

Speaker Adaptation

Steve Renals

Automatic Speech Recognition – ASR Lecture 14
3 March 2016

- **Speaker independent** (SI) systems have long been the focus for research in transcription, dialogue systems, etc.
- **Speaker dependent** (SD) systems can result in word error rates 2–3 times lower than SI systems (given the same amount of training data)
- A **Speaker adaptive** (SA) system... we would like
 - Error rates similar to SD systems
 - Building on an SI system
 - Requiring only a small fraction of the speaker-specific training data used by an SD system

Speaker-specific variation

- Acoustic model

- Speaking styles
- Accents
- Speech production anatomy (eg length of the vocal tract)

Also non-speaker variation, such as channel conditions (telephone, reverberant room, close talking mic) and application domain

Speaker adaptation of acoustic models aims to reduce the mismatch between test data and the models

Speaker-specific variation

- **Acoustic model**
 - Speaking styles
 - Accents
 - Speech production anatomy (eg length of the vocal tract)

Also non-speaker variation, such as channel conditions (telephone, reverberant room, close talking mic) and application domain

Speaker adaptation of acoustic models aims to reduce the mismatch between test data and the models

- **Pronunciation model**: speaker-specific, consistent change in pronunciation
- **Language model**: user-specific documents (exploited in personal dictation systems)

Modes of adaptation

- **Supervised or unsupervised**
 - Supervised: the word level transcription of the adaptation data is known (and HMMs may be constructed)
 - Unsupervised: the transcription must be estimated (eg using recognition output)

Modes of adaptation

- **Supervised or unsupervised**
 - Supervised: the word level transcription of the adaptation data is known (and HMMs may be constructed)
 - Unsupervised: the transcription must be estimated (eg using recognition output)
- **Static or dynamic**
 - Static: All adaptation data is presented to the system in a block before the final system is estimated (eg as used in enrollment in a dictation system)
 - Dynamic: Adaptation data is incrementally available, and models must be adapted before all adaptation data is available (eg as used in a spoken dialogue system)

Approaches to adaptation

- **Model based:** Adapt the parameters of the acoustic models to better match the observed data
 - Maximum a posteriori (MAP) adaptation of HMM/GMM parameters
 - Maximum likelihood linear regression (MLLR) of Gaussian parameters
 - Learning Hidden Unit Contributions (LHUC) for neural networks
- **Speaker normalization:** Normalize the acoustic data to reduce mismatch with the acoustic models
 - Vocal Tract Length Normalization (VTLN)
 - Constrained MLLR (cMLLR) — model-based normalisation!
- **Speaker space:** Estimate multiple sets of acoustic models, characterizing new speakers in terms of these model sets
 - Cluster-adaptive training
 - Eigenvoices
 - Speaker codes

Approaches to adaptation

- **Model based:** Adapt the parameters of the acoustic models to better match the observed data
 - **Maximum a posteriori (MAP) adaptation of HMM/GMM parameters**
 - **Maximum likelihood linear regression (MLLR) of Gaussian parameters**
 - **Learning Hidden Unit Contributions (LHUC) for neural networks**
- **Speaker normalization:** Normalize the acoustic data to reduce mismatch with the acoustic models
 - Vocal Tract Length Normalization (VTLN)
 - **Constrained MLLR (cMLLR) — model-based normalisation!**
- **Speaker space:** Estimate multiple sets of acoustic models, characterizing new speakers in terms of these model sets
 - Cluster-adaptive training
 - Eigenvoices
 - **Speaker codes**

Desirable properties for speaker adaptation

- **Compact:** relatively few speaker-dependent parameters
- **Unsupervised:** does not require labelled adaptation data, or changes to the training
- **Efficient:** low computational requirements
- **Flexible:** applicable to different model variants

