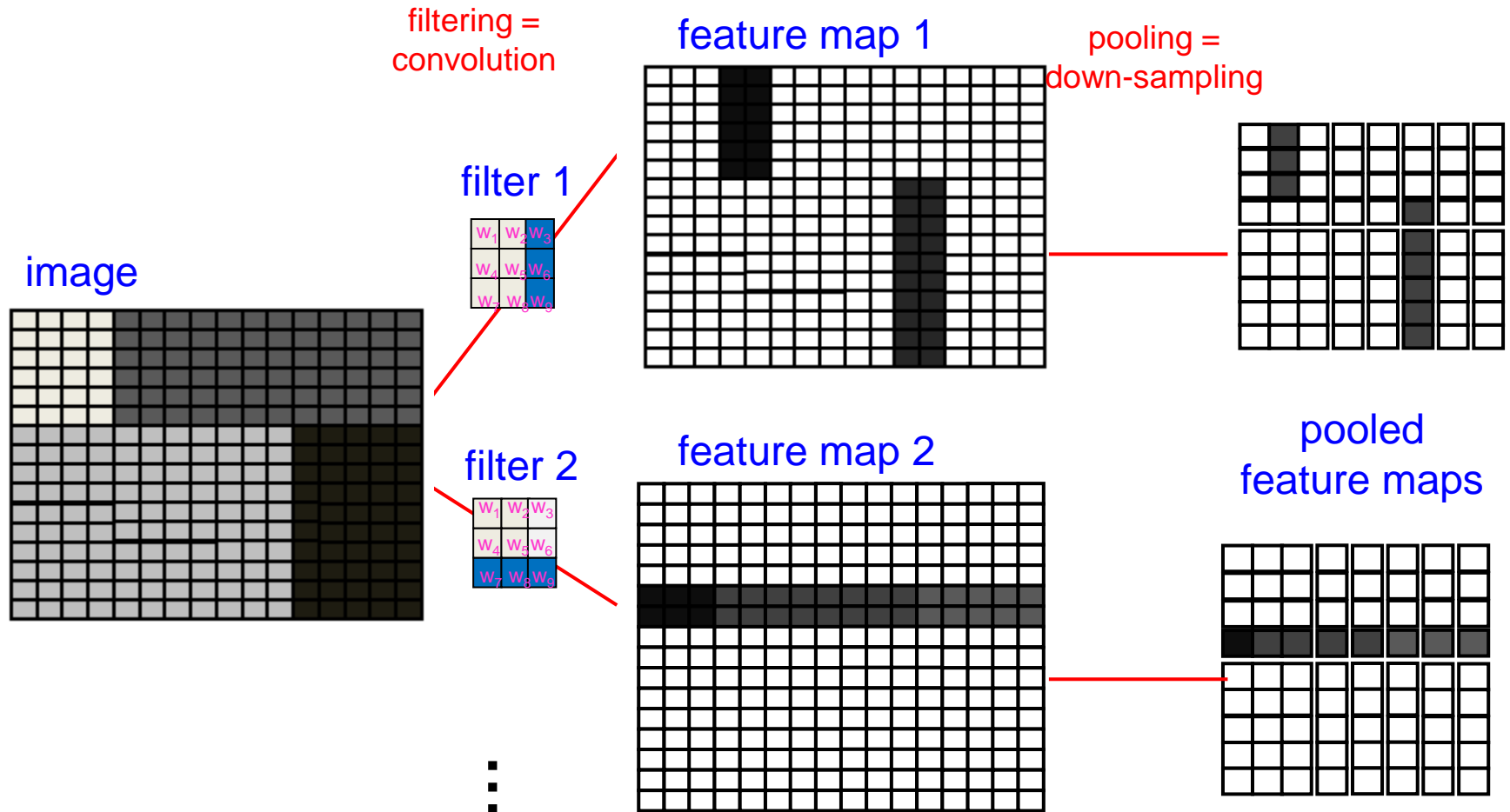
The filter is applied at each position of the image and it can be shown that the **result is maximal** if the **image pattern corresponds to the weight pattern**.

The results form again an image called **feature map** (=activation map) which shows at which position the feature is present.

feature/activation map



Convolutional networks use neighborhood information and replicated local feature extraction



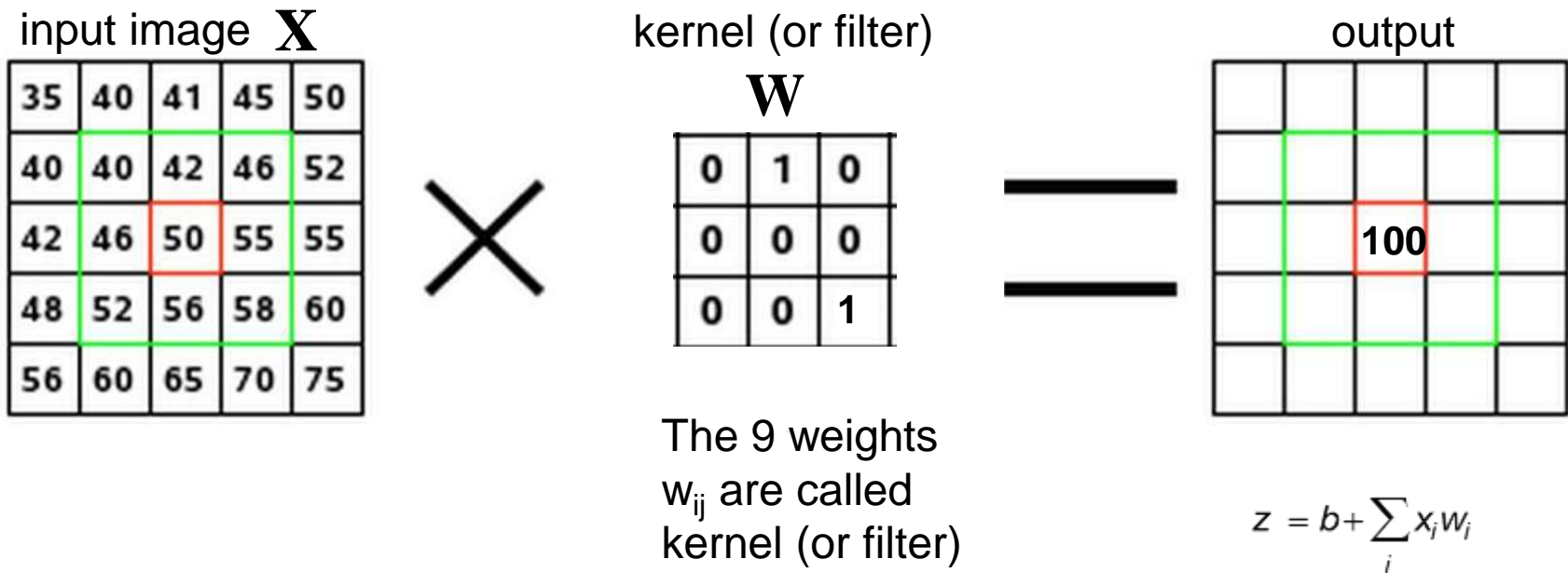
The weights of each filter are randomly initiated and then adapted during the training.

