

# Deep Learning: Review Notes

January 12, 2024

## Lecture 4: Linear Classification

- Soft-max Classifier (Multinomial Logistic Regression):

$$s_i = \frac{e^{o_i}}{\sum_{j=1}^K e^{o_j}}$$

- Loss function:

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \ell(f(X_i, \theta), y_i)$$

- Cross-Entropy loss function: Let  $p_k$  be the predicted probability that the instance belongs to class  $k$ , and  $p = (p_1, \dots, p_K)$ . Then

$$\ell(p, y) = - \sum_k I(y = k) \log p_k$$

- What is the connection of cross-entropy and maximum-likelihood estimator?
- What is the connection of cross-entropy and KL distance?

## Lecture 5: Multi-Layer Perceptron

- Definition of Single layer network (perception model):  $f_w(x) = \sigma(wx + b)$
- Show that perception model can only solve linear Separable problems. Provide some example of non Linear Separable Problems such as XOR.
- Definition of multi-layer network and related terminology: input layer, hidden layers, last layers, active function, neuron (node), Weights and Biases.
- Feedforward Algorithm:

$$f_W(x) = h_L \circ h_2 \circ \dots \circ h_1(x)$$

where  $h_i(x) = a(w_i x + b_i)$

- Universal Approximation Theorem (for two layers and for Width-Bounded ReLU Networks).

## Lecture 8: Back-propagation Algorithm

- Computational Graph and automatic differentiation.
- Back-propagation Algorithm:

$$\frac{\partial \ell(f_W(x), y)}{\partial h_i} = \frac{\partial \ell}{\partial h_L} \cdot \frac{\partial h_L}{\partial h_{L-1}} \cdots \frac{\partial h_{i+1}}{\partial h_i}$$

$$\frac{\partial \ell}{\partial W} = \sum_i \frac{\partial \ell}{\partial h_i} \cdot \frac{\partial h_i}{\partial W}$$

Also, we know that

$$\text{Downstream Gradient} = \text{Local Jacobian} \times \text{Upstream Gradient}$$

## Lecture 9: Optimization

- Stochastic Gradient Descent (SGD), and convergence theorem.
- Classic Robbins Monro Condition:  $\sum_{i=1}^{\infty} \eta_i = \infty, \sum_{i=1}^{\infty} \eta_i^2 < \infty$
- Comparing SGD, GD, and mini-batch
- Momentum and Nesterov Momentum. What is the intuition behind these two ideas?
- AdaGrad, RMSprop, Adam algorithms
- Second order optimization. Why is this impractical
- FBGS (optional)

## Lecture 10: Convolutional Neural Networks

- Convolution operation on images definition.
- Convolution layer, padding, stride, tensors. kernel, downsampling
- Output size formula:  
Input size:  $C_{in} \times H \times W$ , and  $C_{out}$   
Hyperparameters:
  - Kernel size:  $K_H \times K_W$
  - number of filters:  $C_{out}$

- Padding:  $P$ , Stride:  $S$
- filters of size  $C_{in} \times K_H \times K_W$

Number of learnable parameters:

Weight matrix:  $C_{out} \times C_{in} \times K_H \times K_W$

bias:  $C_{out}$

Output size:  $C_{out} \times H' \times W'$

$$H' = (H - K + 2P)/S + 1$$

$$W' = (W - K + 2P)/S + 1$$

Number of multiply operations:  $C_{out} \times H' \times W' \times C_{in} \times K_H \times K_W$

- You do not need to memorize the architecture of different network. But it is good (optional) to know the names of some famous architecture: VGG19, Resnet, GoogleNet, DenseNet

## Lecture 11: Training DNN

- Different activation function: sigmoid, tanh, ReLU, LeakyRelu, ELU, GLU
- Why learning does not happen for saturated neurons? Why sigmoid is not a good choice?
- Compare sigmoid, tanh, Relu as activation function.
- What is the advantage of Leaky ReLU to ReLU?
- Why does initialization important?
- Show that  $var_{input} = var_{output}$  if  $var_w = \frac{1}{\sqrt{n_{in}}}$ ? (LeCun Formula)
- What is Xavier initialization?  $var(w) = \frac{2}{\sqrt{n_{in} + n_{out}}}$
- What is batch Normalization?
- How does batch normalization help?
- What are different regularization? Why does regularization important? (overfitting)
- What is Dropout? How does Dropout helps? (co-adaptation, and ensemble)

## Lecture 12: Word Embedding

- What is the advantage of Word2vec compare to one-hot embedding? (Similar words are embed to closer vectors)
- Skip-gram model and negative sampling idea

## Lecture 13: Langue Models

- What is N-Gram Language Model? What is sparsity problem? What is storage problem?
- Simple RNN model and time-back propagation
- Perplexity and its connection to likelihood
- What is vanishing/exploding gradient? Why it happens in long RNN?
- What is gradient clipping?

## Lecture 14: LSTM

- Be familiar with LSTM: what is forget, input and output gates? How does Cell update? How does cell solve the vanishing problem?
- Why vanishing gradient happens in too deep networks?
- Residual connection and dense connection and highway connection solve vanishing gradient for DNN
- Bi-directional and multilayer RNN

## Lecture 15: Attention

- Machine translation problem as probability problem
- BLEU (optional)
- Multi-layer deep encoder-decoder for Machine translation
- Beam search algorithm
- What is attention? why attention important? (example from machine translation)

## Lecture 16: Transformer

- Issues with RNN models: Linear interaction distance, lack of paralization
- Self Attention: What is Key, Query, value ? Attention can be seen as a soft look up table
- Position embedding. Why do we need it?
- Masking in decoding
- Multi-head Attention (look at the exercise for details)
- Transformer: Decoding, encoding, Decoding-encoding: cross-attention.