



**CORNELL
TECH**

Deep Learning Clinic (DLC)

Lecture 7 - Tricks and Tips

Jin Sun

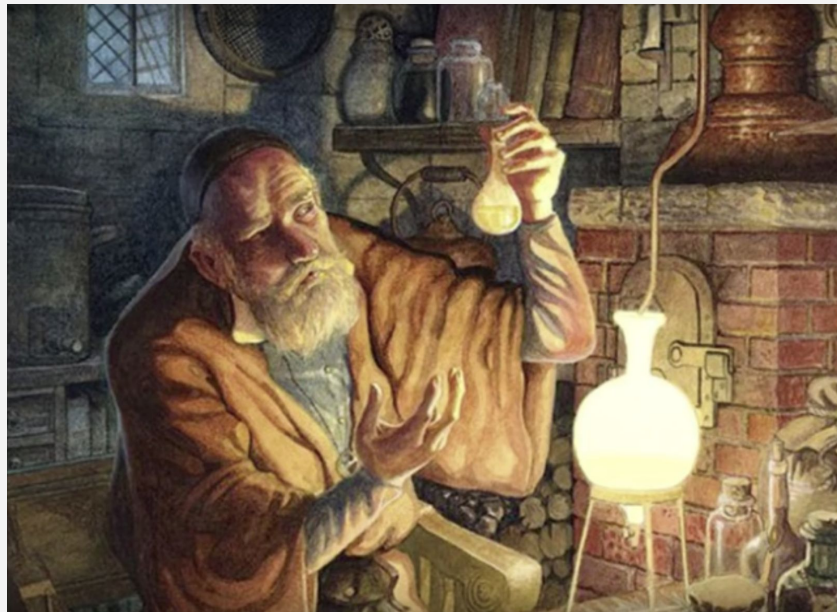
11/9/2018

Today - Tricks and Tips in DL Training

- **Overview**
- Data Preprocessing
- Network Design
 - Activation Functions
 - Weight Initialization
 - Normalization
 - Monitoring the Learning Process
 - Hyperparameters
- Optimization
- Transfer Learning

Overview

Training a neural network is an art!



But there are tricks people found useful.

Overview

Network Design

Activation, preprocessing, initialization, regularization

Training Dynamics

Monitor the changing of train/val loss, parameter updates, hyperparameters

Recap: Training A Network

Training Loop:

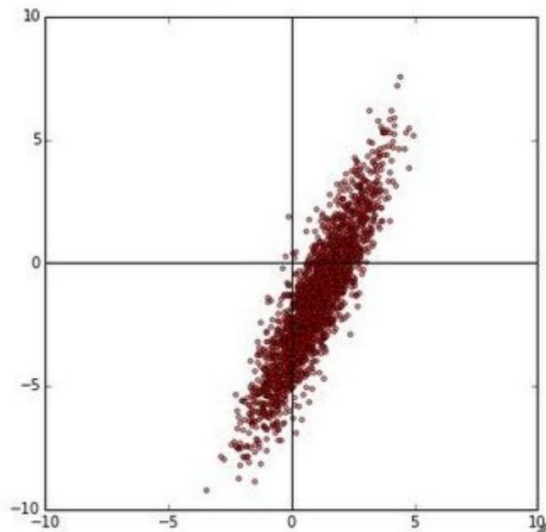
1. Sample a **batch** of data.
2. **Forward** pass the data through the network and calculate the loss.
3. **Backprop** the loss to calculate gradients of the network.
4. **Update** parameters using gradient based optimization method.

Today - Tricks and Tips in DL Training

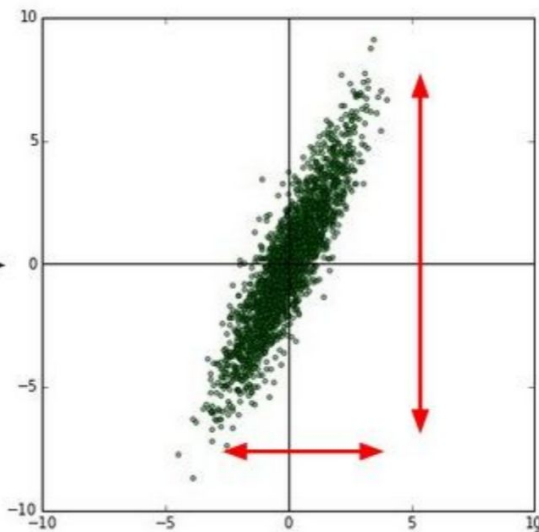
- Overview
- **Data Preprocessing**
- Network Design
 - Activation Functions
 - Weight Initialization
 - Normalization
 - Monitoring the Learning Process
 - Hyperparameters
- Optimization
- Transfer Learning

