



# Convolutional Neural Network

Shadi Albarqouni

Postdoctoral Researcher

shadi.albarqouni@tum.de



# Outline

- 1 Introduction
- 2 Network Architecture
- 3 Training ConvNets
- 4 What we learned?
- 5 ConvNets Debugging
- 6 Transfer Learning
- 7 ConvNets Successes





# Motivation



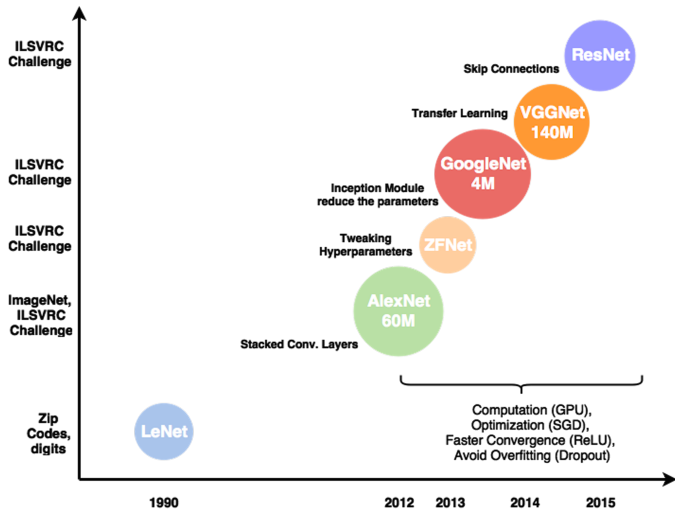
early separable feature Space?

# Object Recognition: Pipeline

Hierarchical and Non-linear feature representation (stacked layers) learned jointly with the classifier



# ConvNets Successes

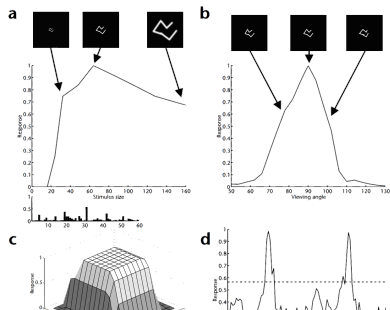
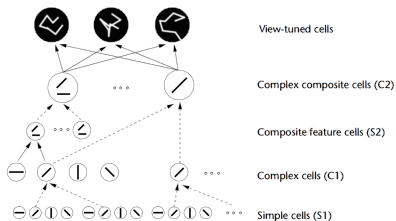


# What is ConvNet?

## Definition (ConvNet)

It is a member of Deep Learning family. It is similar to Artificial Neural Networks (ANN), however, the connectivity pattern between its neuron is inspired by the hierarchical organization of animal visual cortex<sup>a</sup>.

<sup>a</sup>riesenhuber1999hierarchical.



# What's wrong with ANN? (1)

- Hard to Train (over-fitting)
- Careful Initialization
- Huge number of parameters

## Key ideas of ConvNets

- image statistics (shared weights)
- Low-level features supposed to be local (local connectivity)
- High-level features supposed to be coarser (subsampling)

"Convolution + Activation + Pooling = Architecture"



# Network Architecture (1)

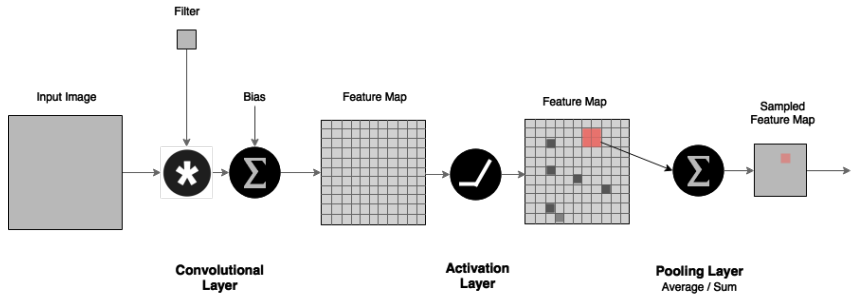


Figure: Symbolic Architecture

Define: receptive field, stride, depth and width of the network.

# Network Architecture (2)



# Notation

We follow the notations appeared in MatConvNet<sup>1</sup>,

- $X$  is the input data,  $X = \{x_1, x_2, \dots, x_N\} \in \mathbb{R}^{H \times W \times D \times N}$ .
- $N$  is the number of input instances/samples.
- $H$  is the height of an image  $x_{i \in N}$ .
- $W$  is the width of an image  $x_{i \in N}$ .
- $D$  is the channels/depth of an image/volume  $x_{i \in N}$ .
- $Y$  is the desired output,  $Y = \{y_1, y_2, \dots, y_N\} \in \mathbb{R}^{c \times N}$

## Objective

Build a model  $f$  that for a given input  $x$  can predict the output  $\hat{y}$ :

$$\hat{y} = f(x; \omega),$$

where  $\omega$  is the model parameter.



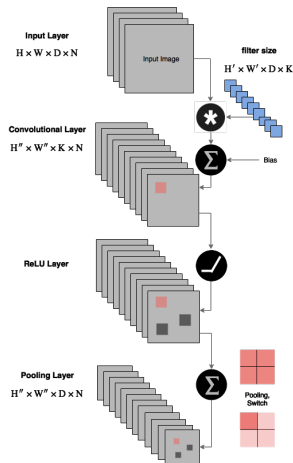
<sup>1</sup>vedaldi2015matconvnet.



