## **Acoustic Modeling: Tied-state HMMs & DNN-based models**





Instructor: Preethi Jyothi

Lecture 7



#### **Recall: Acoustic Model**



# **Triphone HMM Models**

- Each phone is modelled in the context of its left and right neighbour phones •
- Number of triphones that appear in data  $\approx$  1000s or 10,000s •
- - Hundreds of millions of parameters! •
- across triphone models!

 Pronunciation of a phone is influenced by the preceding and succeeding phones. E.g. The phone [p] in the word "peek" : p iy k" vs. [p] in the word "pool" : p uw l

If each triphone HMM has 3 states and each state generates *m*-component GMMs  $(m \approx 64)$ , for d-dimensional acoustic feature vectors ( $d \approx 40$ ) with  $\Sigma$  having  $d^2$  parameters

Insufficient data to learn all triphone models reliably. What do we do? Share parameters



# **Parameter Sharing**

- Sharing of parameters (also referred to as "parameter tying") can be done at any level:
  - Parameters in HMMs corresponding to two triphones are said to be tied if they are identical



 More parameter tying: Tying variances of all Gaussians within a state, tying variances of all Gaussians in all states, tying individual Gaussians, etc.

# **1. Tied Mixture Models**

- All states share the same Gaussians (i.e. same means and covariances)
- Mixture weights are specific to each state Triphone HMMs (No sharing)





Triphone HMMs (Tied Mixture Models)

# 2. State Tying

• generate acoustically similar data



Triphone HMMs (No sharing)





Triphone HMMs (State Tying)

# Observation probabilities are shared across states which

#### **Tied state HMMs**



Image from: Young et al., "Tree-based state tying for high accuracy acoustic modeling", ACL-HLT, 1994

Four main steps in building a tied state HMM system:

- Create and train 3-state monophone HMMs with single Gaussian observation probability densities
- Clone these monophone distributions to initialise a set of untied triphone models. Train them using Baum-Welch estimation. Transition matrix remains common across all triphones of each phone.
- ... etc 3. For all triphones derived from the same monophone, cluster states whose parameters should be tied together.
  - 4. Number of mixture components in each tied state is increased and models re-estimated using BW

## **Tied state HMMs**



Image from: Young et al., "Tree-based state tying for high accuracy acoustic modeling", ACL-HLT, 1994

each

moo

Four main steps in building a tied state HMM system:

- Create and train 3-state monophone HMMs with single Gaussian observation probability densities
- 2. Clone these monophone distributions to initialise a set of untied triphone models. Train them using Baum-Welch estimation. Transition matrix remains common across all triphones of each phone.

etc 3. For all triphones derived from the same monophone, cluster states whose parameters should be tied toge
4. Num Which states should be tied in

together? Use decision trees.

## **Decision Trees**

#### Classification using a decision tree: Begins at the root node: What property is satisfied? Depending on answer, traverse to different branches



## **Decision Trees**

- Given the data at a node, either declare the node to be a leaf or find another property to split the node into branches.
- Important questions to be addressed for DTs:
  - 1. How many splits at a node? Chosen by the user.
  - 2. Which property should be used at a node for splitting? One which decreases "impurity" of nodes as much as possible.
  - 3. When is a node a leaf? Set threshold in reduction in impurity

## **Tied state HMMs**



Image from: Young et al., "Tree-based state tying for high accuracy acoustic modeling", ACL-HLT, 1994

each

moo

Four main steps in building a tied state HMM system:

- Create and train 3-state monophone HMMs with single Gaussian observation probability densities
- 2. Clone these monophone distributions to initialise a set of untied triphone models. Train them using Baum-Welch estimation. Transition matrix remains common across all triphones of each phone.

etc 3. For all triphones derived from the same monophone, cluster states whose parameters should be tied toge
4. Num Which states should be tied in

together? Use decision trees.

#### How do we build these phone DTs?

#### 1. What questions are used?

Linguistically-inspired binary questions: "Does the left or right phone come from a broad class of phones such as vowels, stops, etc.?" "Is the left or right phone [k] or [m]?"

#### How do we build these phone DTs?

#### 1. What questions are used?

right phone [k] or [m]?"

2. What is the training data for each phone state,  $p_i$ ? (root node of DT)

Linguistically-inspired binary questions: "Does the left or right phone come from a broad class of phones such as vowels, stops, etc.?" "Is the left or

# **Training data for DT nodes**

- Align training data,  $x_i = (x_{i1}, ..., x_{iT_i})$  i=1...N where  $x_{it} \in \mathbb{R}^d$ , against a set of triphone HMMs
- Use Viterbi algorithm to find the best HMM state sequence corresponding to each  $x_i$
- Tag each  $x_{it}$  with ID of current phone along with left-context and right-context



of the 3-state HMM corresponding to the triphone b/aa/g

• For a state *j* in phone *p*, collect all  $x_{it}$ 's that are tagged with ID  $p_i[?/?]$ 

sil/b/aa b/aa/g aa/g/sil

 $x_{it}$  is tagged with ID aa<sub>2</sub>[b/g] i.e.  $x_{it}$  is aligned with the second state

#### How do we build these phone DTs?

#### 1. What questions are used?

right phone [k] or [m]?"