Model-based adaptation: The MAP family

- **Basic idea** Use the SI models as a prior probability distribution over model parameters when estimating using speaker-specific data
- Theoretically well-motivated approach to incorporating the knowledge inherent in the SI model parameters
- Maximum likelihood (ML) training sets the model parameters λ to maximize the likelihood $p(\mathbf{X} | \lambda)$
- Maximum a posteriori (MAP) training maximizes the posterior of the parameters given the data:

$$p(\lambda | \mathbf{X}) \propto p(\mathbf{X} | \lambda)p_0(\lambda)$$

$p_0(\lambda)$ is the prior distribution of the parameters

- The use of a prior distribution, based on the SI models, means that less data is required to estimate the speaker-specific models: we are not starting from complete ignorance

Recall: ML estimation of GMM/HMM

- The mean of the m th Gaussian component of the j th state is estimated using a weighted average

$$\boldsymbol{\mu}_{mj} = \frac{\sum_n \gamma_{jm}(n) \mathbf{x}_n}{\sum_n \gamma_{jm}(n)}$$

- Where $\sum_n \gamma_{jm}(n)$ is the component occupation probability
- The covariance of the Gaussian component is given by:

$$\boldsymbol{\Sigma}_{mj} = \frac{\sum_n \gamma_{jm}(n) (\mathbf{x}_n - \boldsymbol{\mu}_{jm})(\mathbf{x}_n - \boldsymbol{\mu}_{jm})^T}{\sum_n \gamma_{jm}(n)}$$

MAP estimation

- What is $p_0(\boldsymbol{\lambda})$?
- Conjugate prior: the prior distribution has the same form as the posterior. There is no simple conjugate prior for GMMs, but an intuitively understandable approach may be employed.
- If the prior mean is $\boldsymbol{\mu}_0$, then the MAP estimate for the adapted mean $\hat{\boldsymbol{\mu}}$ of Gaussian is given by:

$$\hat{\boldsymbol{\mu}} = \frac{\tau \boldsymbol{\mu}_0 + \sum_n \gamma(n) \mathbf{x}_n}{\tau + \sum_n \gamma(n)}$$

- τ is a *hyperparameter* that controls the balance between the ML estimate of the mean, its prior value. Typically τ is in the range 2–20
- \mathbf{x}_n is the adaptation vector at time n
- $\gamma(n)$ the probability of this Gaussian at this time
- As the amount of training data increases, so the MAP estimate converges to the ML estimate

- **Basic idea** The main drawback to MAP adaptation is that it is local
- Only the parameters belonging to Gaussians of observed states will be adapted
- Large vocabulary speech recognition systems have about 10^5 Gaussians: most will not be adapted
 - Structural MAP (SMAP) approaches have been introduced to share Gaussians
 - The **MLLR** family of adaptation approaches addresses this by assuming that transformations for a specific speaker are systematic across Gaussians, states and models
- MAP adaptation is very useful for domain adaptation:
 - Example: adapting a conversational telephone speech system (100s of hours of data) to multiparty meetings (10s of hours of data) works well with MAP

The MLLR family

- **Basic idea** Rather than directly adapting the model parameters, estimate a transform which may be applied to the Gaussian means and covariances
- Linear transform applied to parameters of a set of Gaussians: adaptation transform parameters are shared across Gaussians
- This addresses the locality problem arising in MAP adaptation, since each adaptation data point can affect many of (or even all) the Gaussians in the system
- There are relatively few adaptation parameters, so estimation is robust
- Maximum Likelihood Linear Regression (MLLR) is the best known linear transform approach to speaker adaptation

MLLR: Maximum Likelihood Linear Regression

- MLLR is the best known linear transform approach to speaker adaptation
- Affine transform of mean parameters

$$\hat{\boldsymbol{\mu}} = \mathbf{A}\boldsymbol{\mu} + \mathbf{b}$$

If the observation vectors are d -dimension, then \mathbf{A} is a $d \times d$ matrix and \mathbf{b} is d -dimension vector

- If we define $\mathbf{W} = [\mathbf{bA}]$ and $\boldsymbol{\eta} = [1\boldsymbol{\mu}^T]^T$, then we can write:

$$\hat{\boldsymbol{\mu}} = \mathbf{W}\boldsymbol{\eta}$$

- In MLLR, \mathbf{W} is estimated so as to maximize the likelihood of the adaptation data
- A single transform \mathbf{W} can be shared across a set of Gaussian components (even all of them!)