Data Preprocessing

original data

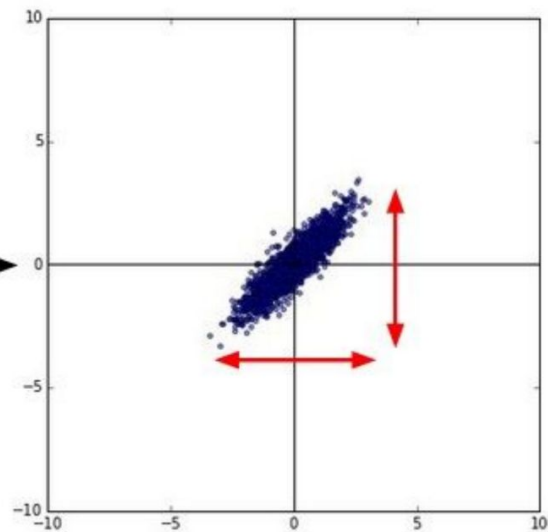


zero-centered data



```
X -= np.mean(X, axis = 0)
```

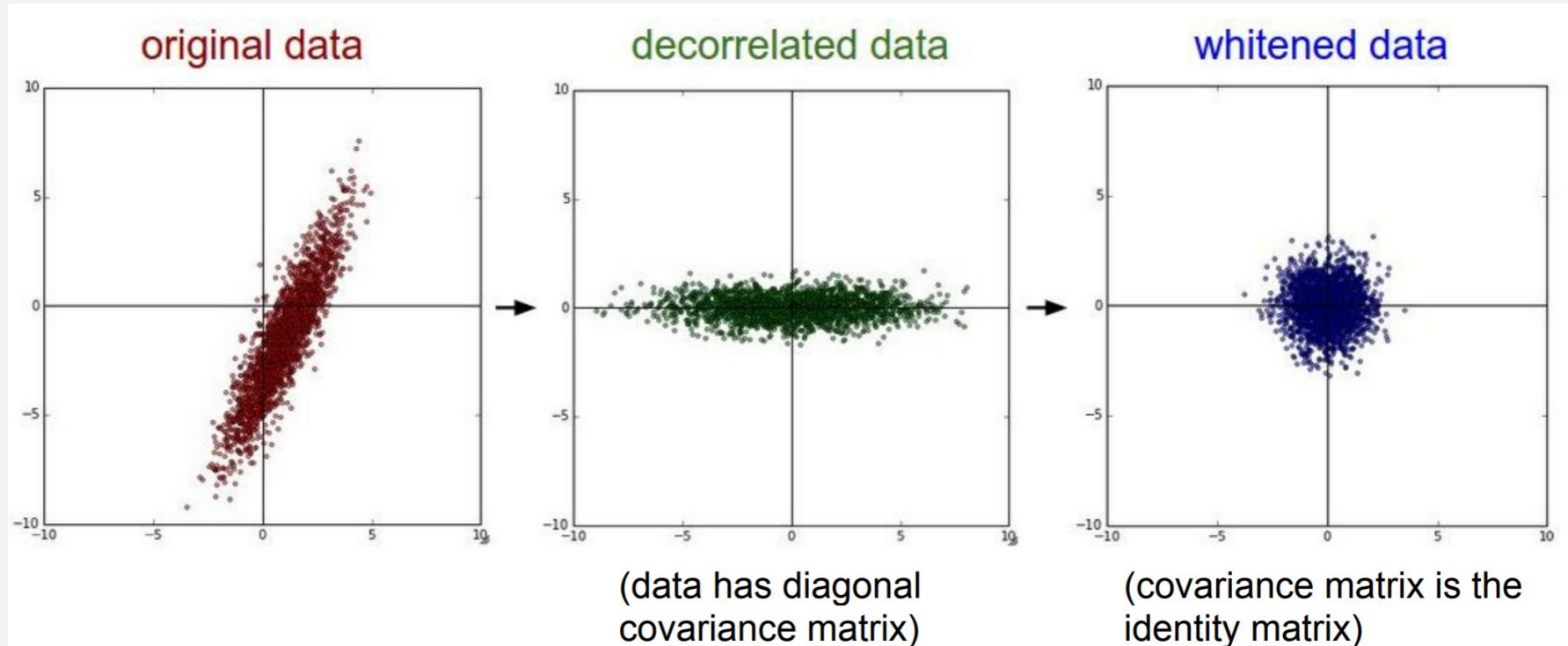
normalized data



```
X /= np.std(X, axis = 0)
```

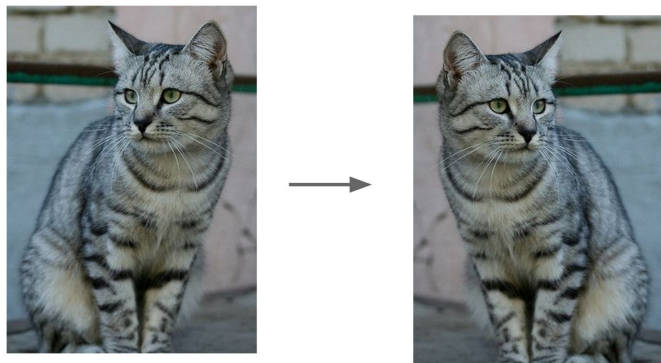
In practice, only do zero-centering for images.

Data Preprocessing - PCA and Whitening

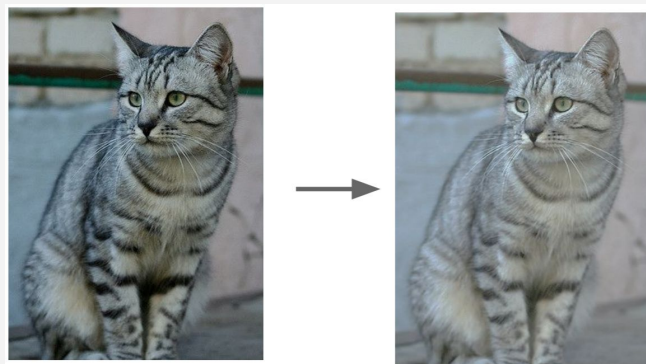


Data Augmentation

Horizontal Flip:

