

Deep Learning for Automatic Speech Recognition – Part II

Xiaodong Cui

IBM T. J. Watson Research Center Yorktown Heights, NY 10598

Fall, 2018





Outline

- A brief revisit of sampling, pitch/formant and MFCC
- DNN-HMM (hybrid) acoustic modeling
- Advanced acoustic modeling
 - Convolutional neural networks (CNNs)
 - Recurrent neural networks (RNNs) and long short-term memory (LSTM) networks
 - Convolutional long short-term memory (CLSTM) networks





Pitch and Formants









Mel-Frequency Filter Bank and MFCC

 $m = 2595 \log_{10} \left(1 + \frac{f}{100} \right)$



- multi-resolution on frequency bins
- use DCT to replace IDFT for better dimension decorrelation and energy clustering
- typical order around 13 for speech recognition
- increased order for speaker recognition/verification (typically 19-22)



Input Speech Features for DNNs

- Commonly-used hand-crafted features
 - MFCCs
 - Mel-frequency filter bank (FBank)
 - Speaker adapted features (FMLLR)
 - Appended speaker embedding vectors (i-vectors)
- Let DNNs learn the features
 - Power spectra
 - Raw audio signals



Two Streams of DNN Acoustic Models

- Hybrid DNNs
 - Commonly referred to as DNN-HMM or CD-DNN-HMM
 - Use GMM-HMM alignments as labels
 - Use dictionary and language model for decoding.
- End-to-end (E2E) DNNs (will be addressed in next lecture)
 - Directly deal with sequence-to-sequence mapping problem with unequal consequence lengths
 - Do not need alignments, dictionary and language model in principle.
 - ► Two E2E architectures:
 - Connectionist Temporal Classification (CTC)
 - Encoder-Decoder Attention models



- Given the text label, how to find the best underlying state sequence?
 - same as decoding except the label is known
 - Viterbi algorithm (dynamic programming)
- Often referred to as Viterbi alignments in speech community
- Widely used in deep learning ASR to generate the target labels for the training data

Context-Dependent DNNs





DNN-HMMs: An Typical Training Recipe

- Preparation
 - 40-dim FBank features with ± 4 adjacent frames (input dim = 40x9)
 - use an existing model to generate alignments which are then converted to 1-hot targets from each frame
 - create training and validation data sets
 - estimate priors p(s) of CD states
- Training
 - set DNN configuration (multiple hidden layers, softmax output layers and cross-entropy loss function)
 - initialization
 - optimization based on back-prop using SGD on the training set while monitoring loss on the validation set
- Test
 - push input features from test utterances through the DNN acoustic model to get their posteriors
 - convert posteriors to likelihoods
 - Viterbi decoding on the decoding network
 - measure the word error rate



Training A Hybrid DNN-HMM System





- Bio-inspired feed-forward neural networks*
 - individual cortical neurons in visual cortex respond to stimuli in a restricted region of space known as the receptive field.
 - receptive fields of different neurons partially overlap such that they tile the visual field.
 - response of an individual neuron to stimuli within its receptive field can be approximated mathematically by a convolution operation.
- LeNet-5
 - ▶ Y. LeCun, L. Bottou, Y. Bengio, P. Haffner (1998). "Gradient-based learning applied to document recognition," Proc. of the IEEE.
 - A 7-layer CNN that outperformed other techniques on a standard handwritten digit recognition task.
- Advantages:
 - Local (sparse) connectivity
 - Weight sharing
 - translation-invariant

*adapted from wikipedia



· Convolutional layers with nonlinearity

$$\begin{split} z_{ij}^k &= \sigma((W^k \ast x)_{ij} + b^k) \\ W(m,n) \ast x(m,n) &= \sum_{u=-\infty}^{+\infty} \sum_{v=-\infty}^{+\infty} x(u,v) W(u-m,v-n) \\ W(m,n) &\neq 0 \quad \text{only for } 0 \leq m \leq M, \quad 0 \leq n \leq N \end{split}$$

- · Pooling layers
 - subsampling
 - typically max-pooling or average-pooling
- Fully connected layers







































A Brief Summary of Some CNN Terminology

- feature map
- local receptive field (filter or kernel)
- padding
- stride
- convolution
- pooling
 - max-pooling
 - average-pooling
 - pooling along different axes
 - resulting different resolutions

Tensor Representation



- weights W_{mn}^{pq} connecting each unit of the *p*-th feature map at layer l-1 and the *q*-th feature map of layer l with the local receptive filter w(m, n).
- number of parameters in the weights that connect the two convolutional layers is (ignore biases here)

$$|W^{pq}_{mn}| = M \times N \times P \times Q$$

where M and N are the dimensions of the local receptive filters and P and Q are the numbers of feature maps in layer l-1 and layer l.

relation of input and output dimensionality

$$O = \frac{W - K + 2P}{S} + 1$$

where W is the input dimensionality, ${\cal O}$ the output dimensionality, K the filter size, ${\cal P}$ the padding and S the stride.