Convolutional networks

Building Block 1: Convolution

What is convolution?

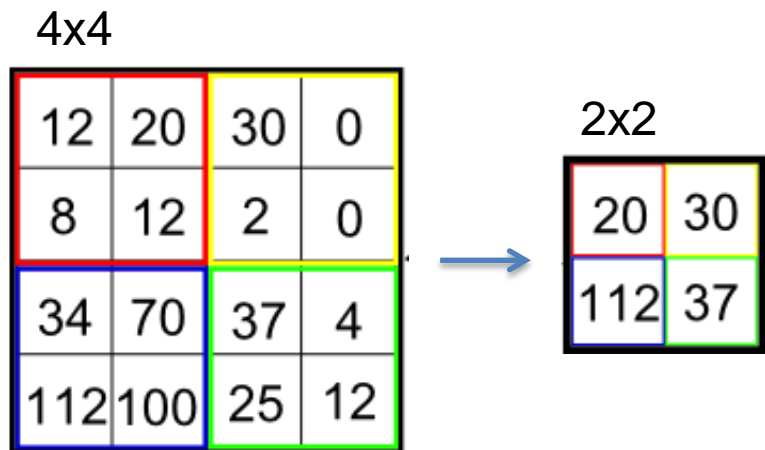
At each position of the kernel we perform an elementwise dot-product between intensity values of the underlying input image



In deep-learning the weights are not fixed, they are learned!

Convolutional networks

Building Block 2: (Max-) pooling



Simply join e.g. 2x2 adjacent pixels in one, e.g. by picking the maximum of all 4 pixel values (max-pooling) or determining the average of the 4 pixel values (average pooling)

Hinton: „The pooling operation used in convolutional neural networks is a big mistake and the fact that it works so well is a disaster“

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	245	250	230	120
170	180	235	145	170	255
190	185	170	165	130	120
255	255	245	190	200	175

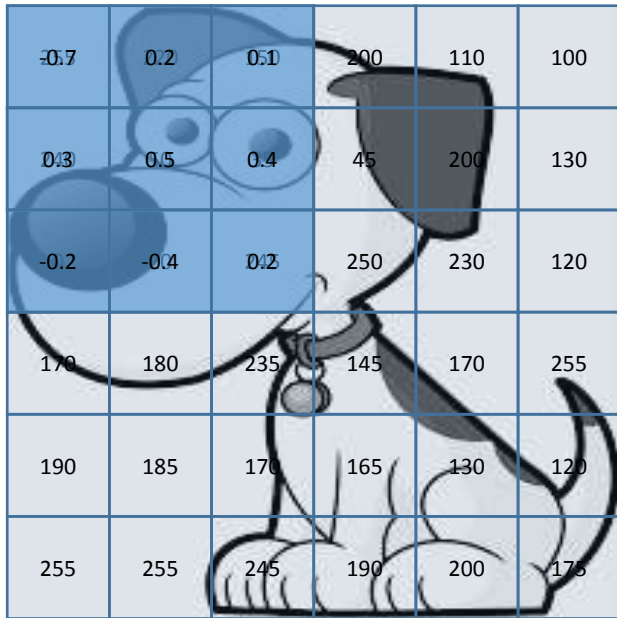
-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

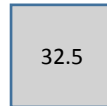
$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1



Feature map
4x4x1



3x3 filter

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	-0.7	0.2	0.1	110	100
240	0.3	0.5	0.4	200	130
0	-0.2	-0.4	0.2	230	120
170	180	235	145	170	255
190	185	170	165	130	120
255	255	245	190	200	175

Feature map
4x4x1

32.5	-105.5
------	--------

3x3 filter

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	-0.7	0.2	0.1	100
240	50	0.3	0.5	0.4	130
0	20	-0.2	-0.4	0.2	120
170	180	235	145	170	255
190	185	170	165	130	120
255	255	245	190	200	175

Feature map
4x4x1

32.5	-105.5	185.5
------	--------	-------

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	-0.7	0.2	0.1
240	50	35	0.3	0.5	0.4
0	20	245	-0.2	-0.4	0.2
170	180	235	145	170	255
190	185	170	165	130	120
255	255	245	190	200	175

Feature map
4x4x1

32.5	-105.5	185.5	54
------	--------	-------	----

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
-0.7	0.2	0.1	45	200	130
0.3	0.5	0.4	250	230	120
-0.2	-0.4	0.2	145	170	255
190	185	170	165	130	120
255	255	245	190	200	175

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5			

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	-0.7	0.2	0.1	200	130
0	0.3	0.5	0.4	230	120
170	-0.2	-0.4	0.2	170	255
190	185	170	165	130	120
255	255	245	190	200	175

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104		

3x3 filter

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	-0.7	0.2	0.1	130
0	20	0.3	0.5	0.4	120
170	180	-0.2	-0.4	0.2	255
190	185	170	165	130	120
255	255	245	190	200	175

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	-0.7	0.2	0.1
0	20	245	0.3	0.5	0.4
170	180	235	-0.2	-0.4	0.2
190	185	170	165	130	120
255	255	245	190	200	175

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
-0.7	0.2	0.1	250	230	120
0.3	0.5	0.4	145	170	255
-0.2	-0.4	0.2	165	130	120
255	255	245	190	200	175

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44			

3x3 filter

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	-0.7	0.2	0.1	230	120
170	0.3	0.5	0.4	170	255
190	-0.2	-0.4	0.2	130	120
255	255	245	190	200	175

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224		

3x3 filter

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	-0.7	0.2	0.1	120
170	180	0.3	0.5	0.4	255
190	185	-0.2	-0.4	0.2	120
255	255	245	190	200	175

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224	38.5	

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	245	-0.7	0.2	0.1
170	180	235	0.3	0.5	0.4
190	185	170	-0.2	-0.4	0.2
255	255	245	190	200	175

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224	38.5	-18

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	245	250	230	120
-0.7	0.2	0.1	145	170	255
0.3	0.5	0.4	165	130	120
-0.2	-0.4	0.2	190	200	175