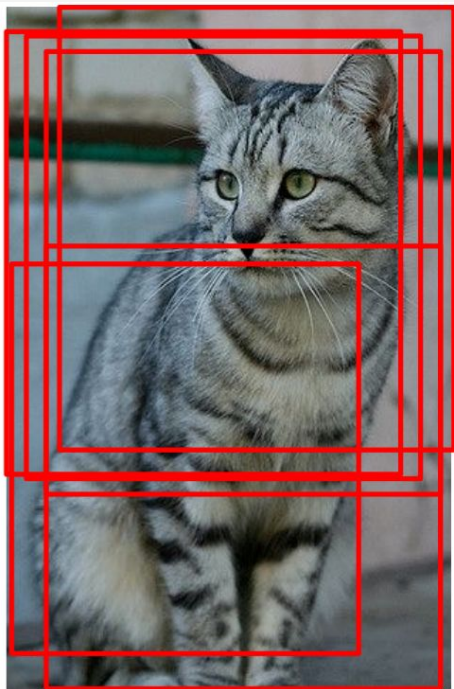


Random Contrast and Brightness

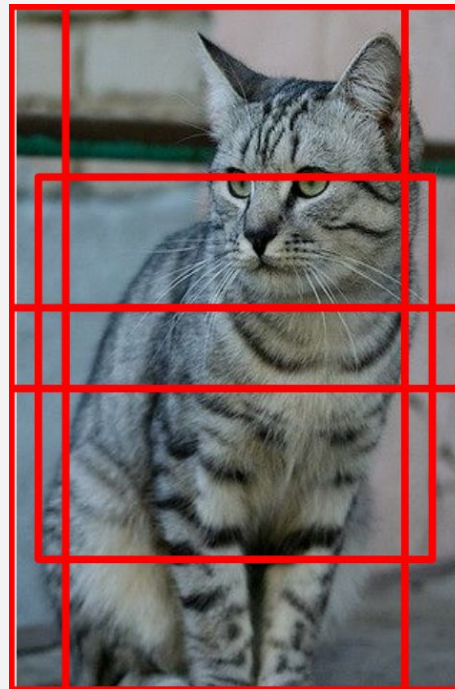


Data Augmentation

Random Crops and Scales (Train)



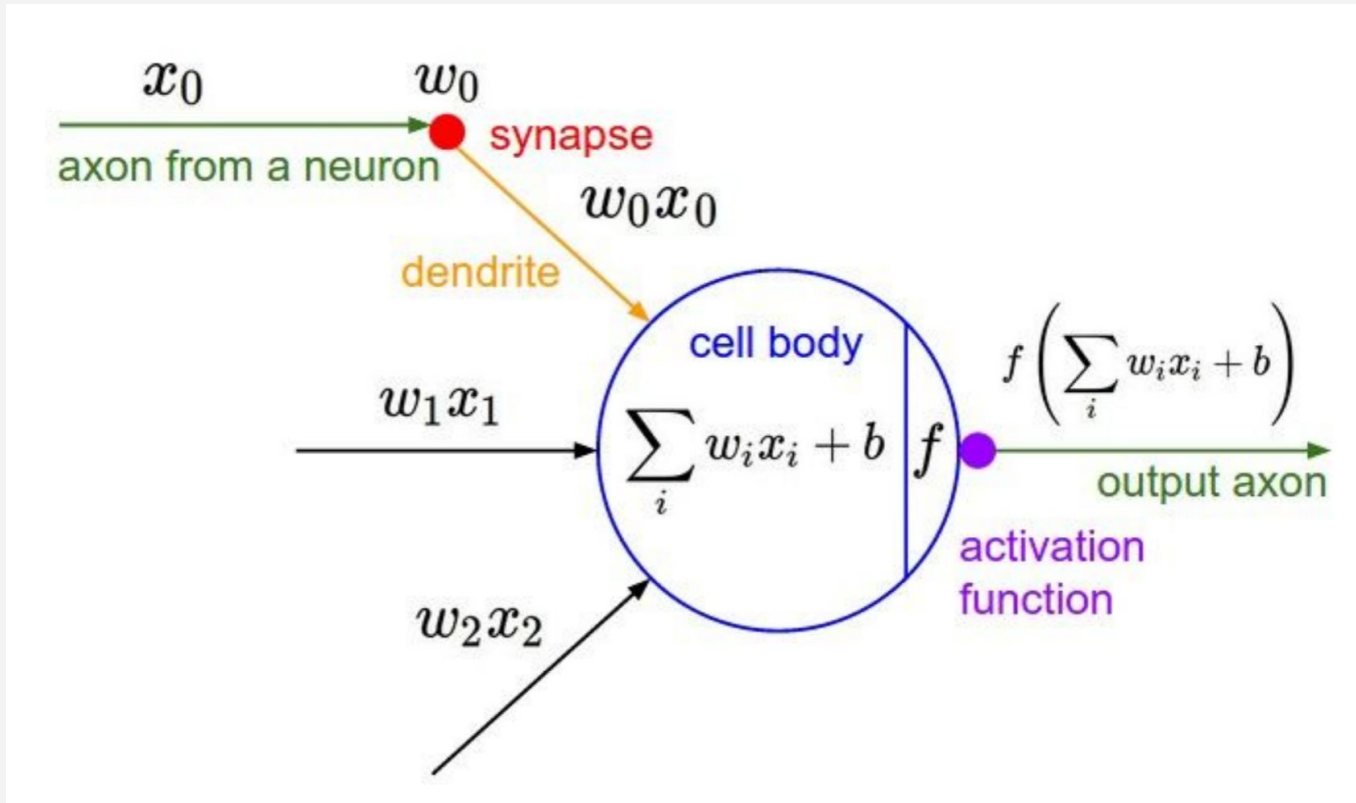
Fixed Crops and Scales (Test)



Today - Tricks and Tips in DL Training

- Overview
- Data Preprocessing
- **Network Design**
 - Activation Functions
 - Weight Initialization
 - Normalization
 - Monitoring the Learning Process
 - Hyperparameters
- Optimization
- Transfer Learning

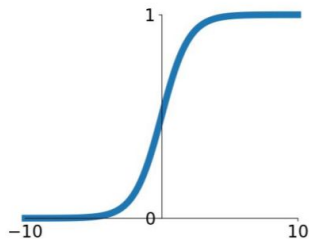
Activation Functions



Activation Functions - the Gallery

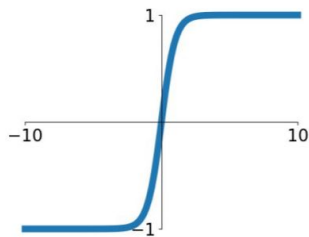
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



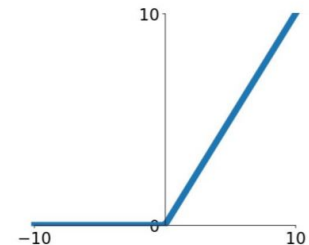
tanh

$$\tanh(x)$$



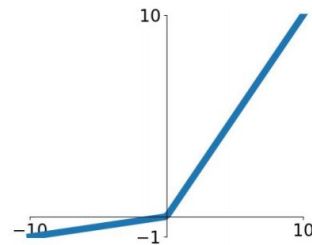
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

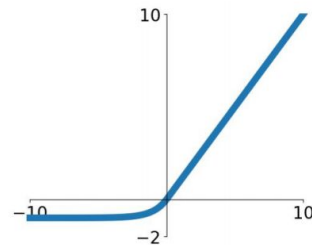


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Activation Functions - Summary

Use **ReLU** as the default

Try Leaky ReLU / Maxout / ELU for alternatives

Try tanh

Sigmoid is usually not recommended.

