

# Macro NN architecture

Machine Learning and Data Mining

Maxim Borisyak

National Research University Higher School of Economics (HSE)

# Outline

# Super inspirational quotes

*Network architecture is more like an art.*

*Behind every is a poorly formulated science.*

# Network architecture

*Neural Network Architecture plays crucial role in Deep Learning.*

Most of the non-trivial architectures:

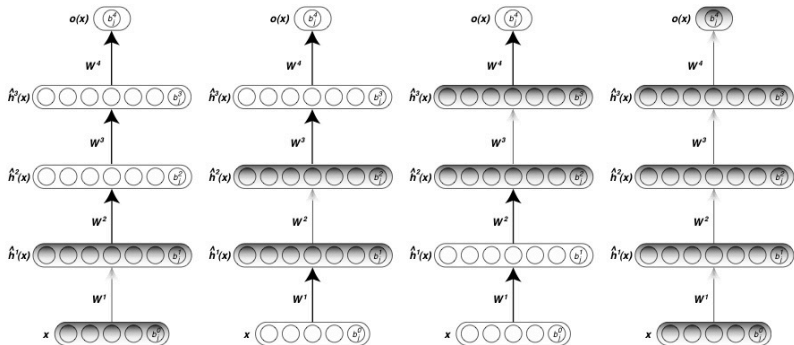
- ❑ derived from common sense;
- ❑ explained by math;
- ❑ demonstrated on some real problems.

# Usual disclaimer

*The following examples are not aimed to be cover major architecture tricks. Just some examples happened to be known by the author.*

# Pretraining

# Layerwise pretraining



# Pretraining

- ❖ layer-wise pretraining:
  - ❖ RBM;
  - ❖ AE;
- ❖ pretraining on simpler but related task.



# Auxiliary losses

# Auxiliary problems

$$\mathcal{L} = \mathcal{L}_{\text{main}} + \lambda_1 \mathcal{L}_1 + \lambda_2 \mathcal{L}_2 + \dots$$

- ❖ solving several objectives with one network:
  - ❖ bringing more information about the solution;
- ❖ auxiliary losses should share the same solution;

