

# Deep learning notes

Erik Pärstrand

April 2018

# 1 Hyper-parameter tuning

## When to use deep learning

The following heuristic can be used to determine if deep learning is appropriate:

1. Little noise, complex structure: use deep learning
2. Lots of noise, little structure: don't use deep learning

In order to determine which type of neural network to use, the following heuristic can be used:

1. No structure: feed forward network
2. Spatial structure: convolutional networks
3. Sequential structure: recurrent network

## Bias and Variance in deep learning

Assume the following performance for a classifier (presume that Bayes error-rate is 0%):

Error	Example 1	Example 2
Train	1%	15%
Validation	11%	16%

Example 1 has low bias, high variance. Example 2 has high bias, low variance. Below is a "basic recipe"/heuristic for training a deep neural network:

1. Has the model high bias? Try: train a larger network, train longer/use better optimization algorithm, try different network architecture/hyper-parameter-search. If the model has low bias, go to step 2.
2. Has the model high variance? Try collect more data, add regularization ( $L_2$ , dropout etc), try different network architecture/hyper-parameter-search. If the model has low bias and low variance, go to step 3.
3. Test if the model works well in practice. See orthogonalization section for that.

## Shallow versus deep networks

Some function can be estimated with shallow networks. However, shallower networks requires exponentially more hidden units to approximate a function compared to deeper ones.

As a rule of thumb, start with one hidden layer. Then increase if needed (i.e. to reduce bias).

## Weight regularization

Two common weight regularizer are  $L_1$  and  $L_2$ , i.e. add  $W^2$  or  $|W|$  to the cost function.  $L_1$  creates a sparse weight matrix (in theory, however usually not the case in practice).  $L_2$  (sometimes called weight decay) is the most commonly used in practice.

## Dropout

Randomly "drop" neurons during training. Intuition about dropout: The network can't rely on certain feature, so it has to spread the weights which leads to shrinking weights. It has similar effect as  $L_2$ -norm, however they are not the same.

One can use lower probability of keeping neurons for larger layers (i.e. more hidden units). Rarely used in input layer. One downside of dropout is that the cost function becomes stochastic.

## Data augmentations

Commonly used in computer vision. For example, flip the pictures horizontally or add random distortion. Three common methods are:

1. Mirroring
2. Flip horizontally
3. Random cropping (choose a random subset of the image)

Additional methods are (not often used in practice often):

1. Rotation
2. Sheering
3. Local wrapping
4. Color shifting (add distortion to the different channels)

## Early stopping

Stop the training when the validation error increases. However, this method complicates orthogonalization (see below).

## Normalizing the features

One method is to ensure that all features have zero mean and unit variance. Makes gradient descent converge faster/better. It is only needed when features are on dramatically different scales.

### Vanishing/exploding gradients

The gradient becomes very small/big. Problem for deep networks. Caused since the weight matrices are multiplied.

### Weight initialization

Initialize the weights "smart". Can help with vanishing/exploding gradient. However, don't initialize the weights to zero. This can be done for the bias terms though.

Recommended initialization for layer  $l$  with Relu activation (note that here we use variance and not standard deviation):

$$W \sim N\left(0, \frac{2}{n^{[l-1]}}\right)$$

where  $n^{[l-1]}$  is the number of hidden units in the layer before  $l$ . This is called He-initialization.

For tanh layers it's recommended to use:

$$W \sim N\left(0, \frac{1}{n^{[l-1]}}\right)$$

This is called Xavier-initialization.

### Mini-batch gradient descent

In one epoch, we go through the entire training set once. During training, one can split the data into different batches (mini batches).

Guidelines: if the training set has less than 2000 examples, one can use batch-gradient descent. Otherwise, one should use mini-batch gradient descent.

Typical sizes of the mini-batches are 64, 128, 256, ... (powers of two). One should make sure that the mini-batch data fits in memory for CPU/GPU, otherwise it will drastically decrease the training speed.

### Exponential moving average

$$v_t = \beta + (1 - \beta)\theta_t$$

$v_t$  approximately average over the last  $\frac{1}{1-\beta}$  values.

If one uses the initialization  $v_0 = 0$ , it will be a bias estimate. It can be fixed by  $\tilde{v}_t = \frac{v_t}{1-\beta^t}$ .

### Gradient descent with momentum

Adds a moment term to the gradient descent for the weights and bias. Typically, one uses  $\beta = 0.9$ . Usually works better than regular gradient descent.

### RMSprop

RMSProp automatically adjusting the step size so that the step is on the same scale as the gradients - as the average gradient gets smaller, the coefficient in the gradient descent update gets bigger to compensate.