Weight Initialization

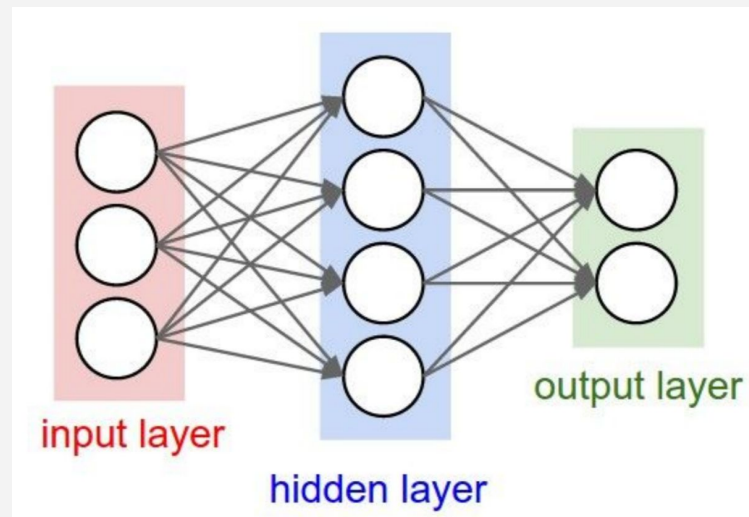
Do not use constant weights.

Small Gaussian random numbers not working well for deep networks..

Recommended:

Xavier initialization, and its variants.

Still an active research area.



Batch Normalization

Idea:

Make a batch of activations to have zero-mean and unit-variance.

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Works as a layer in neural networks.

Put it after fully connected or convolutional layers, before nonlinearity.

Batch Normalization

Benefit:

Improve gradient flow in the network.

Allows higher learning rate.

Works also like a regularizer.

During Testing:

Use a single fixed mean/std obtained from training.

Monitoring the Learning Process

Some simple things to try and observe:

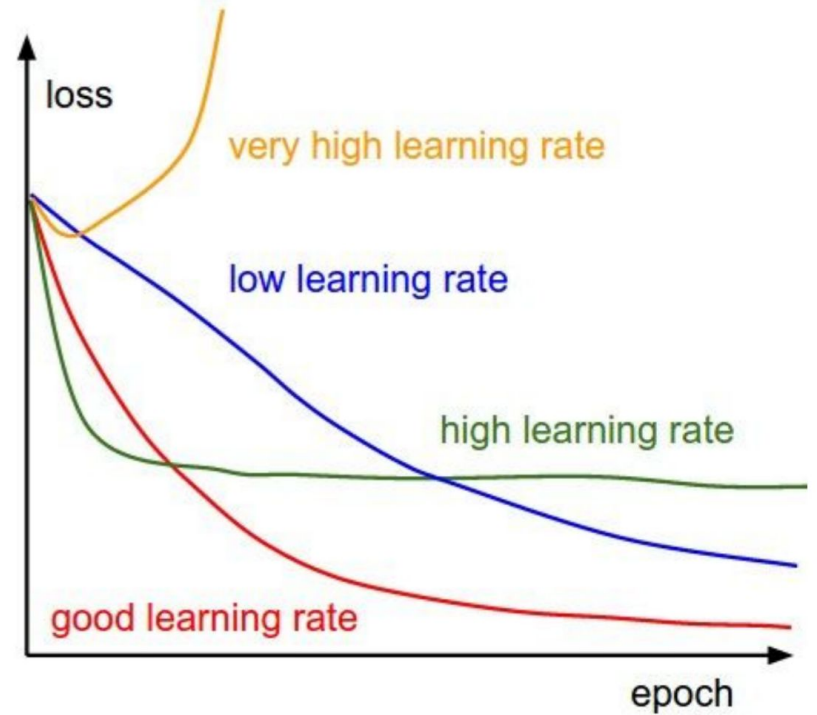
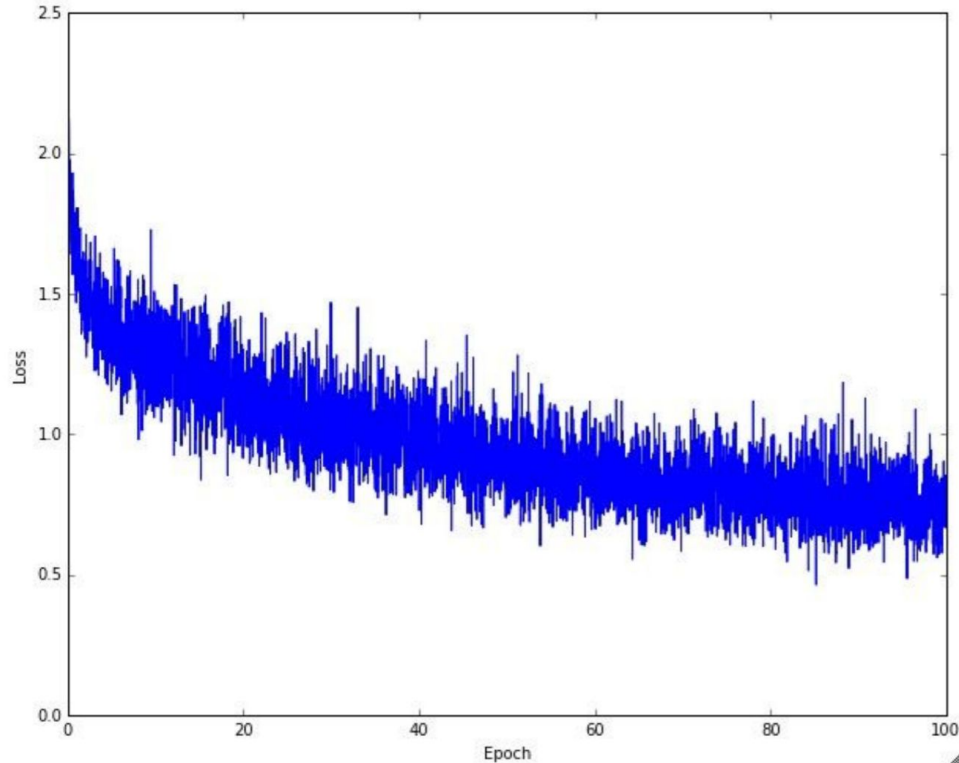
You should be able to overfit your model to a small portion of the data. Do this as a sanity check to make sure the pipeline works.