2. What is the training data for each phone state,  $p_i$ ? (root node of DT)

has *p* as the middle phone

3. What criterion is used at each node to find the best question to split the data on?

give the maximum increase in log likelihood

Linguistically-inspired binary questions: "Does the left or right phone come from a broad class of phones such as vowels, stops, etc.?" "Is the left or

All speech frames that align with the  $j^{th}$  state of every triphone HMM that

Find the question which partitions the states in the parent node so as to

# Likelihood of a cluster of states

- K $\mathcal{L}(S) = \sum \sum$  $i=1 \ s \in S$
- following quantity:

 $\Delta_q = \mathcal{L}(S^q_{\text{ves}})$ 

- the question for which  $\Delta_q$  is the biggest
- with a split falls below a threshold

• If a cluster of HMM states,  $S = \{s_1, s_2, ..., s_M\}$  consists of M states and a total of K acoustic observation vectors are associated with S,  $\{x_1, x_2, \dots, x_K\}$ , then the log likelihood associated with S is:

$$\log \Pr(x_i; \mu_S, \Sigma_S) \gamma_s(x_i)$$

• For a question q that splits S into S<sub>yes</sub> and S<sub>no</sub>, compute the

$$\mathcal{L}(S^q_{\mathrm{no}}) - \mathcal{L}(S)$$

Go through all questions, find  $\Delta_q$  for each question q and choose

• Terminate when: Final  $\Delta_q$  is below a threshold or data associated

## Likelihood criterion



Given a phonetic question, let the initial set of untied states S be split into two partitions S<sub>ves</sub> and S<sub>no</sub>

Each partition is clustered to form a single Gaussian output distribution with mean  $\mu_{Syes}$  and covariance  $\Sigma_{Syes}$ 

Use the likelihood of the parent state and the subsequent split states to determine which question a node should be split on

# **Example: Phonetic Decision Tree (DT)**

One tree is constructed for each state of each phone to cluster all the corresponding triphone states



# For an unseen triphone at test time

- Transition Matrix:
  - •
- State observation densities: •
  - of the decision tree
  - that leaf

All triphones of a given phoneme use the same transition matrix common to all triphones of a phoneme

Use the triphone identity to traverse all the way to a leaf

Use the state observation probabilities associated with

#### That's a wrap on HMM-based acoustic models



#### **DNN-based acoustic models?**





Can we use deep neural networks instead of HMMs to learn mappings between acoustics and phones?



#### **Brief Introduction to Neural Networks**

## **Feed-forward Neural Network**



#### Input Layer

#### **Feed-forward Neural Network Brain Metaphor**



Image from: <a href="https://upload.wikimedia.org/wikipedia/commons/1/10/Blausen\_0657\_MultipolarNeuron.png">https://upload.wikimedia.org/wikipedia/commons/1/10/Blausen\_0657\_MultipolarNeuron.png</a>

#### **Feed-forward Neural Network Parameterized Model**



 $a_5 = g(w_{35} \cdot a_3 + w_{45} \cdot a_4)$ 

If x is a 2-dimensional vector and the layer above it is a 2-dimensional vector h, a fully-connected layer is associated with:

h = xW + b

where w<sub>ij</sub> in W is the weight of the connection between i<sup>th</sup> neuron in the input row and j<sup>th</sup> neuron in the first hidden layer and b is the bias vector

Parameters of the network: all w<sub>ii</sub> (and biases not shown here)

 $W_{45} \cdot (g(W_{14} \cdot a_1 + W_{24} \cdot a_2)))$ 

#### **Feed-forward Neural Network Parameterized Model**



 $a_5 = g(w_{35} \cdot a_3 + w_{45} \cdot a_4)$ 

The simplest neural network is the perceptron:

A 1-layer feedforward neural network has the form:

 $MLP(x) = g(xW_1 + b_1) W_2 + b_2$ 

Perceptron(x) = xW + b

# **Common Activation Functions (g)**

Sigmoid:  $\sigma(x) = 1/(1 + e^{-x})$ 



nonlinear activation functions



# **Common Activation Functions (g)**

Sigmoid:  $\sigma(x) = 1/(1 + e^{-x})$ Hyperbolic tangent (tanh): tanh



$$h(x) = (e^{2x} - 1)/(e^{2x} + 1)$$

# **Common Activation Functions (g)**

Sigmoid:  $\sigma(x) = 1/(1 + e^{-x})$ Hyperbolic tangent (tanh): tanh Rectified Linear Unit (ReLU): F



$$n(x) = (e^{2x} - 1)/(e^{2x} + 1)$$
  
RELU(x) = max(0, x)