# CNN Layers

A CNN Network can be obtained by cascading several layers in a directed acyclic graph (DAG).

- Input Layer
- Convolutional Layer
- Activation Layer
- Pooling Layer
- Fully Connected Layer
- Dropout Layer
- Output Layer



# Input Layer ( $H \times W \times D \times N$ ) (1)

- Data Preprocessing (Mean subtraction, PCA/Whitening)
- Data Augmentation: geometric transformation; rotation and translation, color transformation: illumination, staining ...etc, adding noise.
- Splitting the dataset (training, validation and testing)
- Batch size



# Input Layer ( $H \times W \times D \times N$ )

- 2D inputs
  - Gray ( $D = 1$ )
  - RGB ( $D = 3$ )<sup>2</sup>
- 2.5D inputs
  - Gray ( $D = 3$ )<sup>3</sup>
  - RGB-D ( $D = 4$ )<sup>4</sup>
- 3D inputs
  - Gray ( $D = d$ )<sup>5</sup>

---

<sup>2</sup>eigen2014predicting.

<sup>3</sup>roth2014new.

<sup>4</sup>gupta2014learning.

<sup>5</sup>kamnitsas2015multi.



# Convolutional Layer ( $H'' \times W'' \times K \times N$ ) (1)

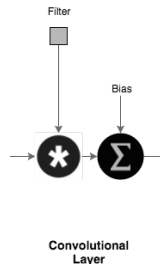


# Convolutional Layer ( $H'' \times W'' \times K \times N$ ) (2)

It computes the convolution of input image  $x$  with a filter  $f$  as follows

$$y_{i,j,k} = b_{i,j,k} + \sum_{h=1}^{H'} \sum_{w=1}^{W'} \sum_{d=1}^D f_{h,w,d,k} \cdot x_{i+h,j+w,d},$$

- input  $x \in \mathbb{R}^{H \times W \times D}$
- filters  $f \in \mathbb{R}^{H' \times W' \times D \times K}$
- biases  $b \in \mathbb{R}^{H'' \times W'' \times K}$
- output  $y \in \mathbb{R}^{H'' \times W'' \times K}$
- stride  $S_{W,H}$  and padding  $P_{W,H}$ ,



# Convolutional Layer ( $H'' \times W'' \times K \times N$ ) (3)



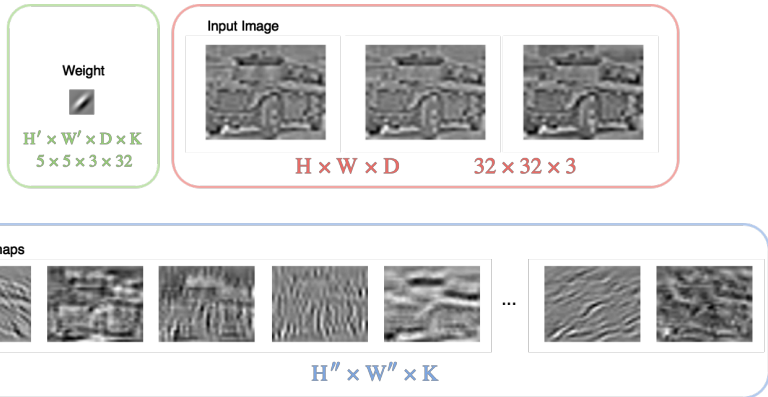
# Convolutional Layer ( $H'' \times W'' \times K \times N$ ) (4)

Example: CIFAR-10 (Convolution,  $5 \times 5 \times 3 \times 32$ )

Keywords: Translation Invariance, few parameters, local consistency



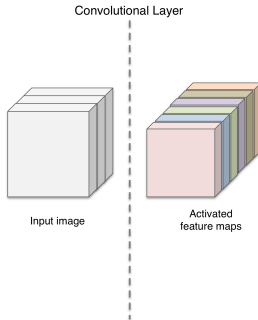
# Convolutional Layer ( $H'' \times W'' \times K \times N$ ) (5)



$$W'' = 1 + \frac{W - W' + (P_L + P_R)}{S_W}, \quad H'' = 1 + \frac{H - H' + (P_U + P_D)}{S_H},$$



# Activation Layer ( $H \times W \times D \times N$ ) (1)

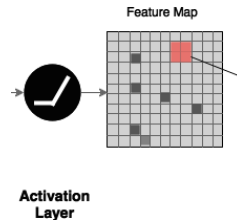


# Activation Layer ( $H \times W \times D \times N$ ) (2)

It computes the **Rectified Linear Unit (ReLU)** of each feature channel  $x$  as follows

$$y_{i,j,d} = \max\{0, x_{i,j,d}\},$$

- input  $x \in \mathbb{R}^{H \times W \times D}$
- output  $y \in \mathbb{R}^{H \times W \times D}$