### Adam

Uses momentum and the RMSprop term in the optimization step. Typical value are  $\beta_1 = 0.9$  (momentum term) and  $\beta_2 = 0.999$  (RMSprop term). The default values usually works well in practice (tune if needed). A good rule of thumb is to use Adam as default optimizer.

### Learning rate decay

One can decay the learning rate a bit after each epoch:

$$\eta = \frac{\eta_0}{1 + DecayRate \times EpochNumber}$$

However, the most common approach is to decay in steps, i.e. decay the learning rate by 10% every 5:th epoch, for example.

### Hyper-parameter tuning

To achieve low bias, learning rate  $\eta$  is usually the most important hyper-parameter. Then the number of hidden units in the layers. Then the number of layers and learning rate decay. Last, one can try to change Adam's hyper-parameters. However, this order depends on the problem.

To achieve low variance, regularization is the most important together with the ones above.

One should use random search instead of grid search. It is a good idea to start with large ranges of the hyper-parameters to see which values makes "sense" (coarse search) and then create smaller ranges from these results (fine search).

If the hyper-parameters are on the same scale, one can sample uniformly. However, if they are not one should sample on log-scale. Assume that we want to investigate the range  $[10^{-9}, 10^{-6}]$ . Sample  $r \in [-9, -6]$  uniformly. Then the hyper-parameters values are set to  $10^r$ .

### **Panda versus caviar approach**

One can baby sit one model (called the panda approach). Good when one doesn't have a lot of computational resources. Check the performance and update the model/hyper-parameters accordingly.

One can also train many models in parallel (called caviar approach). Good if one has a lot of computational resources.

### **Batch-normalization**

Normalize the mean and variance before each activation step (some advocate to normalize after the activation step), i.e.

$$z_{norm}^{(i)} = \frac{z^{(i)} - \mu}{\sigma}$$

$$\hat{z}^{(i)} = \gamma z_{norm}^{(i)} + \Gamma$$

where  $\gamma$  and  $\Gamma$  are learnable parameters through back-propagation. Use  $\hat{z}^{(i)}$  in forward and back-propagation. Useful for covariance shift, i.e. reduce distribution shift for different layers.

It has some regularizing effect, however it should not be used for that.

During training, keep track of  $\mu$  and  $\sigma$  with an EMA from the mini-batches. Use the mean and standard deviation to normalize  $z$  at test time.

## 2 Practical aspect of deep learning

### Orthogonalization

Split a problem into orthogonal parts:

1. Fit the training set well given a cost function. This can be improved by larger networks, better optimizer (Adam etc). If satisfied, move to (2).
2. Fit the validation set well given a cost function. This can be improved by regularization, more data etc. If satisfied, move to (3).
3. Make sure that the model works well on the test set. This can be improved by getting a larger validation set. If satisfied, move to (4).
4. Make sure that the model works well in practice. This can be improved by changing the validation set or cost function.

### Single number evaluation metric

Before creating a model, define and use a single number evaluation metric (accuracy,  $F_1$  etc). Use your validation set and evaluation metric to choose model/direction. This leads to faster iterations of the modelling.

### Satisfying and optimizing metric

Use the evaluation metric to (globally) optimize the model. However, one can have satisfying metric that the model needs to accomplish, for example run-time. Choose the model that has the highest/lowest optimizing metric s.t. the model accomplish the satisfying metrics.

For example, maximizing accuracy s.t. the run-time of the model is less then 1000 ms. If you have  $N$  metrics of interest, optimize one and have the  $N - 1$  as satisfying metrics.

### train/validation/test distribution

Validation and test set should come from the same distribution. However, the training set can come from another one. Choose the validation and test set to reflect the data you expect to get in the future and consider important to do well on.

### Rule of thumb the size of the different sets

If data size  $< 10,000$ : Use 60% for training, 20% for validation and 20% for test.

If data size  $> 1,000,000$ : Use 98% for training, 1% for validation and 1% for test.

The size of the test set should be big enough to give a high confidence in the overall performance of the model. It's recommended to have a train, validation and test set (to get an unbiased estimate of the performance).

If the current metric doesn't reflect your "goal", redefine it. For example, add weights to the different classes when computing the accuracy. Can also add weights to the loss function. Moreover, if doing well on your metric + validation/test set doesn't correspond to doing well on your application, change your metric and/or validation/test set.

Furthermore, if the performance of the validation and test set differ a lot, you have over-fitted the validation set and should get more validation data.

### Comparing to human-level performance

Humans are quite good at a lot of tasks, for example image/speech recognition etc. So as long as the model is worse than humans, you can:

1. Get labeled data from humans
2. Gain insight from manual error analysis
3. Better analysis of bias/variance

Human-level error can be viewed as a proxy for the Bayes error-rate (i.e. lowest possible error).