An Example of CNN Acoustic modeling





Notable Development In CNNs

- AlexNet
 - Alex Krizhevsky, Ilya Sutskever and Geoffrey E. Hinton, "ImageNet Classification with Deep Convolutional Networks," NIPS 2012.
 - 5 conv layers, 3 fully-connected layers, ReLU, max-pooling, dropout, data augmentation, 2-GPU implementation
- ZFNet
 - Matthew D. Zeiler and Rob Fergus, "Visualizing and understanding convolutional networks." ECCV, 2014.
 - similar architecture to AlexNet, reduced window size and stride step in first conv layer.
- VGGNet
 - Karen Simonyan and Andrew Zisserman, "Very deep convolutional neural networks for large-scale image recognition," ICLR 2015.
 - up to 19 conv layers small 3x3 filters with stride 1, padding, 2x2 max-pooling layers with stride 2 every 2,3 or 4 conv layers
- ResNet
 - Kaiming He, Xiangyu Zhang, Shaoqing Ren and Jian Sun, "Deep Residual Learning for Image Recognition," CVPR 2016.
 - 152 layers with 3x3 filters
 - introduced residual learning by a direct by-pass identity link
- DenseNet
 - Gao Huang, Zhuang Liu, Laurens van der Maaten, Kilian Q. Weinberger, "Densely connected convolutional networks," CVPR 2017.
 - do not throw away anything you learn along the way

VGG Acoustic Modeling



T. Sercu, C. Puhrsch, B. Kingsbury and Y. LeCun, "Very Deep Multilingual Convolutional Neural Networks for LVCSR," ICASSP 2016.



Recurrent Neural Networks (RNNs)

- Sequential data is ubiquitous – speech, video, text, stock index ···
- Sequence modeling is crucial – how to model the history?

$$oldsymbol{h}_t = g_t(oldsymbol{x}_t, oldsymbol{x}_{t-1}, oldsymbol{x}_{t-2}, \cdots, oldsymbol{x}_1)$$

• RNNs model sequences with shared parameters (functions)

$$\boldsymbol{h}_t = f(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t; \boldsymbol{\theta})$$



Compact and Unfolded Representations of RNNs



compact

unfolded

$$egin{aligned} m{a}_t &= m{b} + m{W}m{h}_{t-1} + m{U}m{x}_t \ m{h}_t &= ext{tanh}(m{a}_t) \ m{o}_t &= m{c} + m{V}m{h}_t \ m{\hat{y}}_t &= ext{softmax}(m{o}_t) \end{aligned}$$





one-to-many mapping — image captioning.





many-to-one mapping — sentiment analysis.





synchronized many-to-many mapping — speech recognition, video labeling.





asynchronized many-to-many mapping — speech recognition, machine translation.



Back-Propagation Through Time (BPTT)



Forward Propagation:

 $a_t = b + Wh_{t-1} + Ux_t$ $h_t = \tanh(a_t)$ $o_t = c + Vh_t$ $\hat{\boldsymbol{y}}_t = \operatorname{softmax}_K(\boldsymbol{o}_t)$ $\mathcal{L}_t = -\sum_{k=1}^{K} y_{t,k} \log \hat{y}_{t,k}$ $\mathcal{L} = \sum_{t=1}^{T} \mathcal{L}_t$



Back-Propagation Through Time (BPTT) Parameters to be optimized: W, U, V, b and c

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial \mathbf{W}} &= \sum_{t} \sum_{i} \frac{\partial \mathbf{h}_{t}^{i}}{\partial \mathbf{W}} \frac{\partial \mathcal{L}}{\partial \mathbf{h}_{t}^{i}} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{U}} &= \sum_{t} \sum_{i} \frac{\partial \mathbf{h}_{t}^{i}}{\partial \mathbf{U}} \frac{\partial \mathcal{L}}{\partial \mathbf{h}_{t}^{i}} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{b}} &= \sum_{t} \frac{\partial \mathbf{h}_{t}}{\partial \mathbf{b}} \frac{\partial \mathcal{L}}{\partial \mathbf{h}_{t}} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{V}} &= \sum_{t} \sum_{i} \frac{\partial \mathbf{h}_{t}^{i}}{\partial \mathbf{V}} \frac{\partial \mathcal{L}}{\partial \mathbf{h}_{t}^{i}} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{c}} &= \sum_{t} \frac{\partial \mathbf{o}_{t}}{\partial \mathbf{c}} \frac{\partial \mathcal{L}}{\partial \mathbf{o}_{t}} \end{cases}$$

which involves gradients of two internal nodes o_t and h_t .

$$rac{\partial \mathcal{L}}{\partial oldsymbol{o}_t}, \quad rac{\partial \mathcal{L}}{\partial oldsymbol{h}_t}$$