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224	38.5	-18
-60.5			

3x3 filter

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	245	250	230	120
170	-0.7	0.2	0.1	170	255
190	0.3	0.5	0.4	130	120
255	-0.2	-0.4	0.2	200	175

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224	38.5	-18
-60.5	213.5		

3x3 filter

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	245	250	230	120
170	180	-0.7	0.2	0.1	255
190	185	0.3	0.5	0.4	120
255	255	-0.2	-0.4	0.2	175

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224	38.5	-18
-60.5	213.5	52.5	

3x3 filter

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

$$z_j = \sum_i (x_i \cdot w_{ij})$$

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	245	250	230	120
170	180	235	-0.7	0.2	0.1
190	185	170	0.3	0.5	0.4
255	255	245	-0.2	-0.4	0.2

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224	38.5	-18
-60.5	213.5	52.5	37.5

3x3 filter

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

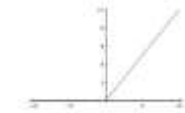
Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	245	250	230	120
170	180	235	-0.7	0.2	0.1
190	185	170	0.3	0.5	0.4
255	255	245	-0.2	-0.4	0.2

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224	38.5	-18
-60.5	213.5	52.5	37.5



Relu

32.5	0	185.5	54
0	104	217.5	31
0	224	38.5	0
0	213.5	52.5	37.5

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	245	250	230	120
170	180	235	-0.7	0.2	0.1
190	185	170	0.3	0.5	0.4
255	255	245	-0.2	-0.4	0.2

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224	38.5	-18
-60.5	213.5	52.5	37.5

Relu

32.5	0	185.5	54
0	104	217.5	31
0	224	38.5	0
0	213.5	52.5	37.5

Maxpool
(2x2x1)

104

3x3 filter

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

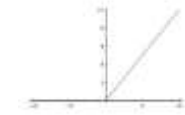
Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	245	250	230	120
170	180	235	-0.7	0.2	0.1
190	185	170	0.3	0.5	0.4
255	255	245	-0.2	-0.4	0.2

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224	38.5	-18
-60.5	213.5	52.5	37.5



Relu

32.5	0	185.5	54
0	104	217.5	31
0	224	38.5	0
0	213.5	52.5	37.5

Maxpool
(2x2x1)

104	217.5
-----	-------

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

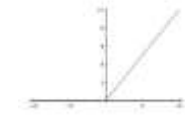
Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	245	250	230	120
170	180	235	-0.7	0.2	0.1
190	185	170	0.3	0.5	0.4
255	255	245	-0.2	-0.4	0.2

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224	38.5	-18
-60.5	213.5	52.5	37.5



Relu

32.5	0	185.5	54
0	104	217.5	31
0	224	38.5	0
0	213.5	52.5	37.5

Maxpool
(2x2x1)

104	217.5
224	

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

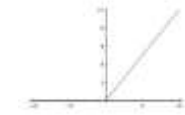
Convolutional neural networks

Input image 6x6x1

255	220	150	200	110	100
240	50	35	45	200	130
0	20	245	250	230	120
170	180	235	-0.7	0.2	0.1
190	185	170	0.3	0.5	0.4
255	255	245	-0.2	-0.4	0.2

Feature map
4x4x1

32.5	-105.5	185.5	54
-105.5	104	217.5	31
-44	224	38.5	-18
-60.5	213.5	52.5	37.5



Relu

32.5	0	185.5	54
0	104	217.5	31
0	224	38.5	0
0	213.5	52.5	37.5

Maxpool
(2x2x1)

104	217.5
224	52.5

-0.7	0.2	0.1
0.3	0.5	0.4
-0.2	-0.4	0.2

3x3 filter

One kernel or filter searches for specific local feature



image patch

0	0	0	0	0	0	30
0	0	0	0	50	50	50
0	0	0	20	50	0	0
0	0	0	50	50	0	0
0	0	0	50	50	0	0
0	0	0	50	50	0	0
0	0	0	50	50	0	0

Pixel representation of the receptive field

filter/kernel: curve detector

0	0	0	0	0	30	0
0	0	0	0	30	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0

*

=6600

Pixel representation of filter

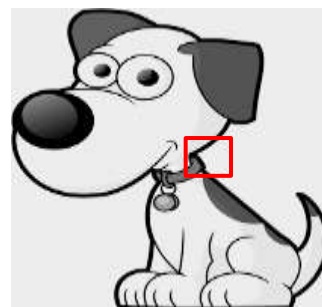


image patch

0	0	0	0	0	0	0
0	40	0	0	0	0	0
40	0	40	0	0	0	0
40	20	0	0	0	0	0
0	50	0	0	0	0	0
0	0	50	0	0	0	0
25	25	0	50	0	0	0