Regression classes

- The number of transforms may obtained automatically
- A set of Gaussian components that share a transform is called a regression class
- Obtain the regression classes by constructing a *regression class tree*
- Each node in the tree represents a regression class sharing a transform
- For an adaptation set, work down the tree until arriving at the most specific set of nodes for which there is sufficient data
- Regression class tree constructed in a similar way to state clustering tree
- In practice the number of regression may be very small: one per context-independent phone class, one per broad class, or even just two (speech/non-speech)

Estimating the transforms

- The linear transformation matrix W is obtained by finding its setting which optimizes the log likelihood
- **Mean adaptation**: Log likelihood

$$L = \sum_r \sum_n \gamma_r(n) \log \left(K_r \exp \left(-\frac{1}{2} (\mathbf{x}_n - \mathbf{W}\boldsymbol{\eta}_r)^T \boldsymbol{\Sigma}_r^{-1} (\mathbf{x}_n - \mathbf{W}\boldsymbol{\eta}_r) \right) \right)$$

where r ranges over the components belonging to the regression class

- Differentiating L and setting to 0 results in an equation for \mathbf{W} : there is no closed form solution if $\boldsymbol{\Sigma}$ is full covariance; can be solved if $\boldsymbol{\Sigma}$ is diagonal (but requires a matrix inversion)
- Variance adaptation is also possible
- See Gales and Woodland (1996), Gales (1998) for details

- Mean-only MLLR results in 10–15% relative reduction in WER
- Few regression classes and well-estimated transforms work best in practice
- Robust adaptation available with about 1 minute of speech; performance similar to SD models available with 30 minutes of adaptation data
- Such linear transforms can account for any systematic (linear) variation from the speaker independent models, for example those caused by channel effects.

Constrained MLLR (cMLLR)

- **Basic idea** use the same linear transform for both mean and covariance

$$\hat{\boldsymbol{\mu}} = \mathbf{A}'\boldsymbol{\mu} - \mathbf{b}'$$
$$\hat{\boldsymbol{\Sigma}} = \mathbf{A}'\boldsymbol{\Sigma}\mathbf{A}'^T$$

- No closed form solution but can be solved iteratively
- Log likelihood for cMLLR

$$L = \mathcal{N}(\mathbf{A}\mathbf{x}_n + \mathbf{b}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \log(|\mathbf{A}|) \quad \mathbf{A}' = \mathbf{A}^{-1}; \mathbf{b}' = \mathbf{A}\mathbf{b}$$

Equivalent to applying the linear transform to the data!
Also called fMLLR (feature space MLLR)

- Iterative solution amenable to online/dynamic adaptation, by using just one iteration for each increment
- Similar improvement in accuracy to standard MLLR

Speaker-adaptive training (SAT)

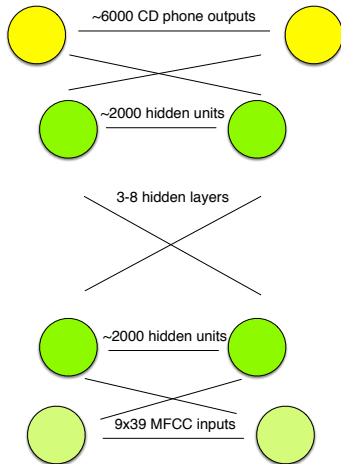
- **Basic idea** Rather than SI seed (canonical) models, construct models designed for adaptation
- Estimate parameters of canonical models by training MLLR mean transforms for each training speaker
- Train using the MLLR transform for each speaker; interleave Gaussian parameter estimation and MLLR transform estimation
- SAT results in much higher training likelihoods, and improved recognition results
- But: increased training complexity and storage requirements
- SAT using cMLLR, corresponds to a type of speaker normalization at training time