Below is an example of a cat classifier:

Error	Example 1	Example 2
Human	1%	7.5%
Train	8%	8%
Validation	10%	10%

In example 1, one should focus on bias reduction. In example 2, one should focus on reducing variance.

Methods for reducing bias: train bigger model, train longer/use better optimization algorithm (RMSprop, Adam etc) or change network architecture/hyper-parameter search.

Methods for reducing variance: get more data, regularization ( $L^2$ , dropout, data augmentation) or change network architecture/hyper-parameter search.



## Error analysis

Manually inspect the errors (on the validation data). Example: should you try to make your cat classifier do better on dogs? Below is an example of a workflow for error analysis (assume that the validation error is 10%):

1. Get  $\sim 100$  mislabeled validation data points.
2. Manually go through the errors, count how many are dogs?
3. If 5/100 are dogs, probably not worth investigating the cause of the error.
4. If 50/100 are dogs, probably worth investigating.

Ideas to improve the cat classifier:

1. "Fix" images of dogs being recognized as cats
2. "Fix" great cats (lions etc) being mis-classified
3. Improve performance on blurry images
4. etc

Document the validation error (in a spread-sheet for example):

Images	Dogs	Great cats	incorrect labeled	...	comments
1	Yes	No	...	...	Pitbull
2	No	Yes	...	...	Lion
...	...	...	...	...	...
% of total errors	8%	43%	10%	...	...

## Cleaning up incorrect labeled data

Deep learning algorithms are quite robust to random errors in the training set. However, less robust to systematic errors. In error analysis, add columns for incorrect labels in the table above. If it makes a significant different, one can try to fit it.

However, one should apply the same process to your validation and test set to make sure they come from the same distribution.

## Build your first system quickly, then iterate

1. Set up your validation and test set and choose a metric
2. Build your first system quickly, then iterate
3. Use bias/variance analysis and error analysis to prioritize the next steps

Less applicable if you have prior knowledge within the field or if there is a large academic literature on the subject.

## Training and testing on different distributions

Example: You have 200,000 pictures of cats from the web (web-crawling) and 10,000 from your app (i.e. the performance you care about). Use the web pictures for training (maybe some data from the app). Use all/most data from the app for the validation and test set.

## Bias variance with mismatching data distributions

Introduce another set, training-validation set. Comes from the same distribution as the training set, but not used for training. Example for the cat classifier (assume Bayes error-rate is 0%):

Error	Example 1	Example 2	Example 3
Train	1%	1%	10%
Train-validation	9%	1.5%	11%
Validation	10%	10%	12%

In example 1, we have a variance problem. In example 2, we have a data-mismatch problem. In example 3, we have a bias problem.

## Addressing data mismatch

Potential routes:

1. Carry out manual error analysis to try to understand the difference between train and validation/test set (i.e. if they have different distributions).
2. Make training data more similar to validation/test data, or collect more data similar to the validation/test set.

Can also try artificial data synthesis. For example, if there is a lot of car noise in the validation/test set, try to add original voice + car noise to generate the synthetic data. However, one should be careful to not overfit the car noise in this example if we use the same noise multiple times.

## Transfer learning

Use pre-trained networks and tune the last/multiple layers for your application. The transfer is from A to B. Make sense when:

1. Task A and B have similar type of input (images etc)
2. You have a lot more data for task A than for task B
3. Low level features from A could be helpful for task B

Rule of thumb regarding how many layers should be re-trained:

If you have a small amount of data, replace the last few layers (i.e. the fully connected layers and output layer). The earlier layers are "freezed", i.e. not updated in gradient descent. Called fine-tuning.

If you have a lot of data, you can re-train the entire network. The earlier layers become initialization. However, the fully connected layers and output layer need to be re-trained from scratch. Called pre-training.

### **End-to-End deep learning**

Let the network come up with the features, i.e. use "raw" data.

Pros:

- Let the data speak
- Less hand-designing of components

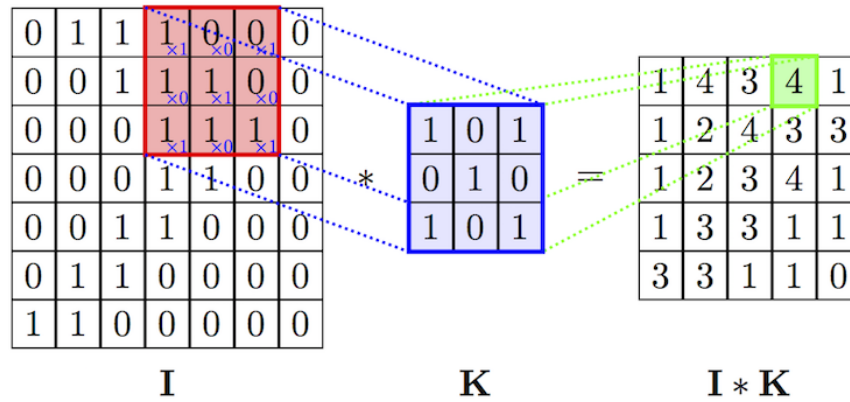
Cons:

- May need a lot of data
- Exclude potentially useful hand-designed components.