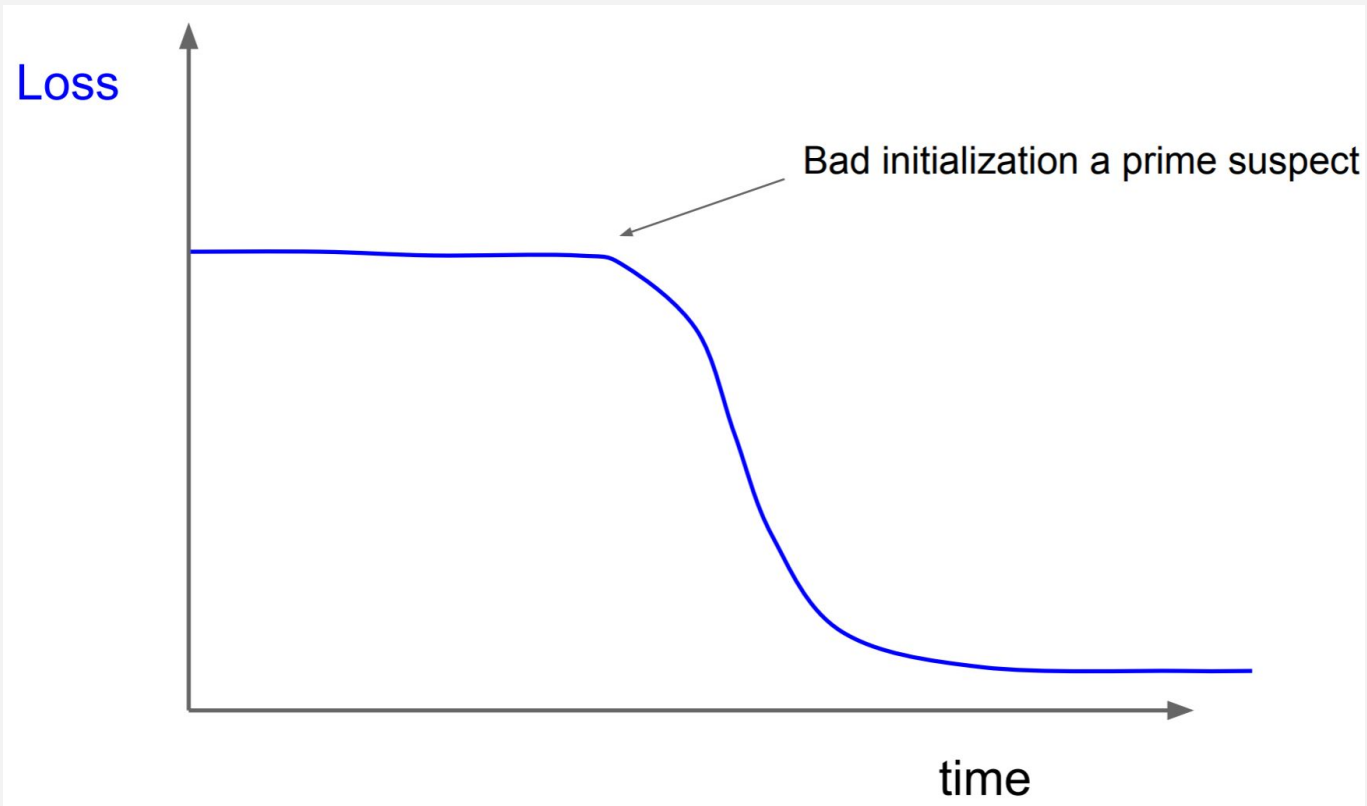
Train loss is not going down: try higher learning rate