Auxiliary losses

Main loss

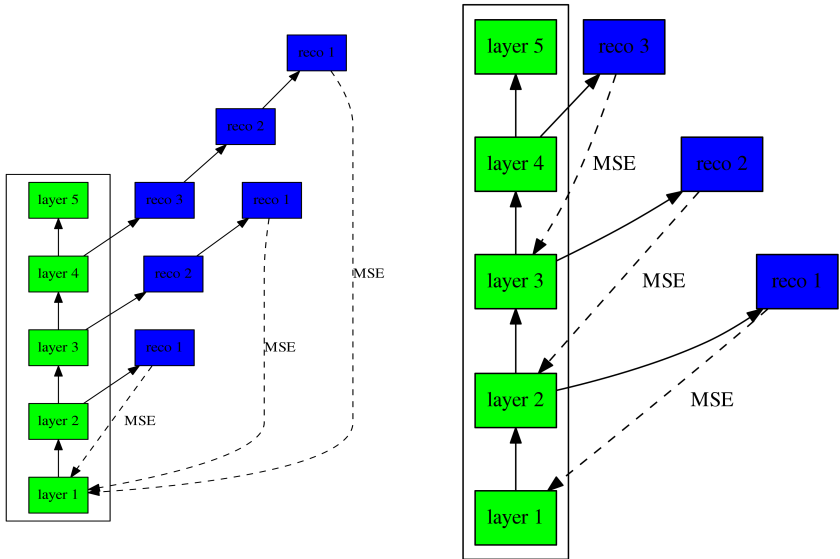


# Auxiliary problems: examples

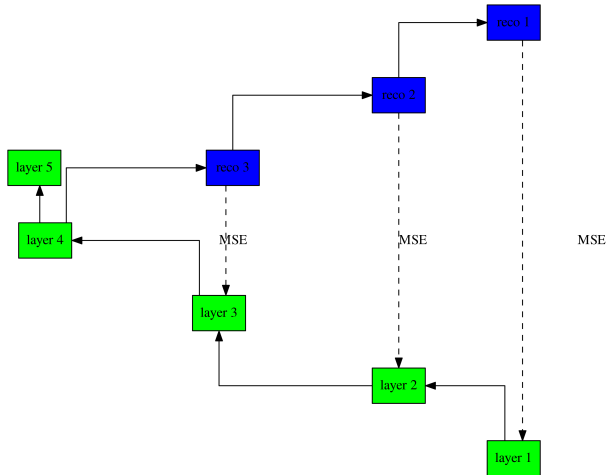
Are the following auxiliary problems reasonable:

- ❑ even vs. odd digit for MNIST;
- ❑ reconstructing initial image for MNIST;
- ❑ producing countour of target objects for detection problems;
- ❑ predicting type of a street-sign for detection problem;
- ❑ predicting super-class for CIFAR-100;
- ❑ predicting faces properies (e.g. smile/anger/neutral, female/male etc) for dimensionality reduction?

# Reconstruction regularisation



# Reconstruction regularization

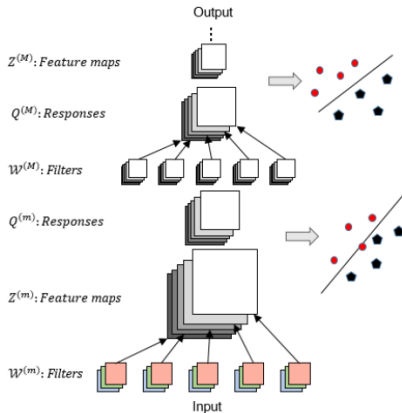


# Reconstruction regularization

- ❖ unsupervised loss may be in conflict with classification loss;
  - ❖ reconstruction generally require higher network capacities;
  - ❖ discriminative features might be lost as unimportant for reconstruction;
- ❖ rarely used in practice.

# Deeply supervised networks

- ❖ try to solve original problem early;
- ❖ improved gradient flow (almost impossible to make it vanish);
- ❖ quite strong regularization effect;
- ❖ no unsupervised vs. supervised conflict.

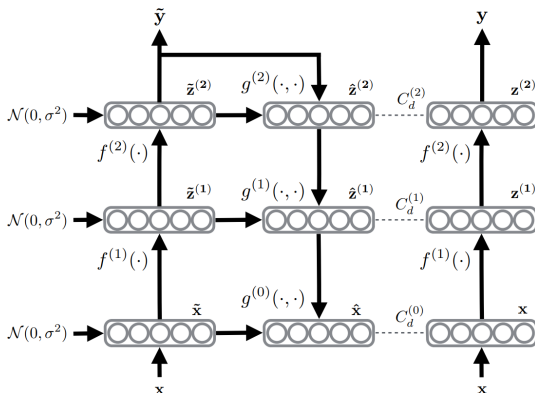


(a) DSN illustration

# Ladder Networks

- replaces reconstruction with denoising:

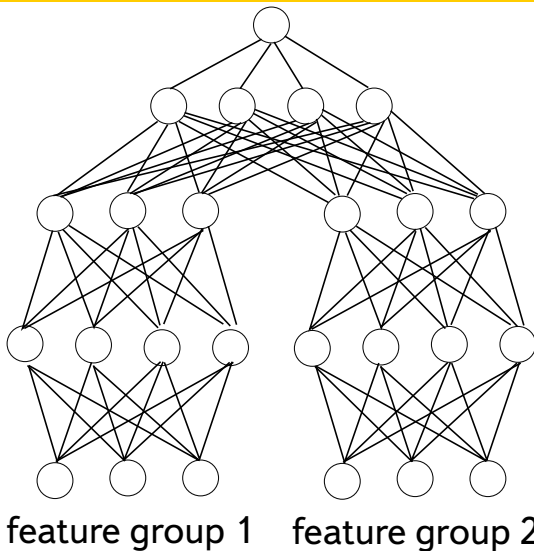
$$\mathcal{L} = \|f(x + \varepsilon) - x\|^2 \rightarrow \min, \varepsilon \sim \mathcal{N}(0, \sigma^2)$$



