CS6421: Deep Neural Networks

Gregory Provan

Spring 2020 Lecture xx: Practical Issues

Based on notes from John Canny, Ismini Lourentzou

- Practical Issues in Deep Networks: Overview
- Initialisation
 - Architecture
 - Activation functions
 - Loss functions
 - Initial parameters
- Data Preprocessing
 Training

Motivation

Building high-performance deep networks is challenging

- Many parameters to set
- Architectures can be difficult to design
- May need data pre-processing

Practical Issues (to be addressed)

- Network initialisation
- Data pre-processing
- Network training
- Inference issues

Practical Issues in Network Initialisation: Overview

- Architecture
- Activation functions
- Loss functions
- Initial parameters
- Data Preprocessing
- Training

General Guidelines

Match application to architecture

- Structured data: CNN
 - Images, text, voice
- Temporal data: RNN/LSTM
 - Time-series analysis (stock forecasting)
- Use appropriate activation functions at each layer
- Weight initialisation
- Use hyper-parameter optimisation
 - Number of layers, size of layers, batch size, stride-length, etc.

- Practical Issues in Network Initialisation: Overview
 - Architecture
 - Activation functions
 - Loss functions
 - Initial parameters
- Data Preprocessing
- Training

Choose the architecture

Example: start with one hidden layer of 50 neurons:



Template

Neural Network



- Map application to architecture
- Always try state-of-the-art architectures first
 - Often available
 - These architectures work for good reasons
- Example: YOLO for vision applications

Zoo of Architectures





- Novel CNN architecture made huge impact
- Improved image recognition accuracy
- Nowadays almost all vision processing is done using deep learning

CNN's Topology



Feature extraction layer or Convolution layer

Detect the same feature at different positions in the input image.



♥ Introduced by LeCun.

 \odot raw image of 32 \times 32 pixels as input.



- \odot 5 × 5 Convolution matrix.
- S2 , S4 : Subsampling layer.
- Subsampling by factor 2.
- F6 : Fully connected layer.





All the units of the layers up to F6 have a sigmoidal activation function of the type:

$$y_j = \varphi(v_j) = A \tanh(Sv_j)$$







$$Y_j = \sum_{i=1}^{84} (F_i - W_{ij})^2, j = 0, ..., 9$$



About 187,000 connections
About 14,000 trainable weights





















- Practical Issues in Network Initialisation: Overview
 - Architecture
 - Activation functions
 - Loss functions
 - Initial parameters
- Data Preprocessing
- Training



Activation functions

Function types

- Sigmoid
- Softmax
- Tanh
- ReLU
- LeakyReLU
- Properties of functions
- Implications for classification, training, etc.

- Practical Issues in Network Initialisation: Overview
 - Architecture
 - Activation functions
 - Loss functions
 - Initial parameters
- Data Preprocessing
- Training



Loss functions and output

Classification

Training examples Rⁿ x {class_1, ..., class_n} (one-hot encoding)

Output Layer Soft-max [map Rⁿ to a probability distribution] $P(y = j \mid \mathbf{x}) = \frac{e^{\mathbf{x}^{\mathsf{T}}\mathbf{w}_j}}{\sum_{k=1}^{K} e^{\mathbf{x}^{\mathsf{T}}\mathbf{w}_k}}$

Cost (loss) function

Cross-entropy

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^{n} \sum_{k=1}^{K} \left[y_k^{(i)} \log \hat{y}_k^{(i)} + \left(1 - y_k^{(i)}\right) \log \left(1 - \hat{y}_k^{(i)}\right) \right]$$

Regression

Rⁿ x R^m



f(x)=x

Mean Squared Error

$$J(\theta) = \frac{1}{n} \sum_{i=1}^{n} (y^{(i)} - \hat{y}^{(i)})^2$$

Mean Absolute Error $J(\theta) = \frac{1}{n} \sum_{i=1}^{n} |y^{(i)} - \hat{y}^{(i)}|$

- Practical Issues in Network Initialisation: Overview
 - Architecture
 - Activation functions
 - Loss functions
 - Initial parameters
- Data Preprocessing
- Training