Train loss explodes: try lower learning rate

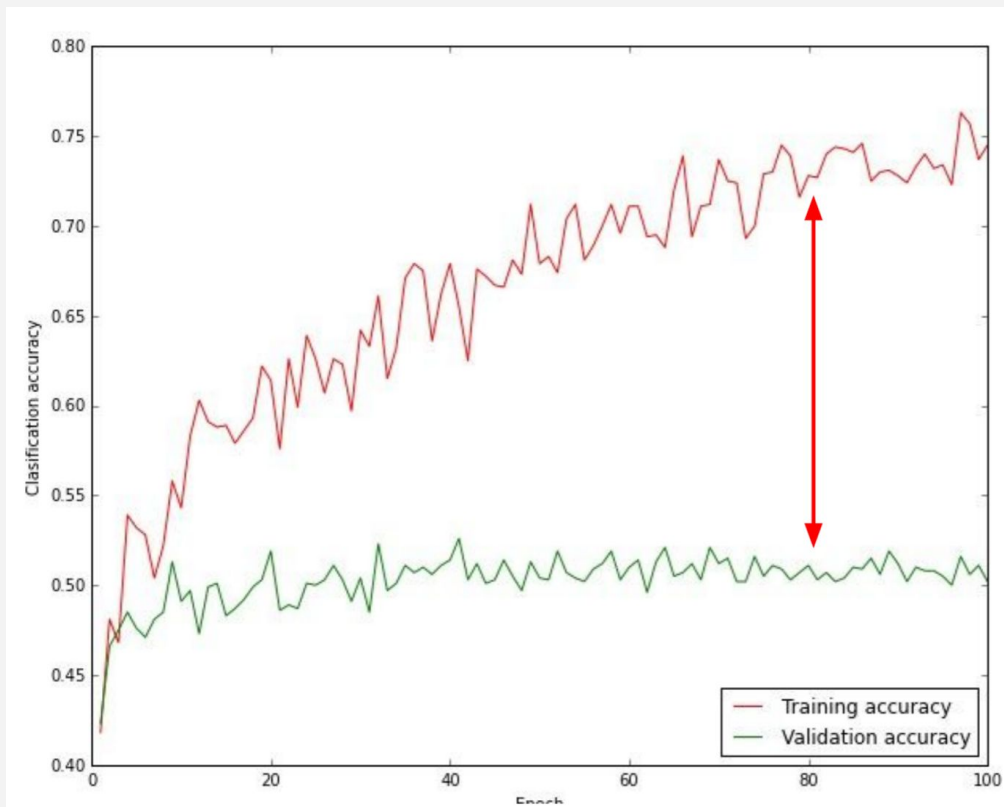
Find a range of learning rate by doing a simple set of experiments.

Learning Rate vs Curves





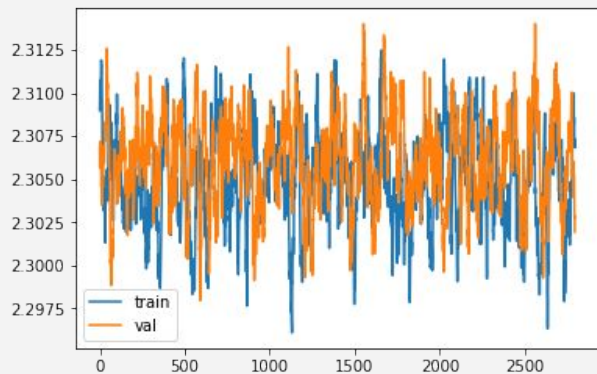
Sign of Overfitting



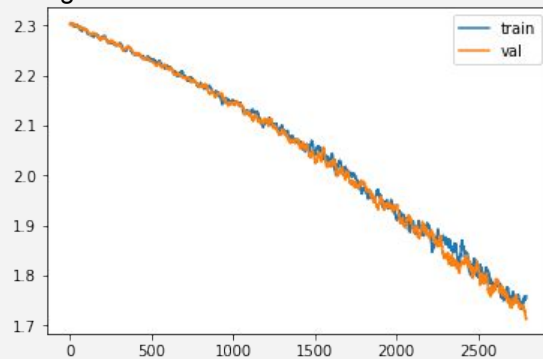
Big gap: Overfitting

No gap: Maybe model can be more complex

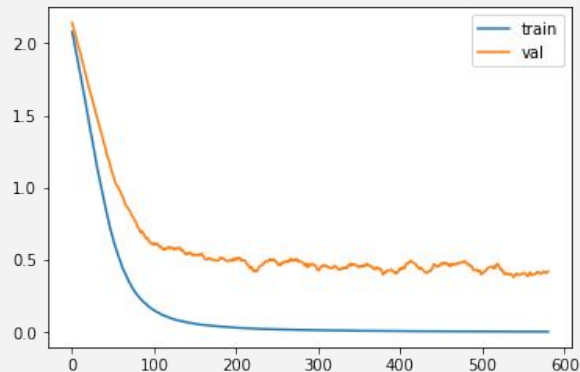
Early Stopping



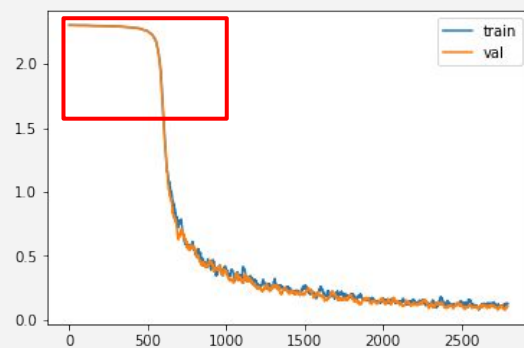
Not learning: gradients not applied to weights



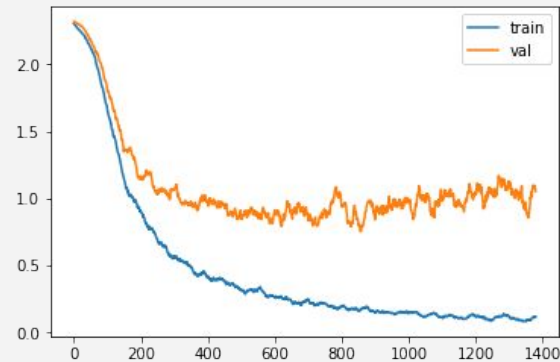
Not converged yet: need longer training



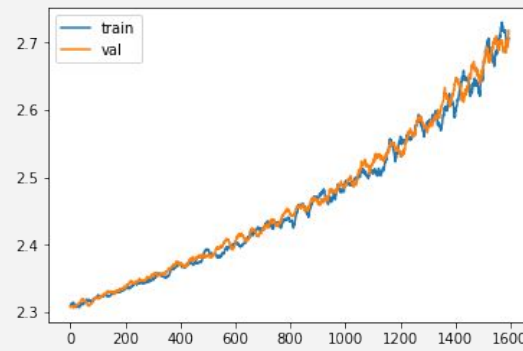
Overfit: model too large/dataset too small
slow start



Slow start: initialization weights too small



More extreme case of overfitting



Applied the negative of gradients

Hyperparameters

Definition:

All the setup variables that not changing during the learning process. For example: number of layers.