Back-Propagation Through Time (BPTT)

$$rac{\partial \mathcal{L}}{\partial oldsymbol{o}_t} = rac{\partial \hat{oldsymbol{y}}_t}{\partial oldsymbol{o}_t} rac{\partial \mathcal{L}_t}{\partial \hat{oldsymbol{y}}_t} rac{\partial \mathcal{L}_t}{\partial \mathcal{L}_t} = \hat{oldsymbol{y}}_t - oldsymbol{y}_t$$

evaluate the gradient of h_t backward starting from the last time step au:

$$\begin{split} \frac{\partial \mathcal{L}}{\partial \boldsymbol{h}_{\mathsf{T}}} &= \frac{\partial \boldsymbol{o}_{\mathsf{T}}}{\partial \boldsymbol{h}_{\mathsf{T}}} \frac{\partial \mathcal{L}_{\mathsf{T}}}{\partial \boldsymbol{o}_{\mathsf{T}}} \frac{\partial \mathcal{L}}{\partial \mathcal{L}_{\mathsf{T}}} = \mathbf{V}^{\mathsf{T}} (\hat{\boldsymbol{y}}_{\mathsf{T}} - \boldsymbol{y}_{\mathsf{T}}) \\ \frac{\partial \mathcal{L}}{\partial \boldsymbol{h}_{t}} &= \frac{\partial \boldsymbol{o}_{t}}{\partial \boldsymbol{h}_{t}} \frac{\partial \mathcal{L}}{\partial \boldsymbol{o}_{t}} + \frac{\partial \boldsymbol{h}_{t+1}}{\partial \boldsymbol{h}_{t}} \frac{\partial \mathcal{L}}{\partial \boldsymbol{h}_{t+1}} \\ &= \mathbf{V}^{\mathsf{T}} (\hat{\boldsymbol{y}}_{t} - \boldsymbol{y}_{t}) + \frac{\partial \boldsymbol{a}_{t+1}}{\partial \boldsymbol{h}_{t}} \frac{\partial \boldsymbol{h}_{t+1}}{\partial \boldsymbol{a}_{t+1}} \frac{\partial \mathcal{L}}{\partial \boldsymbol{h}_{t+1}} \\ &= \mathbf{V}^{\mathsf{T}} (\hat{\boldsymbol{y}}_{t} - \boldsymbol{y}_{t}) + \mathbf{W}^{\mathsf{T}} \cdot \mathbf{diag} \left(1 - \boldsymbol{h}_{t+1}^{2} \right) \cdot \frac{\partial \mathcal{L}}{\partial \boldsymbol{h}_{t+1}} \end{split}$$



Back-Propagation Through Time (BPTT)

$$\begin{split} \frac{\partial \mathcal{L}}{\partial \mathbf{W}} &= \sum_{t} \operatorname{diag} \left(1 - \boldsymbol{h}_{t}^{2} \right) \cdot \frac{\partial \mathcal{L}}{\partial \boldsymbol{h}_{t}} \cdot \boldsymbol{h}_{t-1}^{\mathsf{T}} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{U}} &= \sum_{t} \operatorname{diag} \left(1 - \boldsymbol{h}_{t}^{2} \right) \cdot \frac{\partial \mathcal{L}}{\partial \boldsymbol{h}_{t}} \cdot \boldsymbol{x}_{t}^{\mathsf{T}} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{b}} &= \sum_{t} \operatorname{diag} \left(1 - \boldsymbol{h}_{t}^{2} \right) \cdot \frac{\partial \mathcal{L}}{\partial \boldsymbol{h}_{t}} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{V}} &= \sum_{t} \frac{\partial \mathcal{L}}{\partial \boldsymbol{o}_{t}} \cdot \boldsymbol{h}_{t}^{\mathsf{T}} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{c}} &= \sum_{t} \frac{\partial \mathcal{L}}{\partial \boldsymbol{o}_{t}} \end{split}$$

An LSTM Memory Cell



- input gate:
- forget gate:
- cell candidate:
- cell:
- output gate:
- hidden state output:
- $$\begin{split} \mathbf{i}_t &= \sigma(\mathbf{W}_{ix} \mathbf{x}_t + \mathbf{W}_{ih} \mathbf{h}_{t-1} + \mathbf{b}_i) \\ \mathbf{f}_t &= \sigma(\mathbf{W}_{fx} \mathbf{x}_t + \mathbf{W}_{fh} \mathbf{h}_{t-1} + \mathbf{b}_f) \\ \mathbf{g}_t &= \tanh(\mathbf{W}_{cx} \mathbf{x}_t + \mathbf{W}_{ch} \mathbf{h}_{t-1} + \mathbf{b}_c) \\ \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t \\ \mathbf{o}_t &= \sigma(\mathbf{W}_{ox} \mathbf{x}_t + \mathbf{W}_{oh} \mathbf{h}_{t-1} + \mathbf{b}_o) \\ \end{split}$$









 Sepp Horchreiter and Jürgen Schmidhuber, "Long Short-Term Memory," Neural Computation 9, 1997

sporopusking g(net.)

 a recurrent self-connection with weight 1.0 – constant error carousel (CEC)



- Felix A. Gers, Jürgen Schmidhuber and Fred Cummins. "Learning to Forget: Continual Prediction with LSTM." Neural Computation 12,2000
- A forget gate added to learn to reset



A Bi-directional Multi-layer LSTM for ASR





A CLSTM for ASR