Pixel representation of receptive field

filter/kernel: curve detector

0	0	0	0	0	30	0
0	0	0	0	30	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0

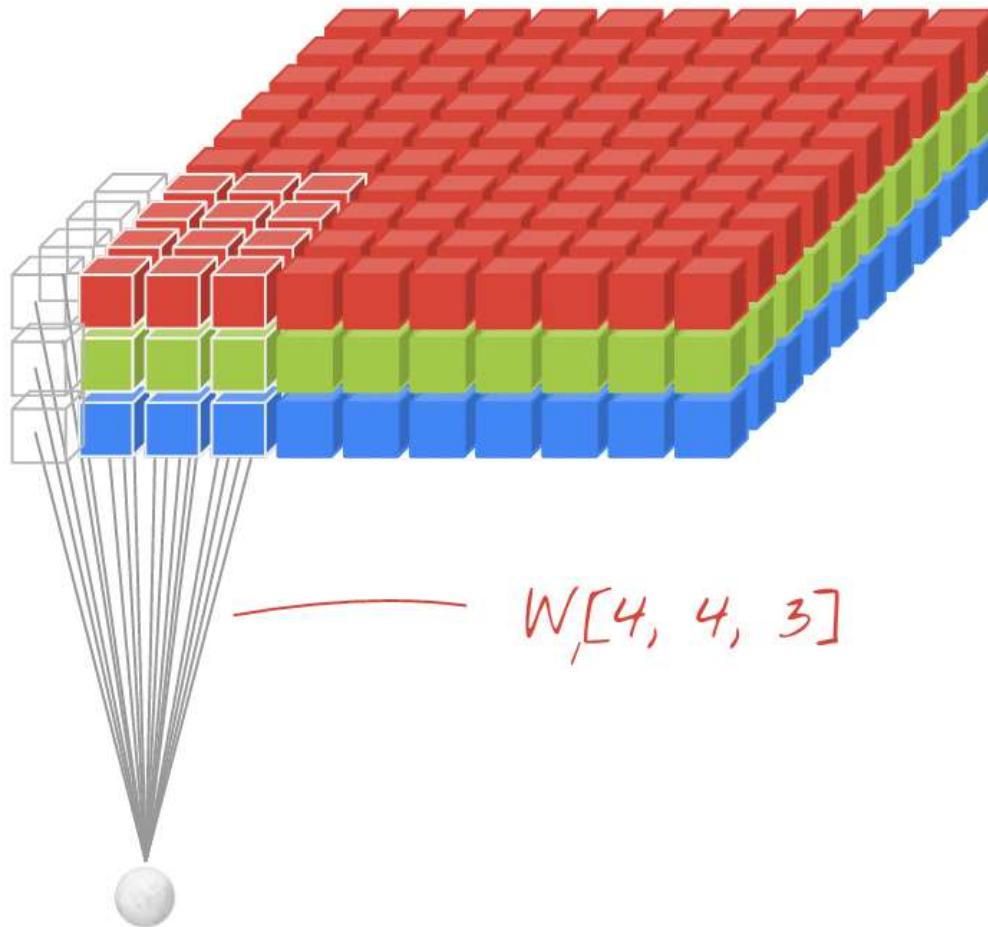
*

=0

Pixel representation of filter

We get a large resulting value if the filter resembles the pattern in the image patch on which the filter was applied.

Animated convolution with 3 input channels



Convolution with 3 input channels and k output channels

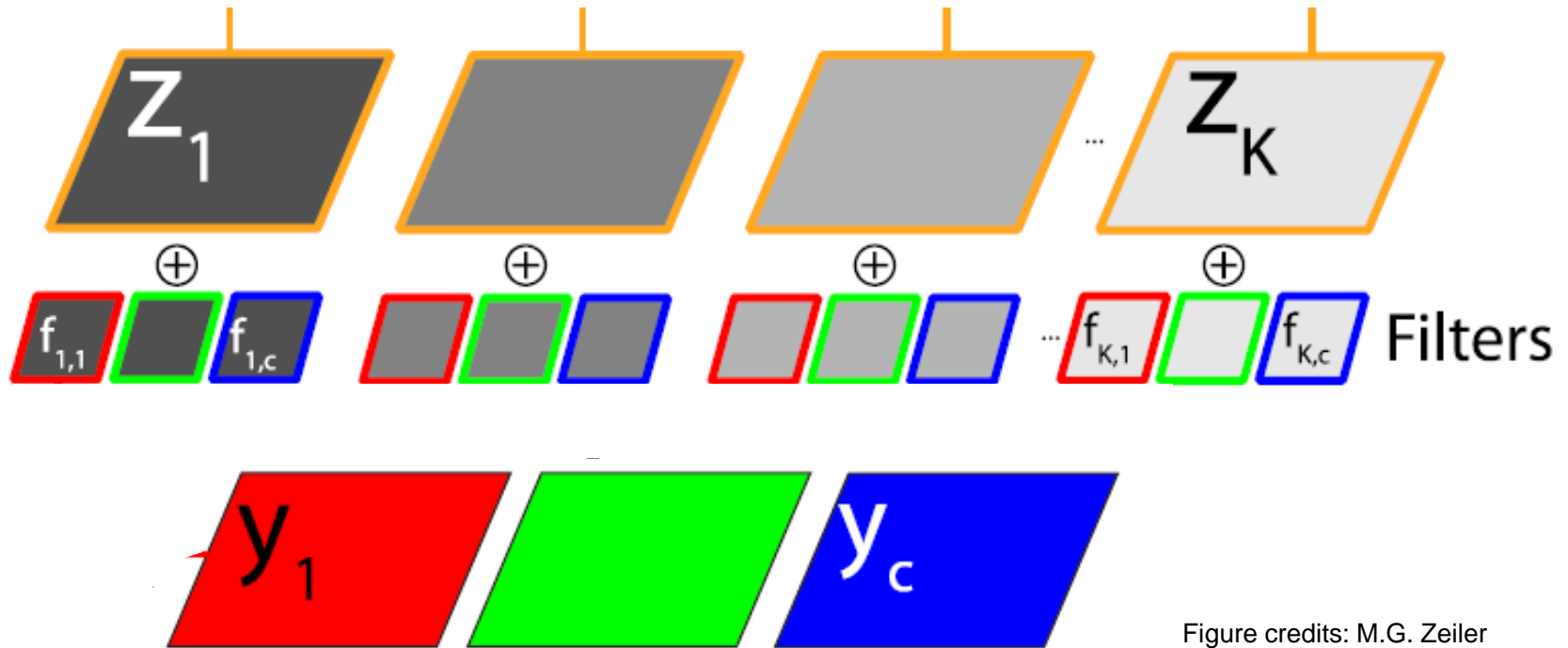


Figure credits: M.G. Zeiler

$$z_k = b + \sum_j \sum_i (y_{ji} \cdot {}^k w_{ji}) = b + \left(\sum_i (y_{1i} \cdot {}^k w_{1i}) + \sum_i (y_{2i} \cdot {}^k w_{2i}) + \sum_i (y_{3i} \cdot {}^k w_{3i}) \right)$$

We can imagine to have for each input channel y_c a channel specific filter f_c and get one output activation map per channel and then we add them element-wise.