Strategy:

Cross-validate different models on train/val set.

Grid search on the parameters (e.g., learning rate, regularizer factors).

Try log-space.

Hyperparameters



Today - Tricks and Tips in DL Training

- Overview
- Data Preprocessing
- Network Design
 - Activation Functions
 - Weight Initialization
 - Normalization
 - Monitoring the Learning Process
 - Hyperparameters
- **Optimization**
- Transfer Learning

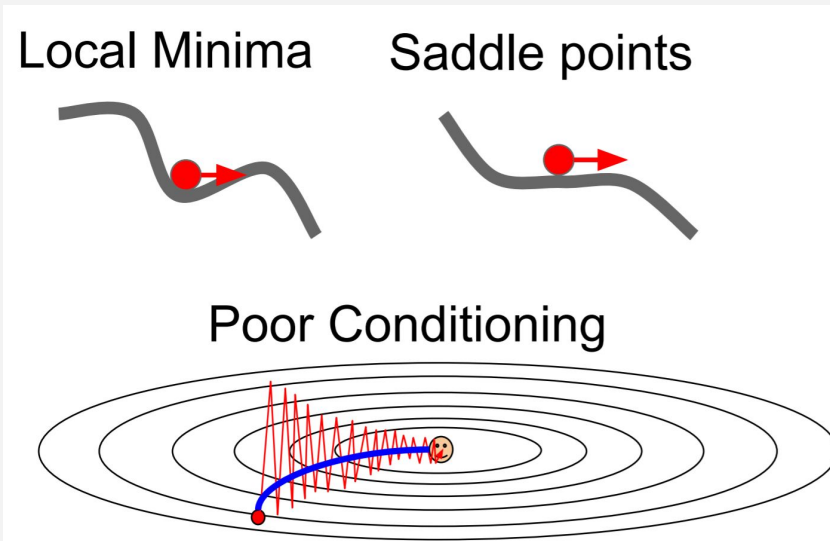
Optimization

Stochastic Gradient Descent (**SGD**)

Usually works with momentum.