### 3 Convolutional networks

#### Convolutional operation

Convolution occurs between two matrices, an image (**I** below) of size  $n \times n$  and a filter (**K** below) of size  $f \times f$ , where  $f$  is usually odd. Below is an example of a convolution operation between **I** and **K** (stride is one):



Traditionally, **K** was constructed by hand. But with deep learning, they are estimated with back propagation.

#### Padding

When doing convolution the edges of the picture gets included in fewer convolutional operations then the central pieces. Therefore, one can add zeros around the edges to address this problem, i.e. zero padding. The picture below is an example when we have added padding of size 1 ( $p=1$ ):

0	0	0	0	0	0
0	35	19	25	6	0
0	13	22	16	53	0
0	4	3	7	10	0
0	9	8	1	3	0
0	0	0	0	0	0

Also, one problem for convolutional networks is that the size of the picture decreases after each convolution. Therefore, zero padding helps to keep the size of the image constant if so needed. If we want the input size to be equal to the output size, set  $p = \frac{f-1}{2}$ . There are two common terms for padding:

- Valid convolution - use no padding
- Same convolution - keep the output size the same as input size (see above).

## Strided convolution

In the convolutional operation section, we used a stride of size 1 ( $s=1$ ), i.e. we moved the convolutional filter one step after each operation. However, one can use a larger number, for example  $s = 2$ . Then we will jump two steps after each operation. The output size with a given stride ( $s$ ), image size ( $n$ ), padding size ( $p$ ) and filter size ( $f$ ):

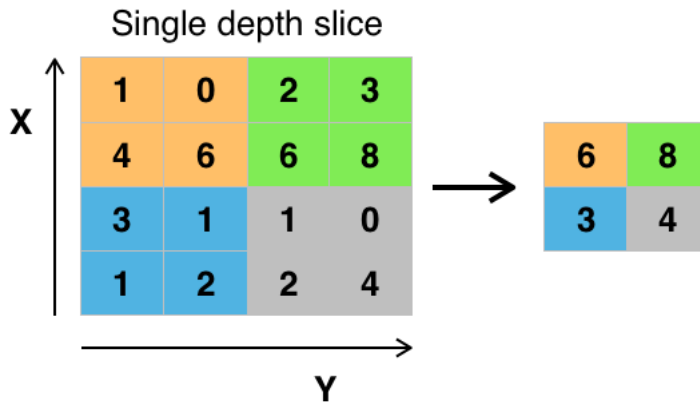
$$\text{output} = \frac{n + 2p - f}{s} + 1$$

## Convolution over volume

Previously, we did convolution over an 2D-area (e.g. grey-scaled image). However, we could also do it over a volume, for example  $n \times n \times 3 * f \times f \times 3$ , i.e. we have stacked 3 filters/channels.

## Pooling

The max-pooling operation is illustrated below (with filter size 2 and stride 2):



If a useful feature is detected, keep it. You can look at it a little like focusing the network's attention on what's most characteristic for the image at hand.

One can also use average pooling, i.e. take the average within the region. Average pooling is sometimes used in deep networks in the end. However, but max-pooling is used more often for most layers.

## Convolution networks advice and intuition

Earlier layers detects lower level features (edges etc) and later layers detects higher level features (faces etc).

In ConvNets, there are three types of layers; convolutional, pooling and fully connected layer. For network architecture, look in the literature for advice. One typical layer sequence is Conv → Pooling → Conv → .. → Fully connected. Also, typical value of the filters and strides are:

$$f = 2, s = 2$$

$$f = 3, s = 2$$

Typically, for ConvNets  $n$  decreases along the layers and the number of channels increases. However, we don't want the size ( $n$ ) to drop too fast.

## Why convolution?

Below are some advantageous properties of convolutional layers:

### Number of parameters

Assume that we have a picture of size  $1000 \times 1000 \times 3$  and that we will use a fully connected layer of size 1000. That would becomes 3 billion parameters in

the first layer, which is intractable. With a convolution layer we would only get 27 (assuming  $3 \times 3 \times 3$  filter).