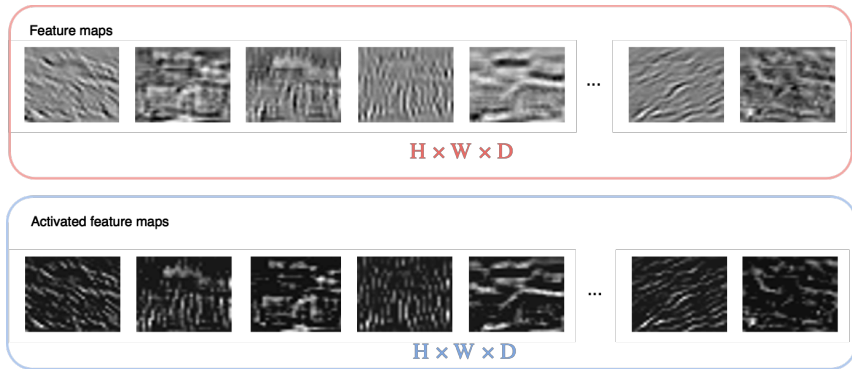
## Activation Layer ( $H \times W \times D \times N$ ) (3)



# Activation Layer ( $H \times W \times D \times N$ ) (4)

Example: CIFAR-10 (ReLU)

Keywords: Simplifies Back-propagation, Makes Learning faster.



# Activation Layer ( $H \times W \times D \times N$ ) (5)

What about other activation functions? Any potential drawbacks?

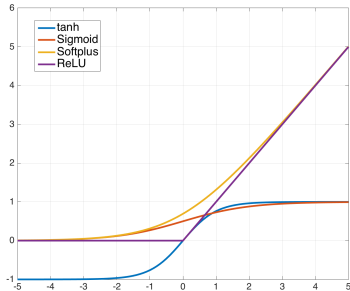


Figure: Activation functions

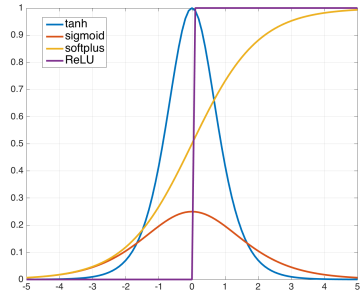


Figure: Activation derivatives

# Pooling Layer ( $H'' \times W'' \times D \times N$ ) (1)



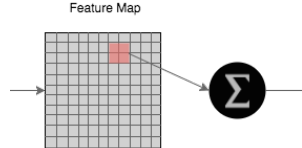
# Pooling Layer ( $H'' \times W'' \times D \times N$ ) (2)

It computes the maximum or average response of each feature channel  $x$  within a 2D patch  $p$  as follows

$$y_{i,j,d} = \max_{1 \leq h \leq H', 1 \leq w \leq W'} x_{i+h, j+w, d},$$

$$y_{i,j,d} = \frac{1}{H'W'} \sum_{1 \leq h \leq H', 1 \leq w \leq W'} x_{i+h, j+w, d},$$

- input  $x \in \mathbb{R}^{H \times W \times D}$
- patch  $p \in \mathbb{R}^{H' \times W'}$
- output  $y \in \mathbb{R}^{H'' \times W'' \times D}$
- stride  $S_{W,H}$  and padding  $P_{W,H}$



**Pooling Layer**  
Average / Sum

## Pooling Layer ( $H'' \times W'' \times D \times N$ ) (3)





# Pooling Layer ( $H'' \times W'' \times D \times N$ ) (4)

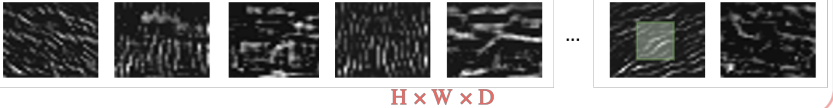
Example: CIFAR-10 (Max. Pooling,  $p = 3 \times 3, S = 2$ )

Keywords: Invariance to small transformation, Larger receptive field

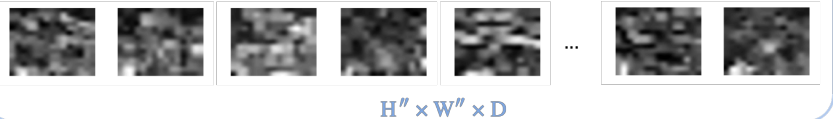


# Pooling Layer ( $H'' \times W'' \times D \times N$ ) (5)

Activated feature maps

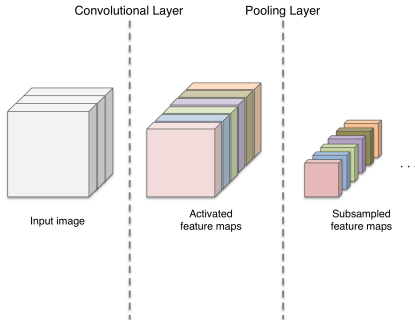


Subsampled feature maps



$$W'' = 1 + \frac{W - W' + (P_L + P_R)}{S_W}, \quad H'' = 1 + \frac{H - H' + (P_U + P_D)}{S_H},$$

# Normalization Layer ( $H \times W \times D \times N$ ) (1)



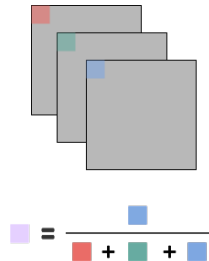
## Normalization Layer ( $H \times W \times D \times N$ ) (2)

It performs a cross-channel normalization at each spatial location as follows

$$y_{i,j,d} = x_{i,j,d} \left( \kappa + \alpha \sum_{d \in D} x_{i,j,d}^2 \right)^{-\beta},$$

where  $\kappa, \alpha, \beta$  are hyperparameters. It is usually called Local Response Normalization (LRN).