Adagrad, RMSProp, Adadelata

Recommended: Adam

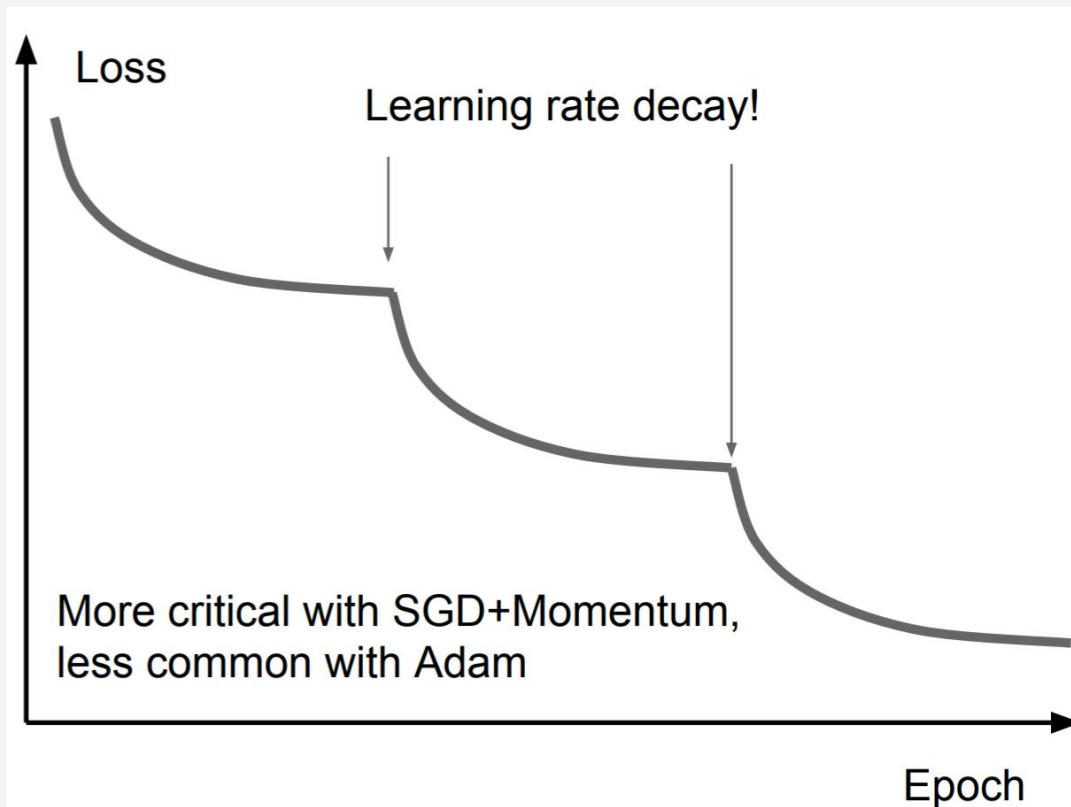


Learning Rate Decay Strategy

Step Decay

Exponential Decay

$1/t$ Decay

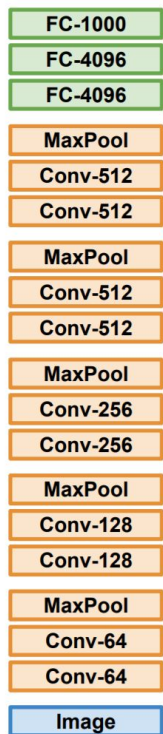


Today - Tricks and Tips in DL Training

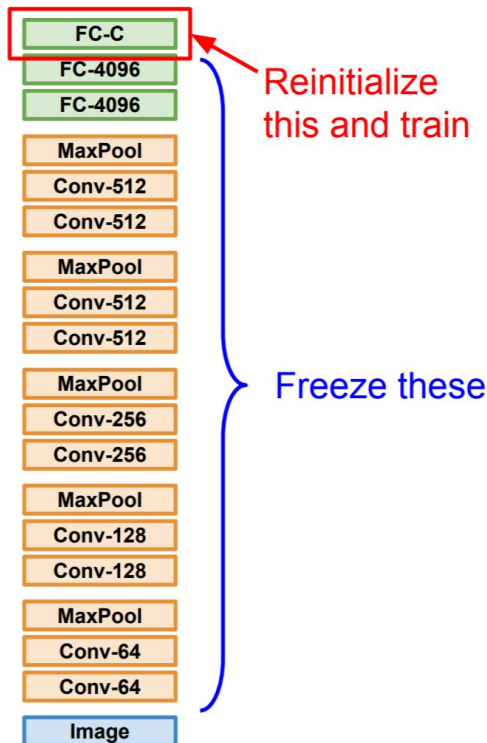
- Overview
- Data Preprocessing
- Network Design
 - Activation Functions
 - Weight Initialization
 - Normalization
 - Monitoring the Learning Process
 - Hyperparameters
- Optimization
- **Transfer Learning**

Transfer Learning

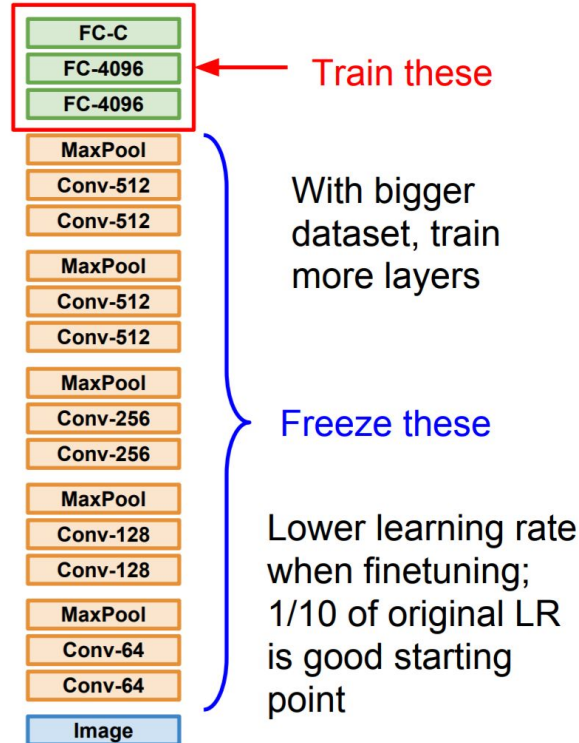
1. Train on Imagenet



2. Small Dataset (C classes)

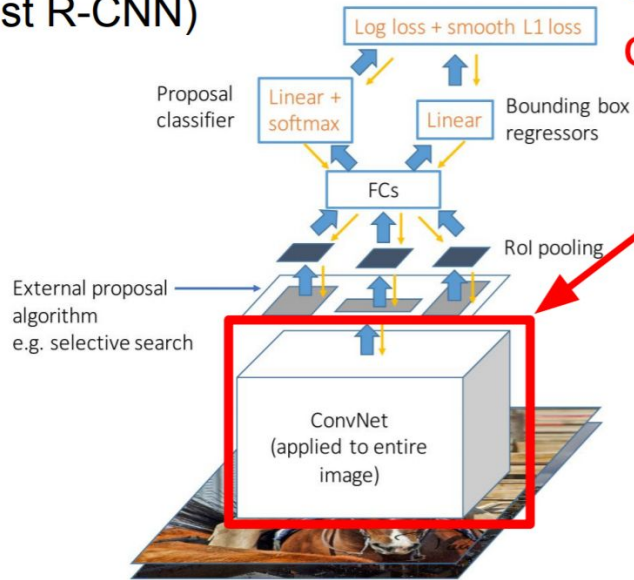


3. Bigger dataset



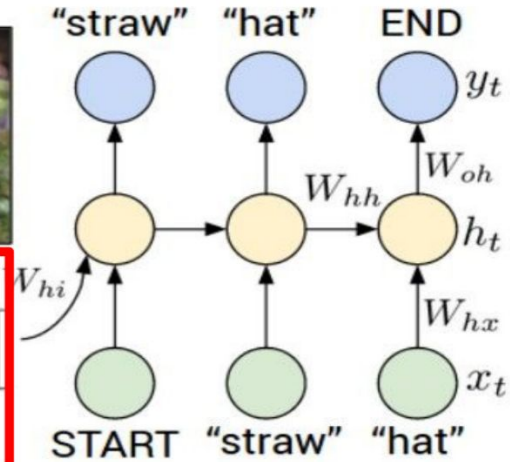
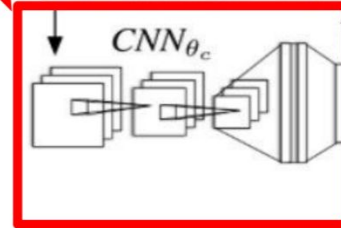
Transfer Learning Is the Norm

Object Detection (Fast R-CNN)



**CNN pretrained
on ImageNet**

Image Captioning: CNN + RNN



Where to Find Pretrained Models

Caffe:

<https://github.com/BVLC/caffe/wiki/Model-Zoo>

Tensorflow:

<https://github.com/tensorflow/models>

Pytorch:

<https://pytorch.org/docs/master/torchvision/>

AlexNet



VGG-16

Input → Conv 1-1 → Conv 1-2 → Pooling → Conv 2-1 → Conv 2-2 → Pooling → Conv 3-1 → Conv 3-2 → Conv 3-3 → Pooling → Conv 4-1 → Conv 4-2 → Conv 4-3 → Pooling → Conv 5-1 → Conv 5-2 → Conv 5-3 → Pooling → Dense → Dense → Dense → Output



Today - Tricks and Tips in DL Training

- Overview
- Data Preprocessing
- Network Design
 - Activation Functions
 - Weight Initialization
 - Normalization
 - Monitoring the Learning Process
 - Hyperparameters
- Optimization
- Transfer Learning