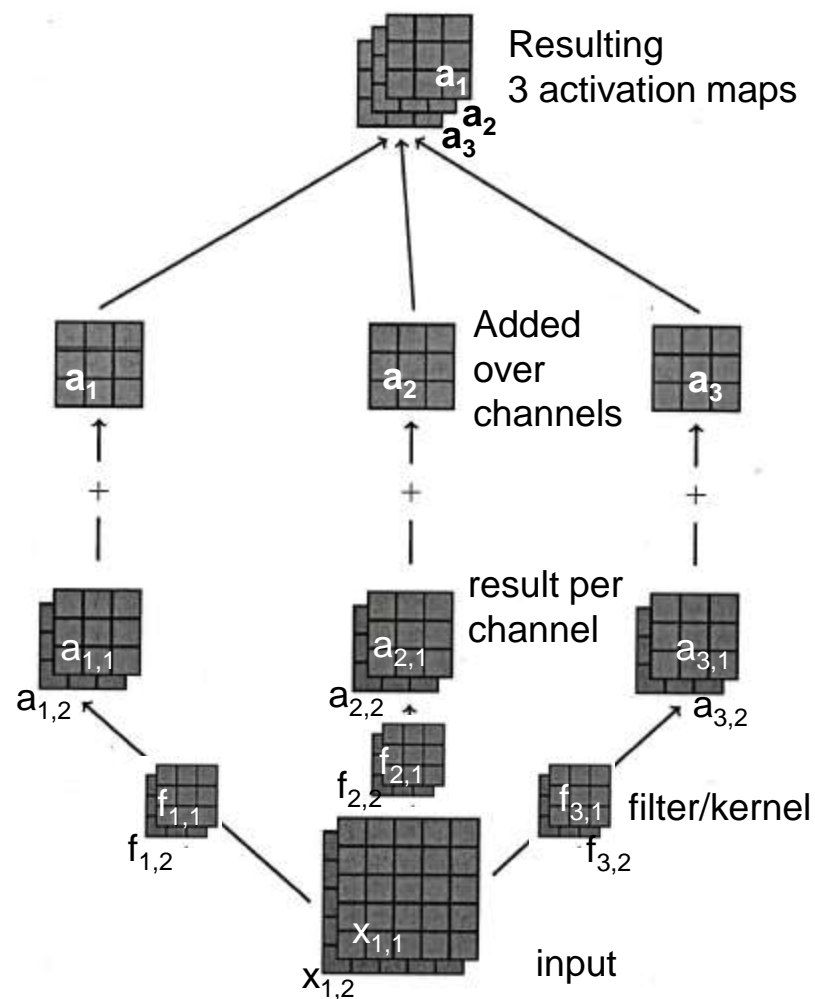
Convolution with 2 input channels and 3 output channels

Here, we need 3-times 2 filters.

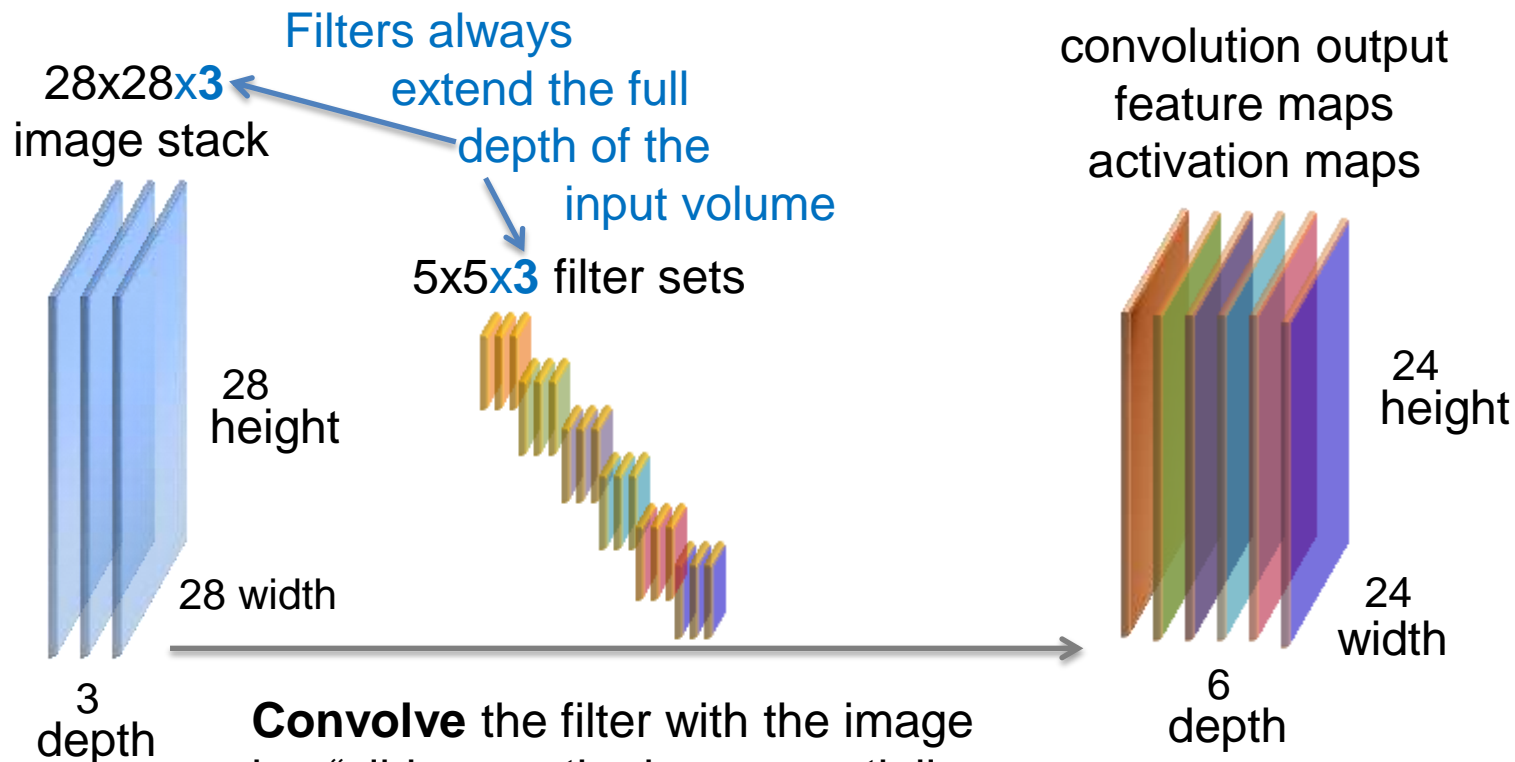
To get the **first output channel** we convolve the **first filter** $f_{1,1}$ of the **first 2-filter-block** with the **first input channel** $x_{1,1}$ and the **second filter** $f_{1,2}$ of the **first 2-filter-block** with the **second input channel** $x_{1,2}$. Then we sum the channel-wise resulting activation maps $a_{1,1}$ and $a_{1,2}$ at each position to get the first output channel a_1 .

To get the **second output channel** we convolve $f_{2,1}$ with the $x_{1,1}$ and $f_{2,2}$ with the $x_{1,2}$. Elementwise adding of the results yields a_2 .

To get the **third output channel** we convolve $f_{3,1}$ with the $x_{1,1}$ and $f_{3,2}$ with the $x_{1,2}$. Elementwise adding of the results yields a_3 .