- input  $x \in \mathbb{R}^{H \times W \times D}$
- output  $y \in \mathbb{R}^{H \times W \times D}$



## Normalization Layer ( $H \times W \times D \times N$ ) (3)



# Normalization Layer ( $H \times W \times D \times N$ ) (4)

Example: CIFAR-10 (LRN,  $\kappa = 0, \alpha, \beta = 1$ )

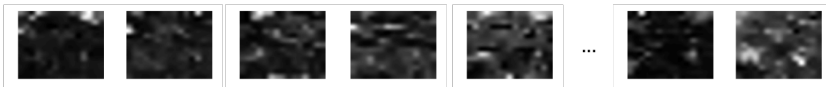
Keywords: Within or Cross feature maps, Before or After Pooling, Have you spotted the mistake in the normalization process?

Subsampled feature maps



$H \times W \times D$

Normalized feature maps



$H \times W \times D$

# Fully Connected Layer ( $1 \times 1 \times K \times N$ ) (1)



## Fully Connected Layer ( $1 \times 1 \times K \times N$ ) (2)

It computes the convolution of input feature maps  $x$  with a filter  $f$  as follows

$$y_{i,j,k} = b_{i,j,k} + \sum_{h=1}^H \sum_{w=1}^W \sum_{d=1}^D f_{h,w,d,k} \cdot x_{i+h,j+w,d},$$

- input  $x \in \mathbb{R}^{H \times W \times D}$
- filters  $f \in \mathbb{R}^{H \times W \times D \times K}$ , we use  $K$  such filters.
- biases  $b \in \mathbb{R}^{1 \times 1 \times K}$
- output  $y \in \mathbb{R}^{1 \times 1 \times K}$
- stride and padding



# Fully Connected Layer ( $1 \times 1 \times K \times N$ ) (3)





# Objective function

What we have presented so far is the feed-forward propagation, however, to minimize our objective function, we need to propagate back the gradients and update the parameters.

The Objective function:

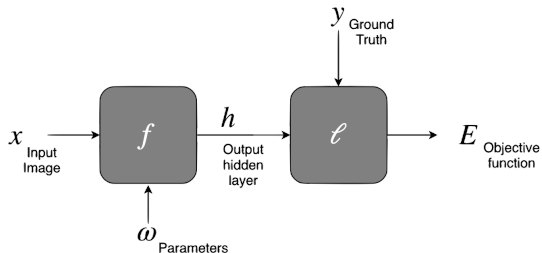
$$\operatorname{argmin}_{\omega_1, \dots, \omega_L} \frac{1}{n} \sum_{i=1}^n \ell(f(x^{(i)}; \omega_1, \dots, \omega_L), y^{(i)})$$

where  $f(x; \omega)$  is the model's output.

**Solver:** Stochastic Gradient Descent (SGD).



# Optimization and Derivatives (1)



Using the chain rule, the partial derivatives can be written as follows:

$$\frac{\partial E}{\partial x} = \frac{\partial E}{\partial h} \frac{\partial h}{\partial x}, \frac{\partial E}{\partial \omega} = \frac{\partial E}{\partial h} \frac{\partial h}{\partial \omega}$$

# Optimization and Derivatives (2)

**Vanilla update.** The weight's update:

$$\omega^{t+1} = \omega^t - \frac{\eta}{n} \sum_{i=1}^n \nabla \ell(x, y; \omega^t),$$

where  $\eta$  is the learning rate.

**Momentum update.** Using the momentum<sup>6</sup>, The weight's update becomes:

$$\omega^{t+1} = \omega^t - \frac{\eta}{n} \sum_{i=1}^n \nabla \ell(x, y; \omega^t) + \alpha \nabla \omega^t,$$

where  $\alpha$  is the momentum.



# Optimization and Derivatives (3)

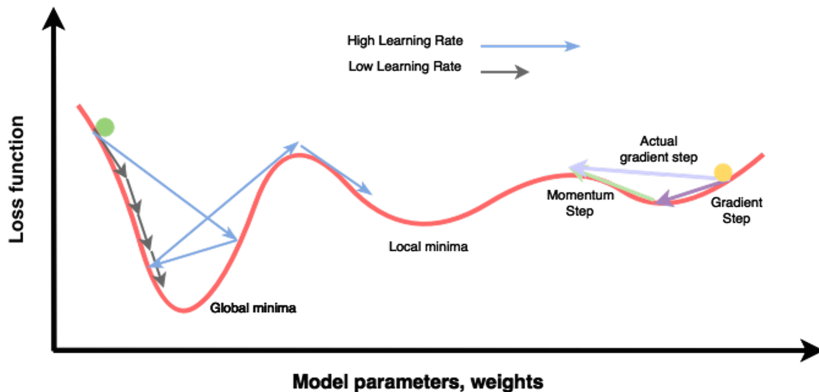


Figure: SGD & Learning Rate<sup>7</sup>

<sup>6</sup>rumelhart1988learning.