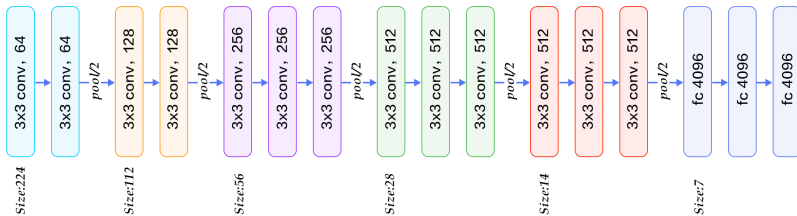


# Network structure

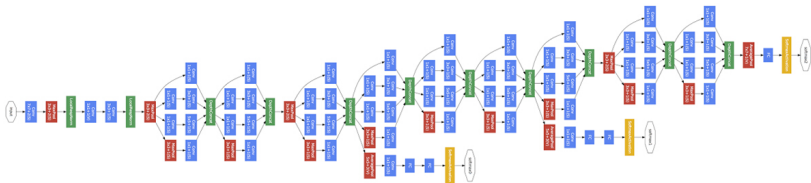
# Tree-like networks



# VGG

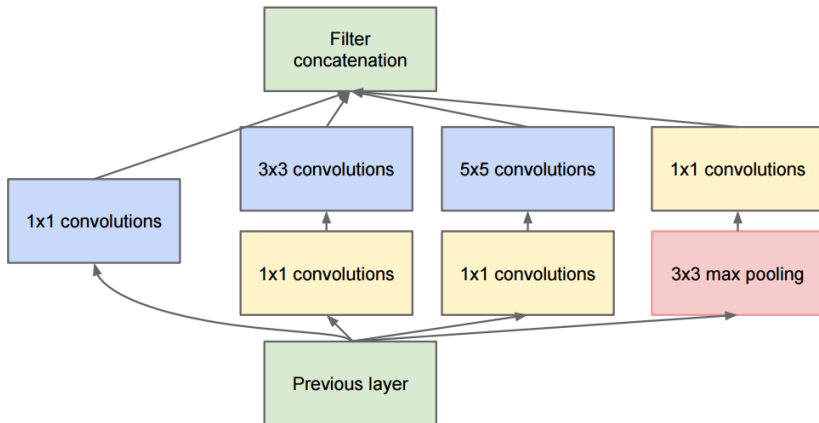


# Inception



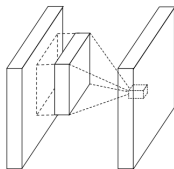
- ❖ blue blocks: conv;
- ❖ red blocks: pool;
- ❖ green blocks: concat;
- ❖ yellow blocks: softmax.

# Inception block

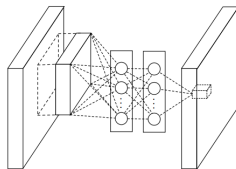


(b) Inception module with dimension reductions

# NIN: conv on steroids



(a) Linear convolution layer



(b) Mlpconv layer

Figure 1: Comparison of linear convolution layer and mlpconv layer. The linear convolution layer includes a linear filter while the mlpconv layer includes a micro network (we choose the multilayer perceptron in this paper). Both layers map the local receptive field to a confidence value of the latent concept.

# ResNet

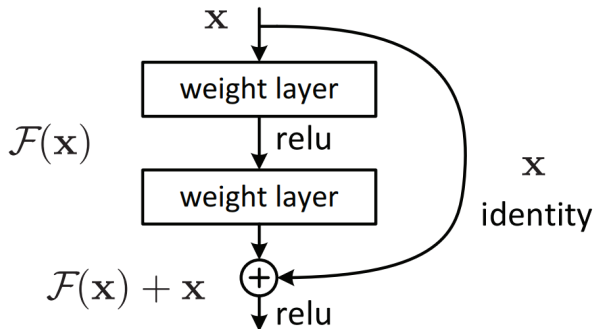
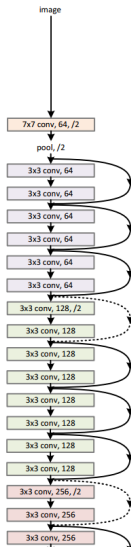


Figure 2. Residual learning: a building block.