Convolution continued

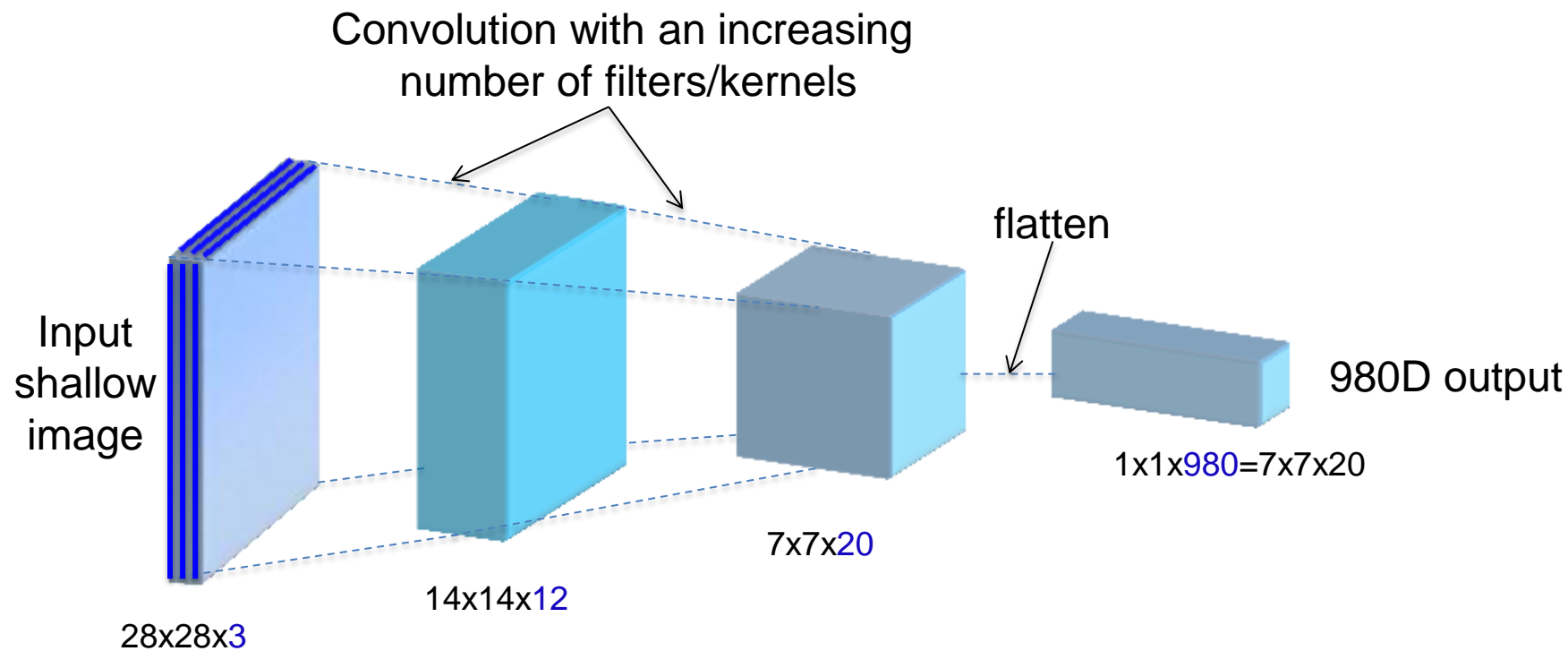


Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products at each position”

Output size in case of
6 5x5x3 filter
no zero-padding
stride=1

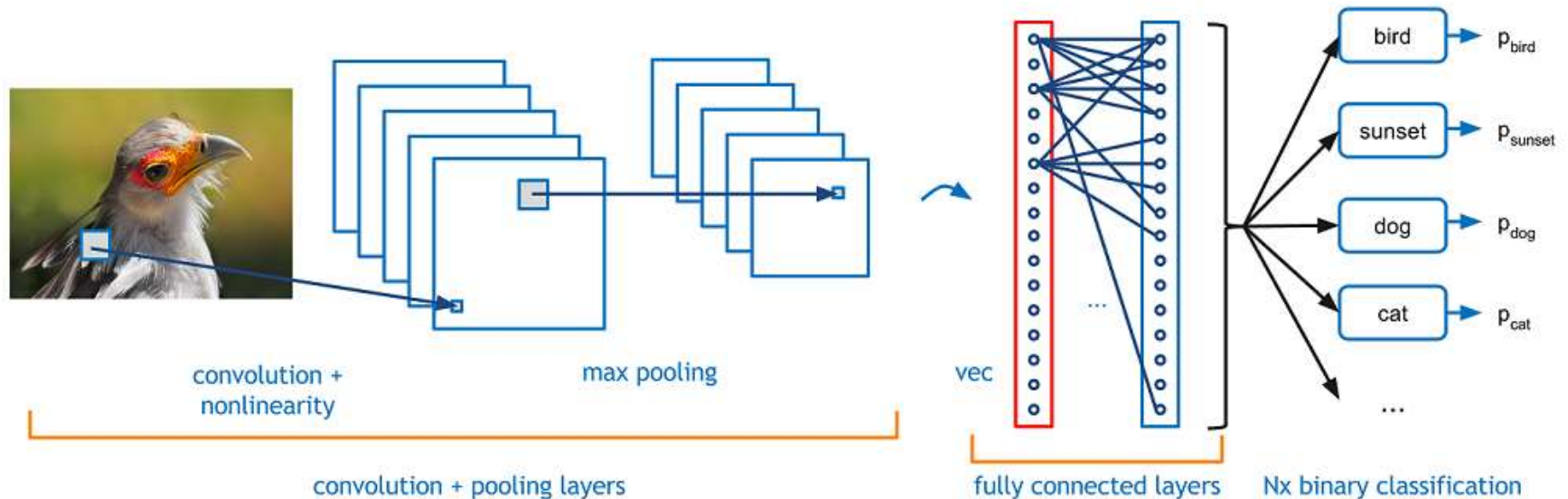
Output: 24x24x6

Typical shape of a classical CNN



Spatial resolution is decreased e.g. via max-pooling while more abstract image features are detected in deeper layers.

A classical CNN has fc layers at the end



In a classical CNN we start with convolution layers and end with fc layers.

The task of the convolutional layers is to extract useful features from the image which might be appear at arbitrary positions in the image.

The task of the fc layer is to use these extracted features for classification.