<sup>7</sup><http://imgur.com/a/Hqolp>

## Loss Layer ( $1 \times 1 \times C \times N$ ) (1)

The loss function  $\ell$ , mainly used in the training phase, could be softmax log-loss for "classification purpose"

$$y = - \sum_{i,j} \left( x_{i,j,c} - \log \sum_{d=1}^D \exp x_{i,j,d} \right),$$

or  $\ell_2$ -norm for "regression purpose" as follows

$$y = \|x_{i,j,c} - x_{i,j,d}\|_2^2,$$

## Loss Layer ( $1 \times 1 \times C \times N$ ) (2)





# Recap



# Low/Mid/High Level Features (1)

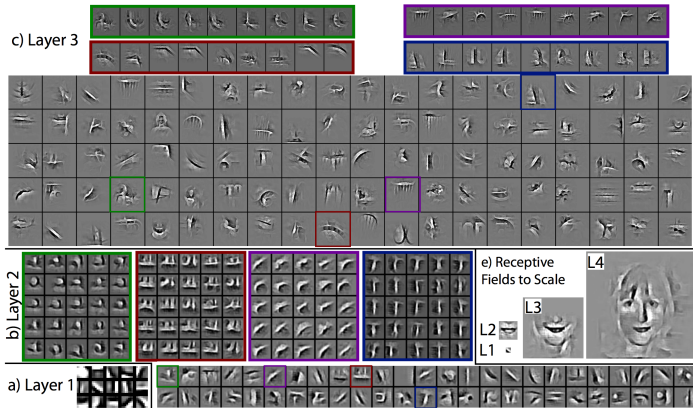


Figure: Low and Mid Level Features, Fig.5 in<sup>8</sup>

# Low/Mid/High Level Features (2)

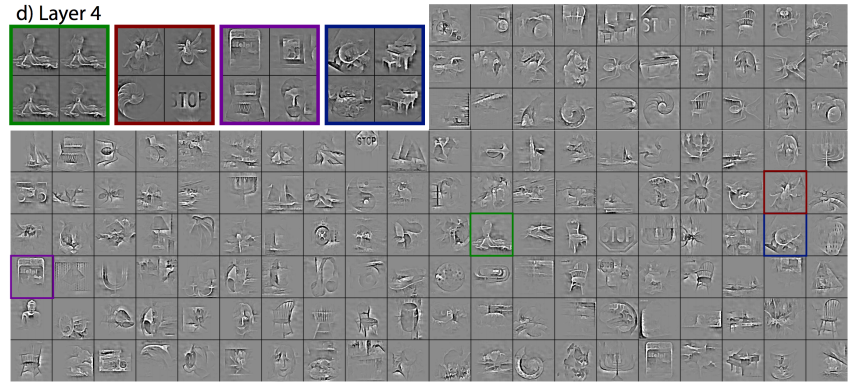


Figure: High Level Features, Fig.5

# Interactive Example

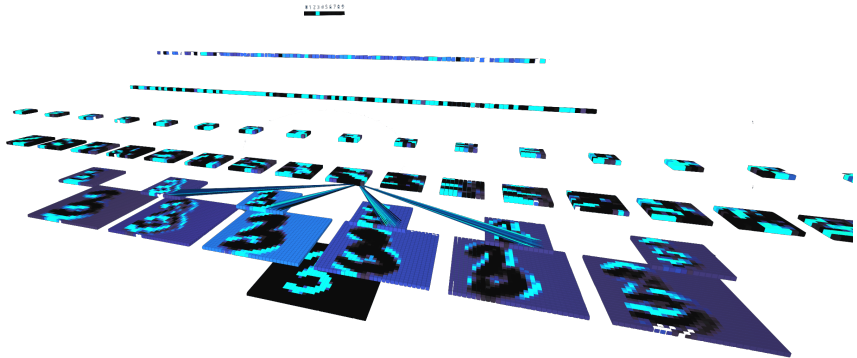
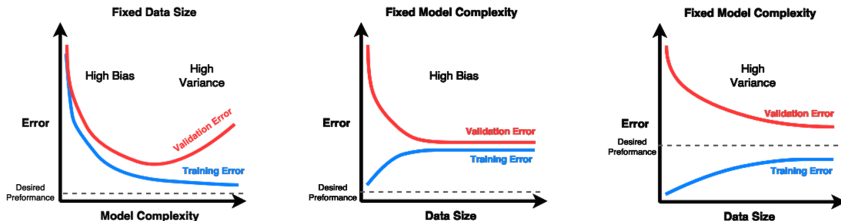


Figure: LeNet5 Architecture, MNIST-10<sup>9</sup>

# Network Training

**Model Check.** Similar to any model-based machine learning, there are two types of error source; 1) Bias and 2) Variance.



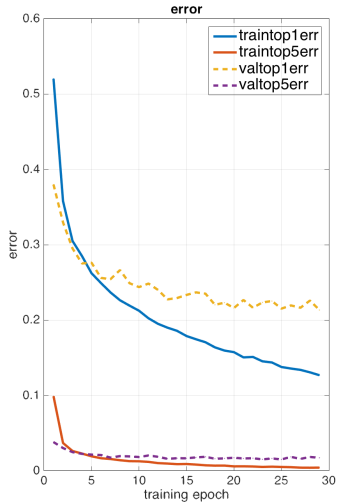
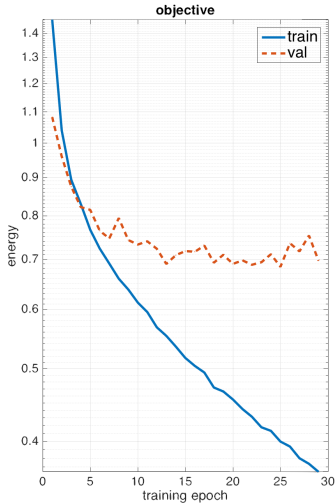
How to fix High Bias? High Variance?<sup>10</sup>

High Variance: Getting more training data (data augmentation), smaller set of features, increase regularization parameter, add more dropout.