#### nonlinear activation functions

# **Optimization Problem**

- To train a neural network, define a loss function  $L(y,\tilde{y})$ : a function of the true output y and the predicted output  $\tilde{y}$
- L(y, $\tilde{y}$ ) assigns a non-negative numerical score to the neural network's output,  $\tilde{y}$
- The parameters of the network are set to minimise L over the training examples (i.e. a sum of losses over different training samples)
- · L is typically minimised using a gradient-based method

## **Stochastic Gradient Descent (SGD)**

#### SGD Algorithm

Inputs: Function NN(x;  $\theta$ ), Training examples,  $x_1 \dots x_n$  and outputs,  $y_1 \dots y_n$  and Loss function L.

```
do until stopping criterion
     Pick a training example x<sub>i</sub>, y<sub>i</sub>
     Compute the loss L(NN(x_i; \theta), y_i)
     Compute gradient of L, \nablaL with respect to \theta
     \theta \leftarrow \theta - \eta \nabla L
```

done

Return: θ

## **Training a Neural Network**

Define the Loss function to be minimised as a node L Goal: Learn weights for the neural network which minimise L  $w \leftarrow w - \eta \partial L / \partial w$ 

How do we efficiently compute  $\partial L/\partial w$  for all w?

Will compute  $\partial L/\partial u$  for every node u in the network!

 $\partial L/\partial w = \partial L/\partial u \cdot \partial u/\partial w$  where *u* is the node which uses *w* 

- Gradient Descent: Find  $\partial L/\partial w$  for every weight w, and update it as

### **Training a Neural Network**

New goal: compute  $\partial L/\partial u$  for every node u in the network

Simple algorithm: Backpropagation

Key fact: Chain rule of differentiation

depend (partially) on another variable *u*, then

- If L can be written as a function of variables  $v_1, \ldots, v_n$ , which in turn
  - $\partial L/\partial u = \sum_i \partial L/\partial v_i \cdot \partial v_i / \partial u$

# Backpropagation

If *L* can be written as a function of variables  $v_1, ..., v_n$ , which in turn depend (partially) on another variable *u*, then  $\partial L/\partial u = \sum_i \partial L/\partial v_i \cdot \partial v_i/\partial u$ 

Consider  $v_1, \ldots, v_n$  as the layer above  $u, \Gamma(u)$ 



 $\partial L/\partial u = \sum_{v \in \Gamma(u)} \partial L/\partial v \cdot \partial v/\partial u$ 



# Backpropagation

**Backpropagation** Base case:  $\partial L / \partial L = 1$ For each *u* (top to bottom): For each  $v \in \Gamma(u)$ : Inductively, have computed  $\partial L / \partial v$ Directly compute  $\partial v / \partial u$ Compute  $\partial L/\partial u$ Compute  $\partial L / \partial w$ where  $\partial L/\partial w = \partial L/\partial u \cdot \partial u/\partial w \ll$ 

#### $\partial L/\partial u = \sum_{v \in \Gamma(u)} \partial L/\partial v \cdot \partial v/\partial u$



#### Forward Pass

First, in a forward pass, compute values of all nodes given an input (The values of each node will be needed during backprop)

Where values computed in the forward pass may be needed

# **History of Neural Networks in ASR**

- Neural networks for speech recognition were explored as early as 1987
- Deep neural networks for speech
  - Beat state-of-the-art on the TIMIT corpus [M09]
  - Significant improvements shown on large-vocabulary systems [D11]
  - Dominant ASR paradigm [H12]

[H12] G. Hinton, et al., "Deep Neural Networks for Acoustic Modeling in Speech Recognition", IEEE Signal Processing Magazine, 2012.

[M09] A. Mohamed, G. Dahl, and G. Hinton, "Deep belief networks for phone recognition," NIPS Workshop on Deep Learning for Speech Recognition, 2009. [D11] G. Dahl, D. Yu, L. Deng, and A. Acero, "Context-Dependent Pre-Trained Deep Neural Networks for Large-Vocabulary Speech Recognition," TASL 20(1



## What's new?

- Why have NN-based systems come back to prominence? •
- Important developments •
  - Vast quantities of data available for ASR training
  - Fast GPU-based training
  - Improvements in optimization/initialization techniques •
  - Deeper networks enabled by fast training •
  - Larger output spaces enabled by fast training and • availability of data