# ResNet

34-layer residual





# ResNet

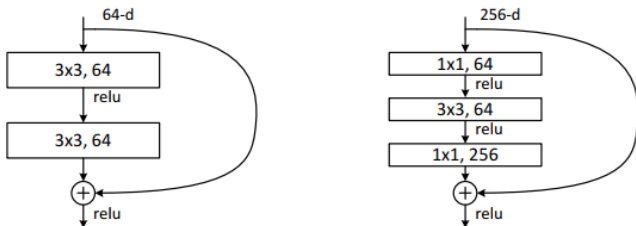


Figure 5. A deeper residual function  $\mathcal{F}$  for ImageNet. Left: a building block (on  $56 \times 56$  feature maps) as in Fig. 3 for ResNet-34. Right: a “bottleneck” building block for ResNet-50/101/152.

# Highway networks

Feed-forward networks:

$$y = H(x, W_H)$$

Residual connection:

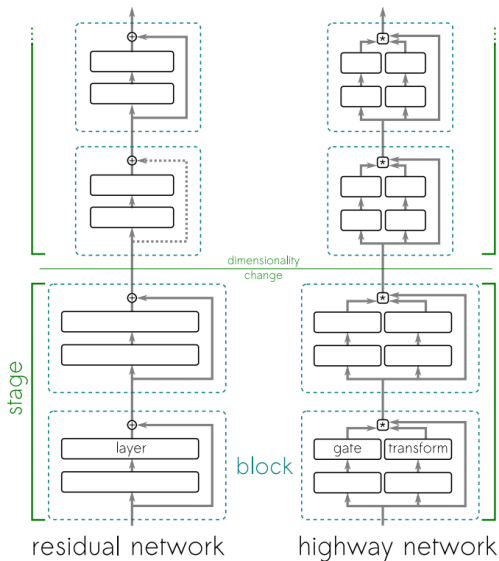
$$y = H(x, W_H) + x$$

Highway connection:

$$y = T(x, W_T)H(x, W_H) + C(x, W_C)x;$$

- ❖  $x, y$  - input, output;
- ❖  $H(x, W_H)$  - some transformation, e.g. convolution;
- ❖  $T(x, W_T), C(x, W_C) \in [0, 1]$  - gates (*transform* and *carry*).