High Bias: Getting larger set of features, deeper architecture.

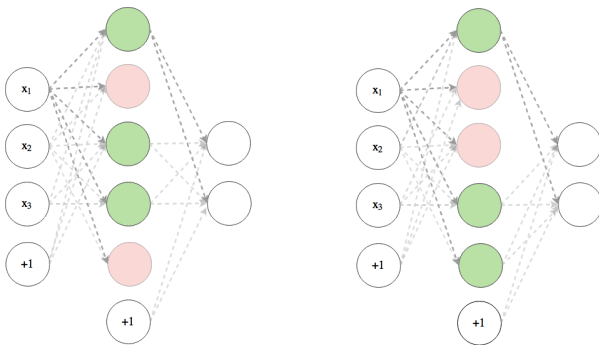
# Network Training

Example: Monitoring the training of tiny VGG model (30 Epochs)



# Dropout Layer $(1 \times 1 \times C \times N)^{11}$

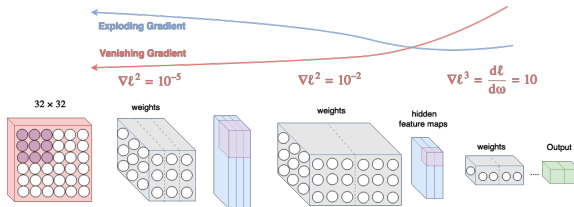
The dropout layer acts as a regularizer for the network to avoid overfitting. It is simply "dropping out" some activation units and setting them to zero during the training phase. It is similar to train thinner networks and do averaging.



# Network Debugging

## Gradient Checks

- One of the major problems with training a CNN deep model is vanishing/exploding gradient<sup>12</sup>.
- Monitor gradient and activation across layers and epochs.
- Try adding Batch Normalization layer, proper weight initialization<sup>13</sup>.



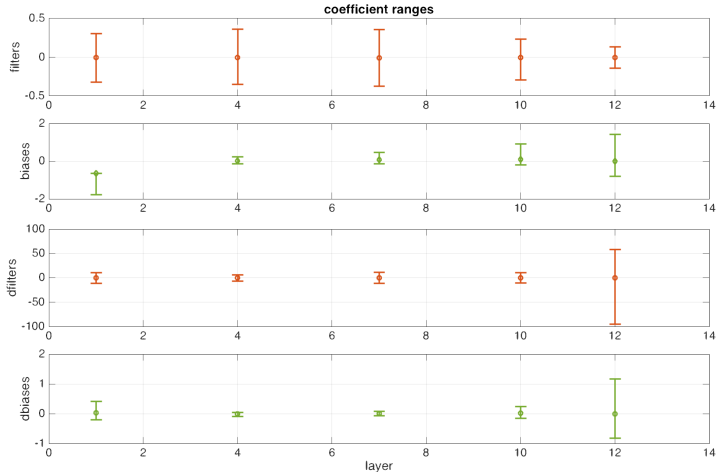
<sup>12</sup>bengio1994learning.

<sup>13</sup>krahenbuhl2015data.



# Network Debugging

Example: Monitoring the gradient of tiny VGG model (Epoch 26)

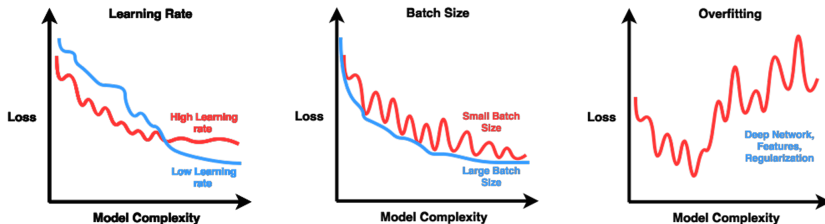


# Network Debugging

## Sanity Checks

- Check if you have an expected loss value (*Hint: Set the regularization parameter to Zero.*)
- Increasing the regularization parameter will increase the loss.
- Overfit a very small subset of data.

## Loss Checks



# Additional Layers

- Deconvolutional Layer<sup>14</sup>
- Batch Normalization<sup>15</sup>
- DropConnect<sup>16</sup>

---

<sup>14</sup>zeiler2011adaptive; zeiler2014visualizing.

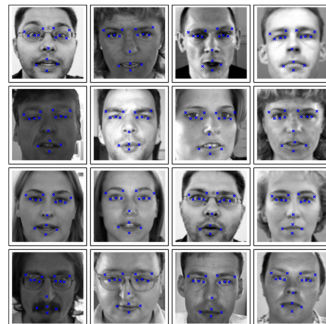
<sup>15</sup>ioffe2015batch.

<sup>16</sup>wan2013regularization.



# Example: Facial keypoints tutorial

- Dataset: Facial Keypoint Detection challenge, Training: 7049 ( $96 \times 96$ ) gray images with 15 keypoints. Testing: 1783 images.
- Loss function: Regression (MSE)
- Parameters: Optimization: nesterov momentum, Learning rate: 0.01, Momentum = 0.9.