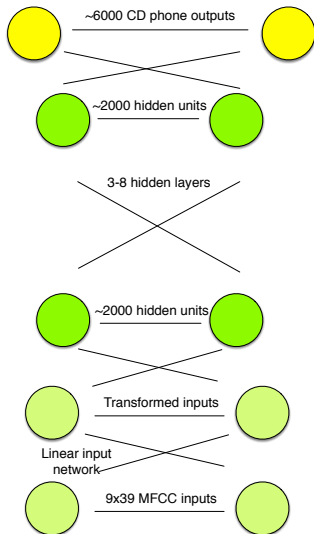
Speaker adaptation in hybrid HMM/NN systems: CMLLR feature transformation

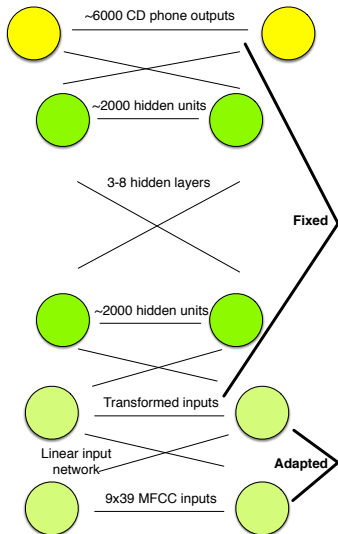
- **Basic idea:** If HMM/GMM system is used to estimate a single constrained MLLR adaptation transform, this can be viewed as a feature space transform
- Use the HMM/GMM system with the same tied state space to estimate a single CMLLR transform for a given speaker, and use this to transform the input speech to the DNN for the target speaker
- Can operate unsupervised (since the GMM system estimates the transform)
- Limited to a single transform (regression class)

Speaker adaptation in hybrid HMM/NN systems: LIN – Linear Input Network

- **Basic idea:** single linear input layer trained to map input speaker-dependent speech to speaker-independent network
- Training: linear input network (LIN) can either be fixed as the identity or (adaptive training) be trained along with the other parameters
- Testing: freeze the main (speaker-independent) network and propagate gradients for speech from the target speaker to the LIN, which is updated — linear transform learned for each speaker
- Requires supervised training data

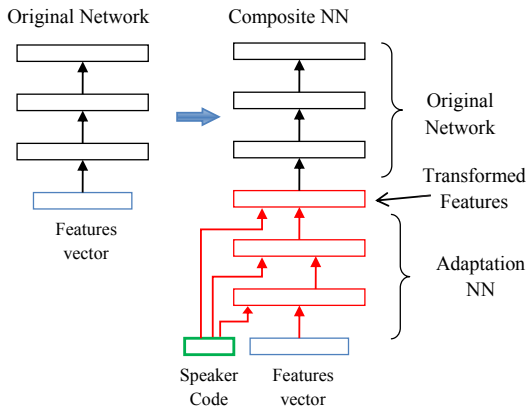






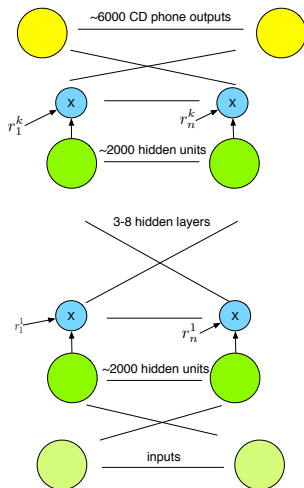
Speaker adaptation in hybrid HMM/NN systems: Speaker codes