### Parameter sharing

A feature detector that's useful in one part of the image is probably useful in another part of the image.

### Sparsity of connections

In each layer, each output value depends only on a smaller number of local inputs.

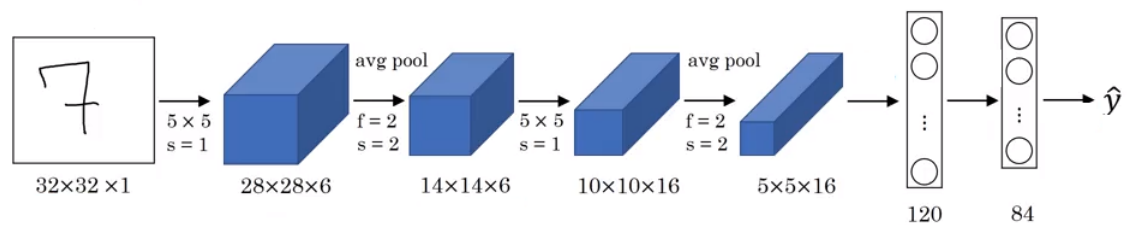
### Translation invariance

It doesn't matter how the object is placed. For example, it classify a cat even if it's rotated.

## Common architectures

### LeNet

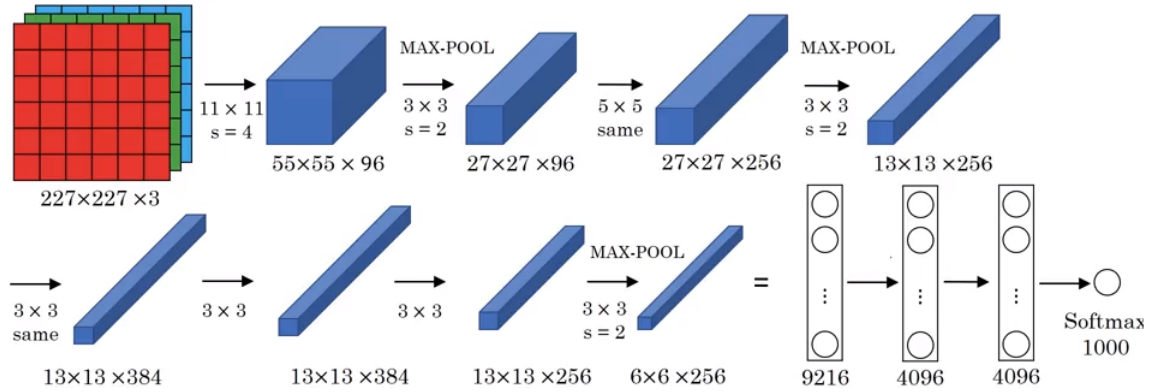
The first convolutional networks. Has around 60k parameters. Deeper in the model the height/width decreases and the number of channels increases. Below is the architecture of it:



### AlexNet

One of the first "new" convolutional networks. Has around 60M parameters. Deeper in the model the height/width decreases and the number of channels increases. Below is the architecture of it:

## AlexNet



## VGG-16

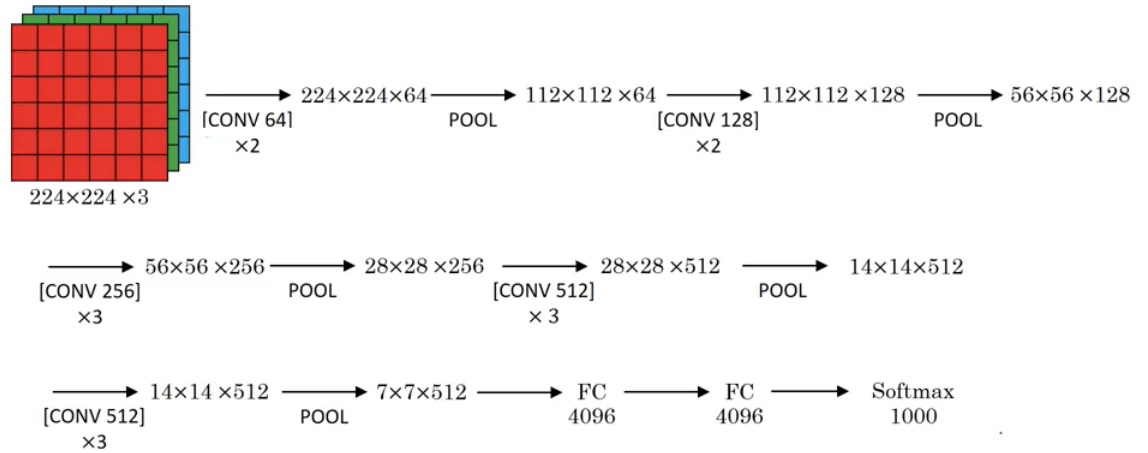
Has around 138M parameters. Deeper in the model the height/width decreases and the number of channels increases. The height/width decreases with a factor of two with each "block" and the number of channels increases by a factor of two with each "block". Below is the architecture of it:



# VGG - 16

CONV =  $3 \times 3$  filter,  $s = 1$ , same

MAX-POOL =  $2 \times 2$ ,  $s = 2$



## ResNet

Use skip connections which allows for very deep networks. Below is residual block:

# Residual block

Diagram illustrating a Residual block structure:

The main path (green underline) consists of:

$$a^{[l]} \rightarrow \text{Linear} \rightarrow \text{ReLU} \rightarrow a^{[l+1]} \rightarrow \text{Linear} \rightarrow \oplus \rightarrow \text{ReLU} \rightarrow a^{[l+2]}$$

A skip connection (purple line) bypasses the main path, adding  $a^{[l]}$  directly to the output of the second Linear layer:

"short cut" / skip connection

Equations for the layers:

$$z^{[l+1]} = W^{[l+1]} a^{[l]} + b^{[l+1]} \quad a^{[l+1]} = g(z^{[l+1]}) \quad z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]}$$

Handwritten notes: "main path" (green), "short cut / skip connection" (purple).

# ResNet

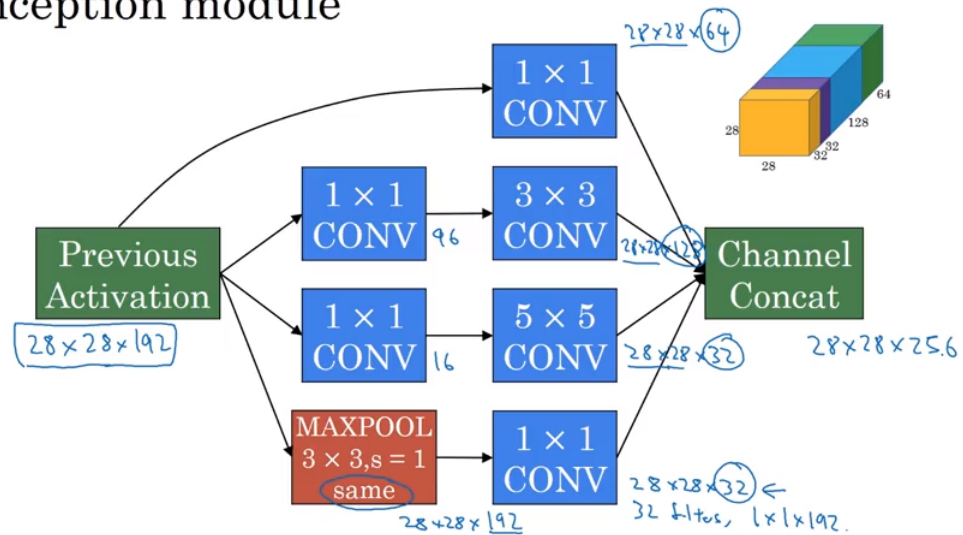
The diagram illustrates the ResNet architecture, comparing a 34-layer plain version (top) with a 34-layer residual version (bottom). The layers are color-coded: orange for the initial 7x7 conv, yellow for 3x3 convs, and pink for 1x1 convs. The residual version includes skip connections (curved arrows) for residual blocks. The layers are as follows:

- 34-layer plain:** Image input → 7x7 conv, 64, /2 → pool, /2 → 3x3 conv, 64 → 3x3 conv, 64 → 3x3 conv, 64 → 3x3 conv, 64 → 3x3 conv, 64 → 3x3 conv, 128, /2 → 3x3 conv, 128 → 3x3 conv, 128 → 3x3 conv, 128 → 3x3 conv, 128 → 3x3 conv, 128 → 3x3 conv, 256, /2 → 3x3 conv, 256 → 3x3 conv, 256 → 3x3 conv, 256 → 3x3 conv, 256 → 3x3 conv, 256 → 3x3 conv, 512, /2 → 3x3 conv, 512 → 3x3 conv, 512 → 3x3 conv, 512 → 3x3 conv, 512 → avg pool → fc 1000.
- 34-layer residual:** Image input → 7x7 conv, 64, /2 → pool, /2 → 3x3 conv, 64 → 3x3 conv, 64 → 3x3 conv, 64 → 3x3 conv, 64 → 3x3 conv, 64 → 3x3 conv, 128, /2 → 3x3 conv, 128 → 3x3 conv, 128 → 3x3 conv, 128 → 3x3 conv, 128 → 3x3 conv, 128 → 3x3 conv, 256, /2 → 3x3 conv, 256 → 3x3 conv, 256 → 3x3 conv, 256 → 3x3 conv, 256 → 3x3 conv, 256 → 3x3 conv, 512, /2 → 3x3 conv, 512 → 3x3 conv, 512 → 3x3 conv, 512 → 3x3 conv, 512 → avg pool → fc 1000.