Note: Image Courtesy of this example at<sup>17</sup>, Facial keypoint challenge<sup>18</sup>.

---

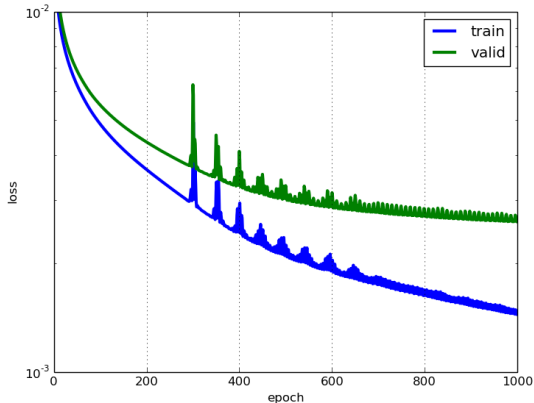
<sup>17</sup>dnaiel\_KP2014.

<sup>18</sup>Kaggle2013.

# Example: Facial keypoints tutorial (Cont.) (1)

One layer network (net1)

- Network: One hidden layer, (9216, 100, 30) units.
- Parameters: Number of Epochs = 400.



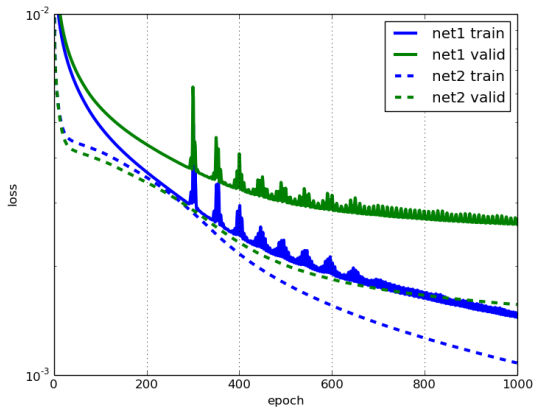
# Example: Facial keypoints tutorial (Cont.) (2)

LeNet5 network (net2)

- Network: Input, (Conv, maxPool)<sup>3</sup> + FC<sup>2</sup>, Output
- Parameters: Number of Epochs = 1000.



## Example: Facial keypoints tutorial (Cont.) (3)

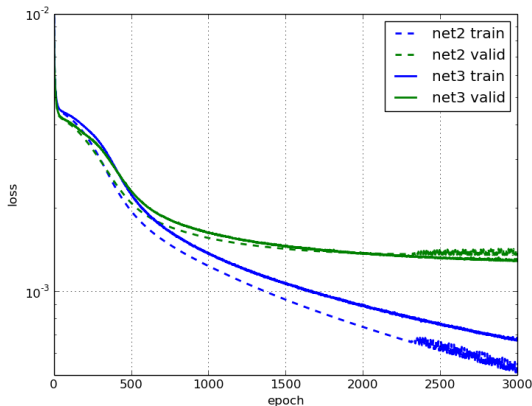


LeNet5 network (net3)



## Example: Facial keypoints tutorial (Cont.) (4)

- Data Augmentation, only flipping 50% of datasets.
- Parameters: Number of Epochs = 3000.





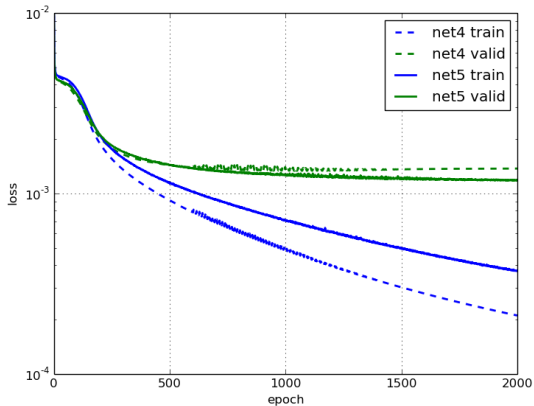
# Example: Facial keypoints tutorial (Cont.) (5)

LeNet5 network (net4, net5)

- Parameters: Learning Rate = 0.03-0.0001, Momentum = 0.9 - 0.999
- with/without Data Augmentation



## Example: Facial keypoints tutorial (Cont.) (6)



# Transfer Learning

**Learning from scratch.** Inspired by some CNN architecture, you can design your own network. However, you need tons of data.

**Transfer Learning**<sup>19</sup>. Once you don't have enough data, you can use the pre-trained CNN models<sup>20</sup> for the following tasks:

- Extract features: The output of the last hidden layer before the softmax can be used as features (CNN Codes) to train a linear SVM classifier.
- Fine-tuning: You may need to propagate back your gradient to update the weights, however, the weights of the first layers can be fixed during the fine-tuning and update the weights of the higher layers.

---

<sup>19</sup>li2015cs231n.

<sup>20</sup>Caffe: <https://github.com/BVLC/caffe/wiki/Model-Zoo>  
MatConvNet: <http://www.vlfeat.org/matconvnet/pretrained/>



# Fine-Tuning Tricks



Pre-trained CNN model

(a)



Pre-trained CNN model

(b)



Pre-trained CNN model

(c)



Pre-trained CNN model

(d)

(a) Fine-tuning, (b) Train from scratch, initialize the weights of the first layers from a pre-trained model, (c) Get the CNN codes and learn a linear SVM, (d) Get the CNN codes from the mid-layers and learn a linear SVM.