- **Basic idea:** Learn a short speaker code vector for each talker

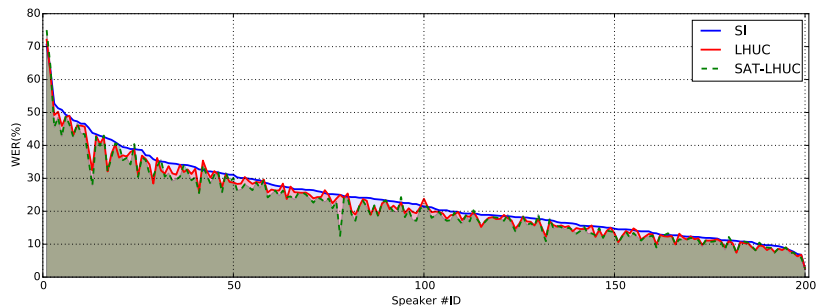


Speaker adaptation in hybrid HMM/NN systems: LHUC – Learning Hidden Unit Contributions

- **Basic idea:** Add a learnable speaker dependent amplitude to each hidden unit
- Speaker independent: amplitudes set to 1
- Speaker dependent: learn amplitudes from data, per speaker

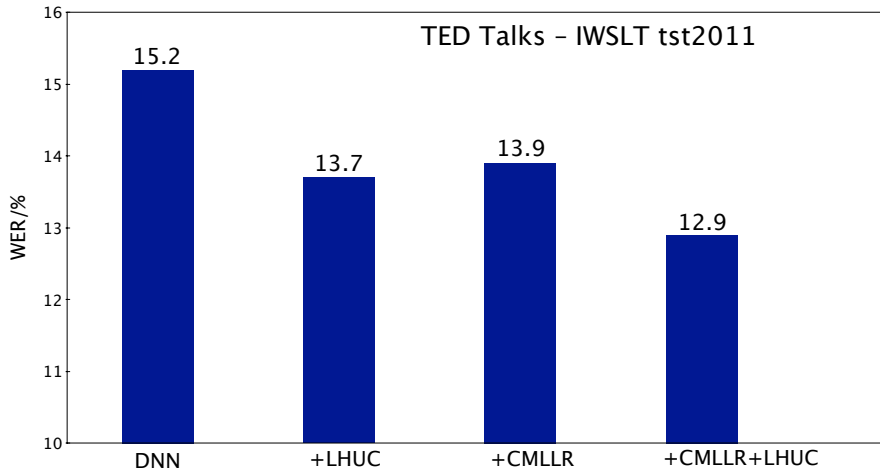


LHUC adaptation by speaker



Results on speakers across AMI, TED, Switchboard corpora

Speaker adaptation in hybrid HMM/NN systems: Experimental Results on TED



Speaker Adaptation

- One of the most intensive areas of speech recognition research since the early 1990s
- HMM/GMM
 - Substantial progress, resulting in significant, additive, consistent reductions in word error rate
 - Close mathematical links between different approaches
 - Linear transforms at the heart of many approaches
- HMM/NN
 - Open research topic
 - GMM-based feature space transforms somewhat effective
 - Direct weight adaptation less effective

- Gales and Young (2007), “The Application of Hidden Markov Models in Speech Recognition”, *Foundations and Trends in Signal Processing*, **1** (3), 195–304: section 5.
<http://mi.eng.cam.ac.uk/~sjy/papers/gayo07.pdf>
- Woodland (2001), “Speaker adaptation for continuous density HMMs: A review”, ISCA ITRW on Adaptation Methods for Speech Recognition.
http://www.isca-speech.org/archive_open/archive_papers/adaptation/adap_011.pdf
- Liao (2013), “Speaker adaptation of context dependent deep neural networks”, ICASSP-2013.
<http://dx.doi.org/10.1109/ICASSP.2013.6639212>
- Abdel-Hamid and Jiang (2013), “Fast speaker adaptation of hybrid NN/HMM model for speech recognition based on discriminative learning of speaker code”, ICASSP-2013.
<http://dx.doi.org/10.1109/ICASSP.2013.6639211>
- Swietojanski and Renals (2014), “Learning Hidden Unit Contributions for Unsupervised Speaker Adaptation of Neural Network Acoustic Models”, SLT-2014. <http://www.cstr.inf.ed.ac.uk/downloads/publications/2014/ps-slt14.pdf>