Fast prototyping CNNs with Keras

For a plain CNN we can use Keras' [Sequential API](#)

```
from keras.models import Sequential
from keras.layers import Dense, Activation, Dropout, BatchNormalization
from keras.layers import Convolution2D, MaxPooling2D, Flatten
import keras
```

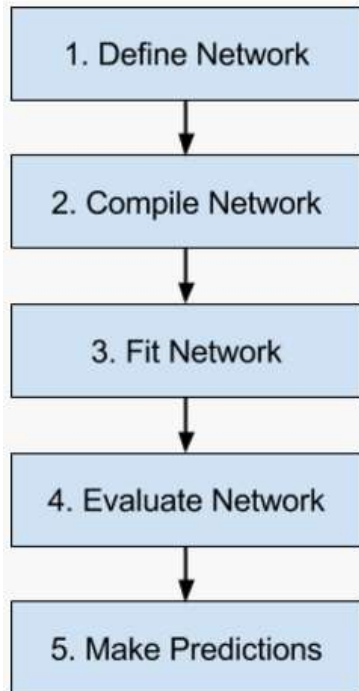
For documentation see:

<https://keras.io/layers/core/>

<https://keras.io/layers/convolutional/>

<https://keras.io/layers/pooling/>

Keras: A high level API with best practice defaults



```
model = Sequential()  
# this applies 32 convolution filters of size 3x3 each.  
model.add(Convolution2D(32, (3, 3), border_mode='same', input_shape=(28, 28, 1)))  
model.add(Activation('relu'))  
model.add(Flatten())  
model.add(Dense(10))  
model.add(Activation('softmax'))
```

Number filter/feature maps (activation maps) in first hidden layers

Will arrange the $32 \times 28 \times 28 = 25088$ neurons, that are located in 32 feature maps of dim 28×28 , in one vector.

```
model.summary()
```

Layer (type)	Output Shape	Param #	Connected to
convolution2d_2 (Convolution2D)	(None, 28, 28, 32)	320	convolution2d_input_2[0][0]
activation_3 (Activation)	(None, 28, 28, 32)	0	convolution2d_2[0][0]
flatten_2 (Flatten)	(None, 25088)	0	activation_3[0][0]
dense_2 (Dense)	(None, 10)	250890	flatten_2[0][0]
activation_4 (Activation)	(None, 10)	0	dense_2[0][0]

Exercise on setting up a simple CNN in keras

Check out the architecture of the CNN described in [“live cnn in browser”](#)

And fill in the pieces to get a Keras code for a model with this architecture

Check out the default parameters and the keras syntax starting from:

<https://keras.io/layers/convolutional/> <https://keras.io/layers/core/> or google...

```
model = Sequential()
model.add(Conv2D(filters= ...,
                 kernel_size=(..., ...),
                 input_shape=(..., ..., ...))
model.add(Activation('...'))
model.add(Conv2D(filters= ...,
                 kernel_size=(..., ...))
model.add(Activation(...))
model.add(MaxPooling2D(pool_size=(..., ...)))
model.add(Dropout(...))
model.add(Flatten())
model.add(Dense(...))
model.add(Activation('...'))
model.add(Dropout(...))
model.add(Dense(...))
model.add(Activation('softmax'))
```

write a digit



Exercise on setting up a simple CNN in keras

Check out the architecture of the CNN described in [“live cnn in browser”](#)
And fill in the pieces to get a Keras code for a model with this architecture

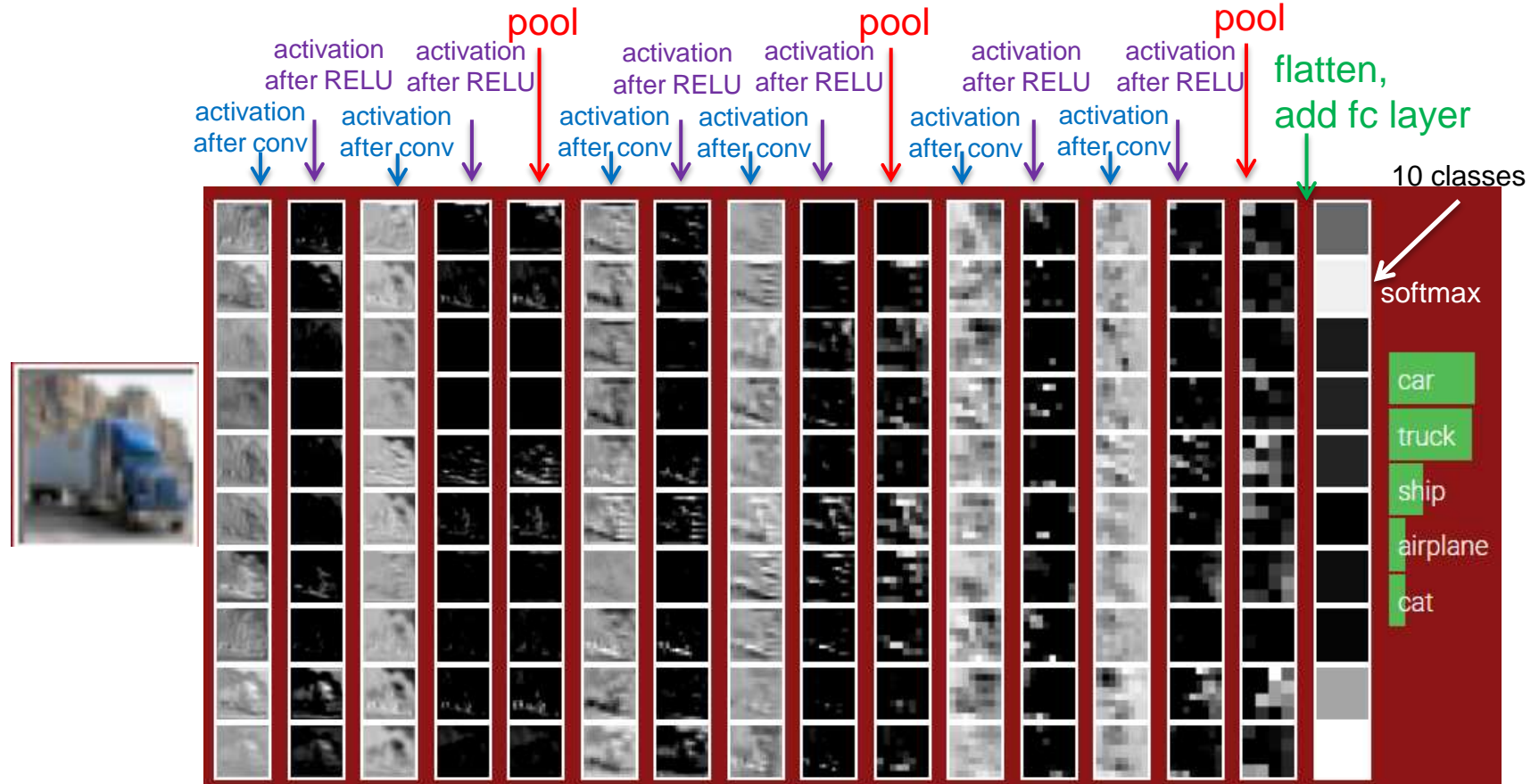
Solution:

```
model = Sequential()
model.add(Conv2D(filters=32,
                  kernel_size=(3, 3),
                  input_shape=(28,28,1))
model.add(Activation('relu'))
model.add(Conv2D(32, (3, 3)))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))
model.add(Flatten())
model.add(Dense(128))
model.add(Activation('relu'))
model.add(Dropout(0.5))
model.add(Dense(num_classes))
model.add(Activation('softmax'))
```





Appearance of activation/feature maps in different layers



Activation maps give insight on the spatial positions where the filter pattern was found in the input **one layer below** (in higher layers activation maps have no easy interpretation) -> only the activation maps in the first hidden layer correspond directly to features of the input image.