Weight initialization

- Typically: initialize to random values sampled from zero-mean Gaussian: $w \sim \mathcal{N}(0, \sigma^2)$
 - Standard deviation matters!
 - Key idea: avoid reducing or amplifying the variance of layer responses, which would lead to vanishing or exploding gradients

Common heuristics:

- $\sigma = 1/\sqrt{n_{in}}$, where n_{in} is the number of inputs to a layer
- $\sigma = 2/\sqrt{n_{\text{in}} + n_{\text{out}}}$ (Glorot and Bengio, 2010)
- $\sigma = \sqrt{2/n_{\text{in}}}$ for ReLU (<u>He et al.</u>, 2015)
- Initializing biases: just set them to 0

More details: http://cs231n.github.io/neural-networks-2/#init

Training

- Practical Issues in Network Initialisation: Overview
 - Architecture
 - Activation functions
 - Loss functions
 - Initial parameters
- Data Preprocessing



- May need to preprocess data
 - Insufficient data of particular types
 - Augmentation
 - Zero centering
 - Missing data
 - Etc.

Data augmentation

- Introduce transformations not adequately sampled in the training data
- Limited only by your imagination and time/memory constraints!
- Avoid introducing obvious artifacts



- Practical Issues in Network Initialisation: Overview
 - Architecture
 - Activation functions
 - Loss functions
 - Initial parameters
- Data Preprocessing
- Training

Training: Many Issues to Consider

Training is non-linear optimisation

Training: optimize network parameters to minimise loss over training set

 $\theta^* = \arg \min_{\theta} \sum_{(x,y) \in (X,Y)} J[y, fL(x, \theta_{1, \dots, L})]$

- Training is extremely complex process
- Tools exist to assist in training
- Hyper-parameter optimisation
 - Define hyper-parameters
 - Monitor and optimise over training process

Overfitting

How to choose hyper-parameters?

- The learning rate η
- Mini-batch size m
- Early stopping
- Learning rate schedules
- -Regularization parameter λ
- Grid search and the automated technique

Overfitting

to



- Trading off performance vs.
 size/parameter-space of network
 - Bias/variance trade-off
- Regularisation parameter λ

$$J_{reg}(\theta) = J(\theta) + \lambda \sum_k \theta_k^2$$

Regularization methods

- Address problem of overfitting
- Methods
 - Weight decay
 - L² normalization
 - L¹ normalization
 - Dropout
 - Artificial expansion of the training data

Regularization



Dropout

- Randomly drop units (along with their connections) during training
- Each unit retained with fixed probability p, independent of other units
- Hyper-parameter p to be chosen (tuned)

Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." Journal of machine learning research (2014)

L2 = weight decay

Regularization term that penalizes big weights,

added to

- Weight decay value determines how dominant regularization is $J_{reg}(\theta) = J(\theta) + \lambda \sum_{k} \theta_k^2$ during gradient computation
- Big weight decay coefficient \rightarrow big penalty for big weights

 \otimes

Early-stopping

- Use validation error to decide when to stop training
- Stop when monitored quantity has not improved after n subsequent epochs
- n is called patience

Tuning hyper-parameters



Chimportant parameter

 $g(x) \approx g(x) + h(y)$

g(x) shown in green h(y) is shown in yellow

Bergstra, James, and Yoshua Bengio. "<u>Random</u> <u>search for hyper-parameter optimization.</u>" Journal of Machine Learning Research, Feb (2012)

Important parameter

Important parameter

"Grid and random search of 9 trials for optimizing function $g(x) \approx g(x) + h(y)$ With grid search, nine trials only test g(x) in three distinct places. With random search, all nine trials explore distinct values of g."

Both try configurations randomly and **blindly** Next trial is independent to all the trials done before

Bayesian optimization for hyper-parameter tuning:

Library available!

Make smarter choice for the next trial, minimize the number of trials

- 1. Collect the performance at several configurations
- 2. Make inference and decide what configuration to try next

Learning rates

Always use smaller rates



From [Duda&Hart:2001]



- Overview of Practical Issues in Deep Networks
- Key issues
 - Initialisation
 - Data Preprocessing
 - Training
- Will examine each major issue