# Highway networks



# Squeeze net

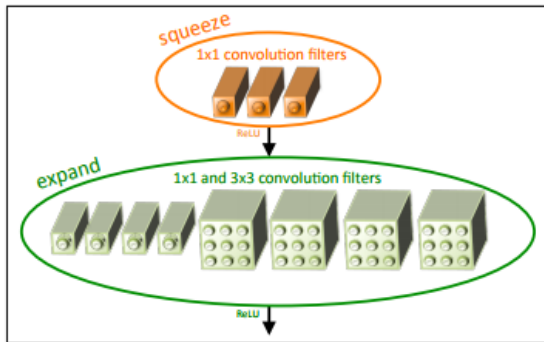
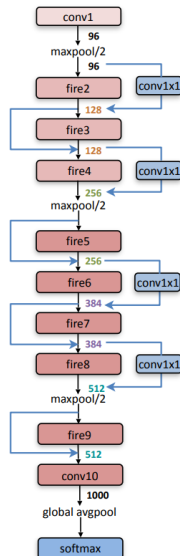
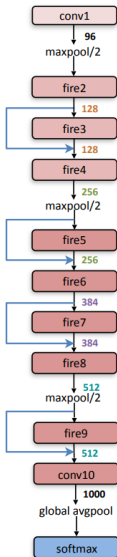
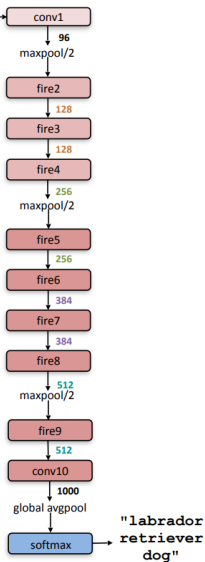
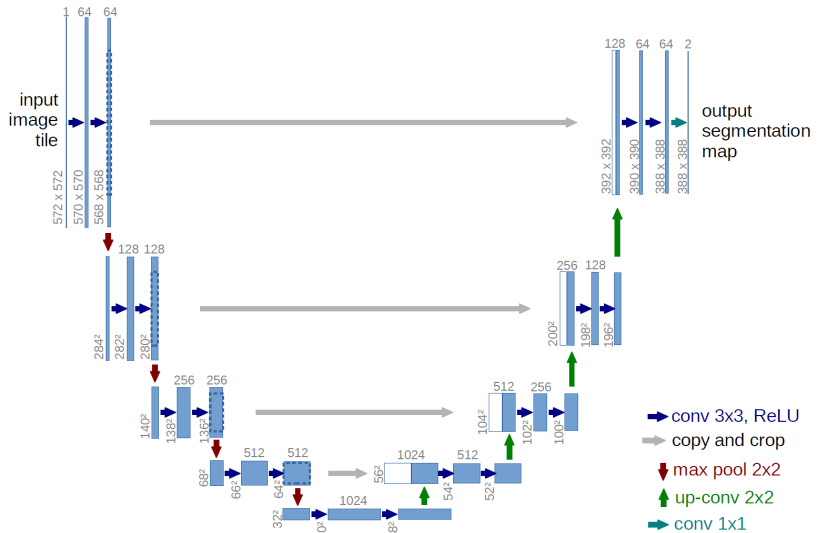


Figure 1: Microarchitectural view: Organization of convolution filters in the **Fire module**. In this example,  $s_{1 \times 1} = 3$ ,  $e_{1 \times 1} = 4$ , and  $e_{3 \times 3} = 4$ . We illustrate the convolution filters but not the activations.

# Squeeze net



# U-net



# Exercise

Suggest an architecture for a face recognition security system:

- ❑ system should be able to grant access to any person with sufficient rights.

Describe:

- ❑ data required;
- ❑ function of the neural network (classification, regression, clusterisation);
- ❑ architecture of the network;
- ❑ training procedure.

# Summary



# Summary

- ❖ network architecture plays crucial role in Deep Learning;
- ❖ additional problems may provide additional information about solution;
- ❖ there are tons of various network architectures.

# References

- ❖ Ronneberger O, Fischer P, Brox T. U-net: Convolutional networks for biomedical image segmentation. In International Conference on Medical Image Computing and Computer-Assisted Intervention 2015 Oct 5 (pp. 234-241). Springer, Cham.
- ❖ Szegedy C, Ioffe S, Vanhoucke V, Alemi AA. Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning. In AAAI 2017 (pp. 4278-4284).
- ❖ Srivastava RK, Greff K, Schmidhuber J. Highway networks. arXiv preprint arXiv:1505.00387. 2015 May 3.
- ❖ Rasmus A, Berglund M, Honkala M, Valpola H, Raiko T. Semi-supervised learning with ladder networks. In Advances in Neural Information Processing Systems 2015 (pp. 3546-3554).
- ❖ Lee CY, Xie S, Gallagher P, Zhang Z, Tu Z. Deeply-supervised nets. In Artificial Intelligence and Statistics 2015 Feb 21 (pp. 562-570).
- ❖ Goodfellow IJ, Warde-Farley D, Mirza M, Courville A, Bengio Y. Maxout networks. arXiv preprint arXiv:1302.4389. 2013 Feb 18.