

## Large-scale deep learning

- Fast CPU Implementations
- GPU Implementations

Operating System perspective, GPU has lower clock speed and less branching capacity, thus intensive and independent computing can use GPU.

- Large-Scale Distributed Implementations

Model parallelism or data parallelism

- Model compression

- Dynamic structure

- The idea behind dynamic structure is that focusing on the part, which influences the results the most (data or model)

- A cascade of classifiers can efficiently separate the complicated task.

- Mixture of experts (demonstrated in PRML)

- Context and attention mechanism: Focus on the weighted part of hidden/input data (same idea with data partition)

- Specialized Implementation

## Computer Vision

Most of the computer vision tasks can be seen as expansion of object recognition and object detection.

- Preprocessing

- Formatting the image to the same scale  $\Rightarrow$  only strictly necessary stuff.

- Data augmentation without changing or losing key information.

Two examples of data preprocessing in CV. (Dataset augmentation omitted)

- Contrast Normalization: Formula reads below

$$X'_{i,j,k} = S \frac{X_{i,j,k} - \bar{X}}{\max\{ \epsilon, \sqrt{\lambda + \frac{1}{3rc} \sum_i \sum_j \sum_k (X_{i,j,k} - \bar{X})^2} \}}$$

Ensuring the denominator not be zero

Scaling factor

$\rightarrow$  normalize the pixels

$\downarrow$  shift factor

$\rightarrow$  The standard variance of pixels

Contrast normalization is for rescaling all the pixels into the same scale, and the  $\lambda, \epsilon, S$  parameters is for ensuring the computational stability and decreasing the effect of noise.

- Global Contrast Norm vs Local Contrast Norm

Focus more on scaling all the pixels within one region for stability.

$\downarrow$   
More sensitive to edge, cause in small windows we often get local invariant.

## Speech Recognition

General method: HMM-GMM  $\rightarrow$  Conv + others  $\rightarrow$  RNN LSTM

Tasks of speech recognition are different from NLP on, it considering converting input radio wave into human understandable format, phonemic or words.

## Natural Language Processing (NLP)

Many NLP applications are based on language models that define a probability distribution over sequence of words, characters, or bytes in natural language.

- n-grams

Taking use of chain rule:  $P(x_1, \dots, x_t) = P(x_1, \dots, x_{t-1}) \prod_{t=n}^T P(x_t | x_{t-n}, \dots, x_{t-1})$

$\rightarrow$  Modeling a probability of a whole sequence .

Here we don't consider application specification, we focus on how modeling, a probability of sequence/sentence is calculated by recurrently applying chain rule with every n words.

One way to view n-gram is to regard it as  $kNN$  cause the transition probability is actually calculated based on counting on training dataset.

- Neural Language Models

Embedding word into a reasonable subspace, unlike other tasks, data in NLP often comes discrete and non-digital. (Note that "embedding" is a different task from sequential modeling)

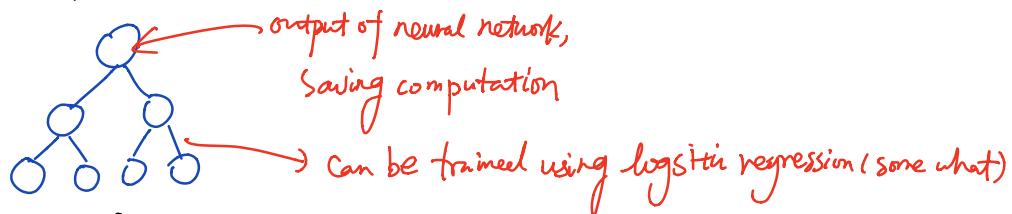
- How we deal with high dimensional output

- Taking use of a shorted list  $L$  with its tail  $T$ .

$L$  (shorter list) reduce the dimension of model and Tail ( $T$ ) is used as n-gram of extra output (rare vocabularies)

- Hierarchical softmax

Using a tree structure to representing vocabularies and when inferencing, output actually searches the best-matched path from root node to the final leaf (words predictions)



- Importance Sampling

Aimed to use Sampling method to reduce the computational cost of  $\frac{\partial \log P(y|c)}{\partial \theta}$

$$\frac{\partial \log P(y|c)}{\partial \theta} = \frac{\partial}{\partial \theta} \log \frac{e^{ay}}{\sum_i e^{ai}} = \frac{\partial y}{\partial \theta} - \sum_i p(y=i|c) \frac{\partial a_i}{\partial \theta}$$

Here is a form of expectation with can be sampling by weights of importance sampling  $\sum_i p(y=i|c) \frac{\partial a_i}{\partial \theta} \doteq \frac{1}{m} \sum_i w_i \frac{\partial a_i}{\partial \theta}$

- Noise-Contrastive Estimation and Ranking Loss
- Combining NLM with n-gram
  - Ensembling :  $Z_{mr}$  independent
  - Adding extra input to output links to calculate n-gram.
- Neural Machine Translation
  - Populating encoder-decoder style model
  - Attention mechanism
 
$$c = \sum_t a_t h_t \rightarrow \text{Can also be input} \quad \left[ \begin{array}{l} \text{This yields each } t \text{ has an attention} \\ \text{re-weighting each } h_t \text{ / input's contribution} \end{array} \right] \rightarrow \text{of particular location of } h.$$

If we select desired  $h_t$  without reweighting, the problem is not differentiable.
  - Word embedding can be seen as data pre-processing for these sequential based model, for better performance.

### Other Applications

- Recommending System (supervised most)
  - collaborative filtering (Bi-linear model or SVD)
  - content-based recommender (For dealing with new-seen data)
  - exploitation<sup>2nd</sup> vs exploration<sup>1st</sup>
- Study within current knowledge → Expanding our knowledge scope
- knowledge representation, reasoning and question answering
  - learning relations ( $S = \{(A, B), (C, D), (A, C)\}$ , tuples)
  - Training data preparation → relational database
- Model structure → Inspired by NLP