## Inception network

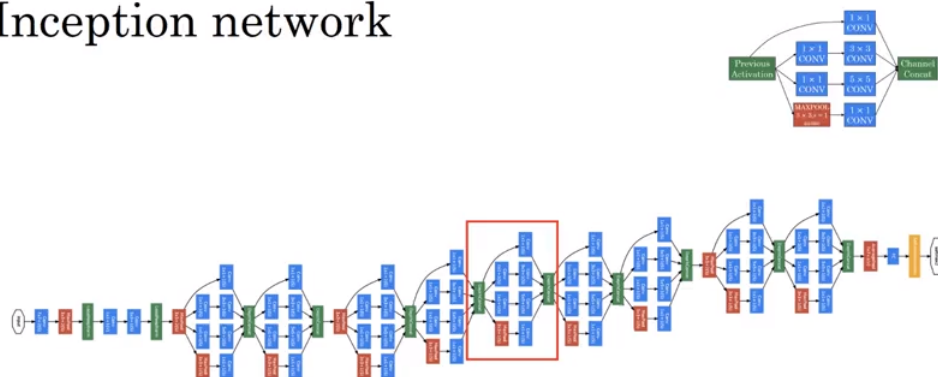
"Network in Network". Uses  $1 \times 1$  filter to decrease the number of channels.  
One inception module is:

### Inception module



The inception network is:

### Inception network



## Object detection

In object detection, one tries to identify and localize multiple objects, i.e. set bounding box around the objects.

One simple method is to train a conv-net on small (cropped) images. Then use a sliding windows on the "big" image. One typically start with a small window and increase it sequentially. However this is not efficient (in terms of computational time).

One improvement is to share the computation between windows (i.e. convolutional implementation of the sliding window). The output becomes a tensor. However the bounding boxes are usually not good. One improvement to this is the YOLO-algorithm.

### Intersection over union

Measures overlap between two bounding boxes. Defined as:

$$\text{IoU} = \frac{\text{Area of intersection}}{\text{Area of union}}$$

One usually claims that there is overlap if  $\text{IoU} > 0.5$ .

### None-max suppression

Find highest probability within each box. Suppresses other boxes that have high IoU.

### YOLO-algorithm

YOLO reframe object detection as a single regression problem, straight from image pixels to bounding box coordinates and class probabilities.

A single convolutional network simultaneously predicts multiple bounding boxes and class probabilities for those boxes. YOLO trains on full images and directly optimizes detection performance. This unified model has several benefits over traditional methods of object detection, such:

- Speed
- Reasons globally about the image when making predictions
- Learns generalizable representations of objects

The system divides the input image into an  $S \times S$  grid. If the center of an object falls into a grid cell, that grid cell is responsible for detecting that object.

Each grid cell predicts  $B$  bounding boxes and confidence scores for those boxes. Formally we define confidence as  $P(\text{Object}) * \text{IoU}$ . If no object exists in that cell, the confidence scores should be zero. Otherwise we want the confidence score to equal the intersection over union (IoU) between the predicted box and the ground truth.

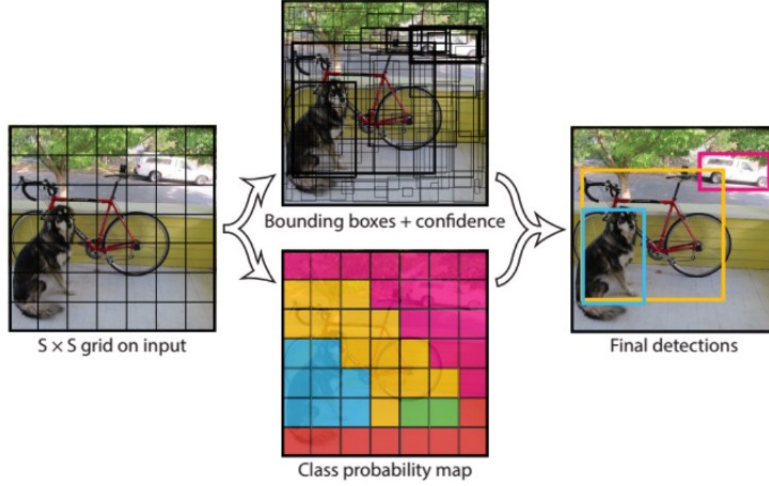
Each bounding box consists of 5 predictions:  $x$ ,  $y$ ,  $w$ ,  $h$ , and confidence. The  $(x, y)$  coordinates represent the center of the box relative to the bounds of the grid cell. The width and height are predicted relative to the whole image.