# Hyper-parameters: additional topics

- Optimization solver<sup>21</sup>.
- Learning Rate Schedule<sup>22</sup>: The more intuitive way to choose the learning rate is to set it high in the beginning (large step and faster), and then lower it down after some epochs (small step and slower), i.e.  $\eta = \frac{\eta_0}{n_{iter} + \kappa}$  or  $\eta = \eta_0 e^{-\kappa n_{iter}}$ .
- Momentum<sup>23</sup>
- Batch Size: between 10 and few hundreds.
- Weight Initialization<sup>24</sup>.

---

<sup>21</sup>bottou2012stochastic; ngiam2011optimization.

<sup>22</sup>bengio2012practical.

<sup>23</sup>sutskever2013importance.

<sup>24</sup>wagner2013learning.



# Parameters exploding!

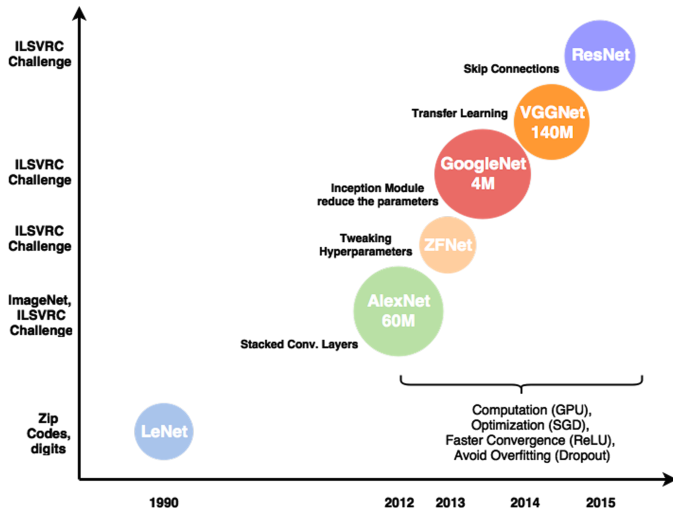
Let's compute how many parameters we have in LeNet5 network (c.f. Fig. ??):

(Conv. + maxPool.) <sup>3</sup>		
conv5-6	$1 \times 5 \times 5 \times 6 + 6$	156
conv5-16	$6 \times 5 \times 5 \times 16 + 16$	2416
conv5-120	$16 \times 5 \times 5 \times 120 + 120$	48120
FC <sup>2</sup>		
conv1-84	$120 \times 1 \times 1 \times 84 + 84$	10164
conv1-10	$84 \times 1 \times 1 \times 10 + 10$	850

61706 parameters, quite good number! Have a look at the number of parameters for the recent networks in the next slide.



# ConvNets Successes



# Case Studies (AlexNet) (1)

Architecture: (Conv + Max. Pooling)<sup>3</sup> + FC<sup>3</sup>

DataSet: ILSVRC-1000, 1.2 million training images (RGB 256<sup>2</sup>), 50,000 validation, and 150,000 testing images.

Pre-Processing: DeMean the training images. Training Param:  
LR = 0.01, M = 0.9, batch = 128





# Case Studies (AlexNet) (2)

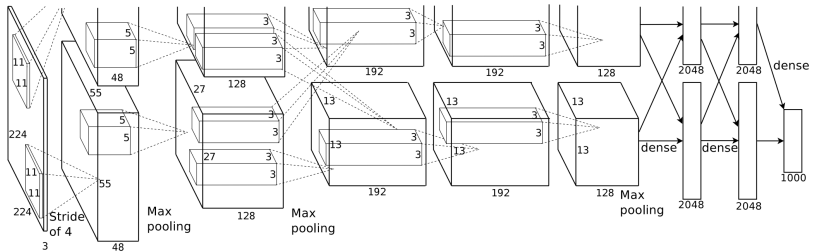


Figure: AlexNet. Fig.2

# Case Studies (AlexNet) (3)

- ReLU converges faster than  $\tanh(\cdot)$ , 6 times faster at 25% error rate.
  - LRN is applied after ReLU, hyper-parameters are determined using validation set ( $\kappa = 2, \alpha = 10^{-4}, \beta = 0.75$ ). Error rate: 13%  $\rightarrow$  11%.
  - Overlapping Pooling. Error rate  $\rightarrow -0.4\%$
  - Data Augmentation including translation, flipping and altering color intensities. Error rate  $\rightarrow -1\%$ .
  - Dropout at rate 0.5 (Over-fitting)
- Note: Images Courtesy of this case study at<sup>25</sup>.

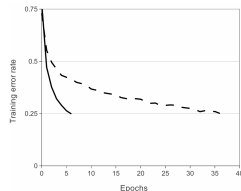


Figure 1: A four-layer convolutional neural network with ReLUs (solid line) reaches a 25% training error rate on CIFAR-10 six times faster than an equivalent network with tanh neurons (dashed line). The learning rates for each network were chosen independently to make training as fast as possible. No regularization of any kind was employed. The magnitude of the effect demonstrated here varies with network architecture, but networks with ReLUs consistently learn several times faster than equivalents with saturating neurons.

Thanks!

# Questions