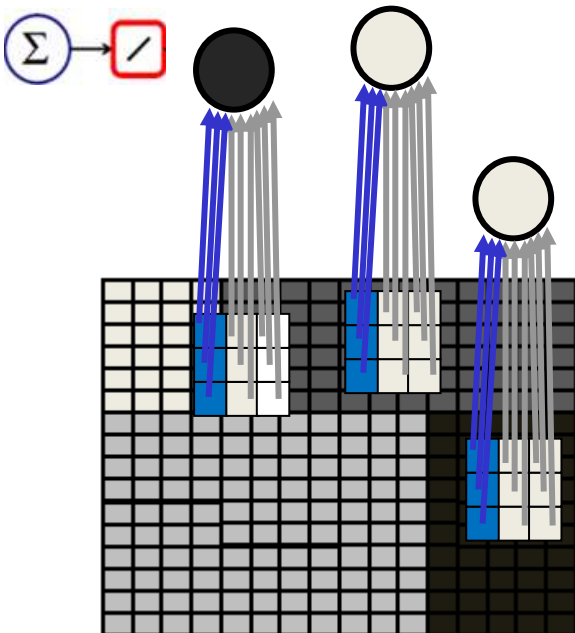
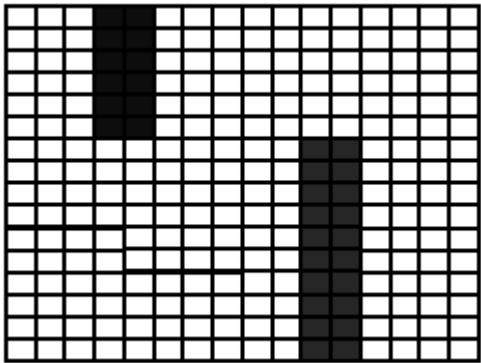
Exercise: use CNN for mnist classification

- Work through the instructions in 07 and 08 CNN [Exercises in day4](#) and use the ipython notebooks that are referred to.



Standardizing data is more important in CNNs than fcNN

feature/activation map



input

Since we share weights in CNNs – only one filter is required per feature map - the **weights in a filter should be appropriate for each patch of the input image**, i.e. yielding inputs to the activation function that are party in their sweet spot.

To ensure that different image patches or even different **pixels yield comparable ranges** of values as input to the filter we need to **standardize the input pixel-wise** (we also restrict the variation to get small activations corresponding to uncertain classifications in the beginning of the training).

```
print(X_train.shape)
print(X_val.shape)
```

```
(4000, 28, 28, 1)
(1000, 28, 28, 1)
```

```
# here we center and standardize the data per pixel
# calculate mean, std over all training images at each pixel position
X_mean = np.mean( X_train, axis = 0)
X_std = np.std( X_train, axis = 0)

X_train = (X_train - X_mean ) / (X_std + 0.0001)
X_val = (X_val - X_mean ) / (X_std + 0.0001)
```

Wrapping up today's story

Why going from fully connected NN to convolutional NN?

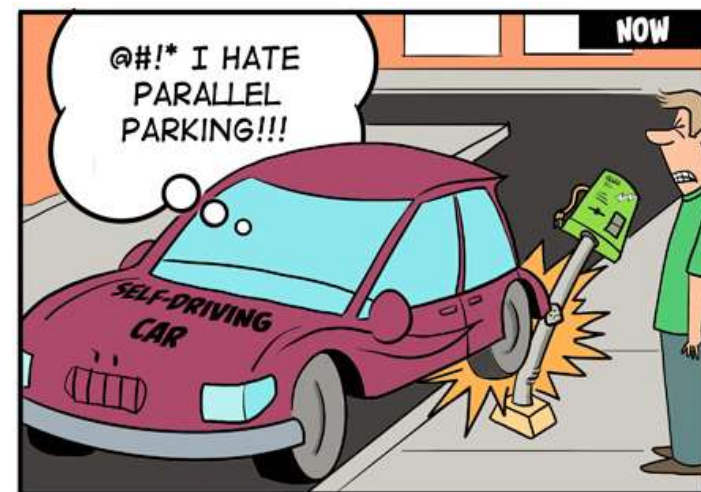
- Most images are much too complicated to be captured by simple template matches of whole shapes (1 layer fc NN).
- Many fully connected layer correspond to many rigid templates.
- We need to learn the features that a image is composed of.
- The classification should not depend on the location of the object in the image
- We want to exploit the information that is contained in the neighborhood structure of pixels in a image.
- Next week we will see how to get insight to the hierarchical features learned by the CNN.

What kind of tasks can be tackled by CNNs?

Convolutional Neural Nets are used for **detecting patterns** in images, videos, sounds and texts.

Where are CNNs used already?

- Recommendation at Spotify, Amazon ...
(<http://benanne.github.io/2014/08/05/spotify-cnns.html>)
- Google, Facebook for image interpretation
e.g. PlaNet—Photo Geolocation
(<http://arxiv.org/abs/1602.05314>)
- Who else is using CNNs?
(<https://www.quora.com/Apart-from-Google-Facebook-who-is-commercially-using-deep-recurrent-convolutional-neural-networks>)



...

Your project to be presented in week 8 ?

Homework: Do some real stuff

Team-up for your first real DL project:

Develop a DL model to solve this task:

For a given image, decide which out of 8 celebrities is on the image.

Data:

For each of the 8 celebrities you get 250 images in the training data set, 50 images in the validation data set and 50 images in test data set.

Special challenge:

The images come from the OXFORD VGG Face dataset. The images were derived from the internet and automatically labeled. The data set contains also mislabeled images or ambiguous images.

Example images:

Label: Steve Jobs (entrepreneur)



Label: Emma Stone (actress)