Finally the confidence prediction represents the IoU between the predicted box and any ground truth box. Each grid cell also predicts  $C$  conditional class probabilities,  $P(\text{Class}_i|\text{Object})$ . These probabilities are conditioned on the grid cell containing an object. We only predict one set of class probabilities per grid cell, regardless of the number of boxes  $B$ .

At test time we multiply the conditional class probabilities and the individual box confidence predictions,

$$P(\text{Class}_i|\text{Object}) * P(\text{Object}) * \text{IOU} = Pr(\text{Class}_i) * \text{IOU}$$

which gives us class-specific confidence scores for each box. These scores encode both the probability of that class appearing in the box and how well the predicted box fits the object.



**Figure 2: The Model.** Our system models detection as a regression problem. It divides the image into an  $S \times S$  grid and for each grid cell predicts  $B$  bounding boxes, confidence for those boxes, and  $C$  class probabilities. These predictions are encoded as an  $S \times S \times (B * 5 + C)$  tensor.

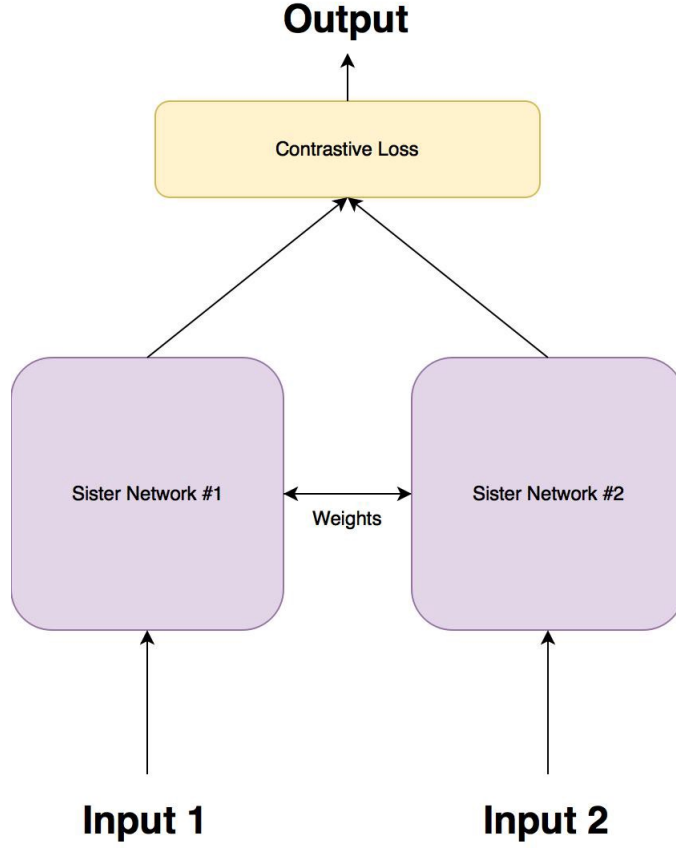
## Face recognition/verification

For face verification, given an input image verify whether the person is the claim one (person).

For face recognition, given a database of  $K$  persons and an input image, check if the image/person is any of these  $K$  persons.

## Siamese network

Siamese neural network is a class of neural network architectures that contain two (or more) identical sub-networks. identical here means they have the same configuration with the same parameters and weights. Parameter updating is mirrored across both sub-networks. Below is an example of the network:



The objective of the network is not to classify input images, but to differentiate between them. So, a classification loss function (such as cross entropy) would not be a good fit. Instead, this network is better suited to use a contrastive function. Intuitively, this function just evaluates how well the network is distinguishing a given pair of images. It's defined as:

$$\frac{(1 - y)}{2}D^2 + \frac{y}{2}(\max(0, m - D))^2$$

where D is defined as the euclidean distance between the outputs of the siamese networks,  $D = \sqrt{(h(X_1) - h(X_2))^2}$ . h is the output of one of the networks,  $X_1$  and  $X_2$  is the input data-pair.

Here, y is either 1 or 0. If the inputs are from the same class, then the its value is 0, otherwise 1.

m is a margin value which is greater than 0. Having a margin indicates that dissimilar pairs that are beyond this margin will not contribute to the loss. This

makes sense, because you would only want to optimise the network based on pairs that are actually dissimilar, but the network thinks are fairly similar.