

Deep Network for Speech Emotion Recognition

—A Study of Deep Learning—



Zhuowei Han

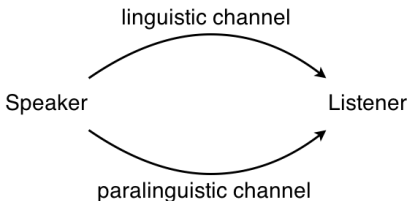
Institut für Signalverarbeitung
und Systemtheorie

Universität Stuttgart

16/04/2015

Speech Emotion Recognition

- Most current work focuses on speech processing based on linguistic information, e.g.: Skype Translator
- More natural human-machine interaction requires paralinguistic information such as age, gender, emotion.
- Speech Recognition / Speaker Identification / Emotion Recognition



Deep Learning

- Deep architecture for extracting complex structure and building internal representations from input
- New research area of machine learning (from shallow to deep structure)
- Widely applied in vision/audition processing, e.g. handwriting recognition (Graves, Alex, et al. 2009), traffic sign classification (Schmidhuber, et al. 2011), text translation (Google, 2014)

Foundations

- Mel Frequency Cepstral Features
- Emotion Recognition Approaches

Conditional Restricted Boltzmann Machine

- Restricted Boltzmann Machine
- CRBM

Conclusion and Outlook

Foundations

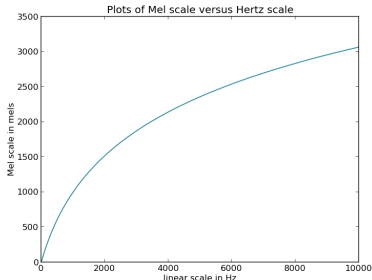
- Mel Frequency Cepstral Features
- Emotion Recognition Approaches

Conditional Restricted Boltzmann Machine

- Restricted Boltzmann Machine
- CRBM

Conclusion and Outlook

- short-term power spectrum
- mel-scale approximate human perception
- widely-used in speech recognition tasks
- Transformation between Mel and Hertz scale



$$f_{mel} = 1125 \ln (1 + f_{Hz}/700)$$

$$f_{Hz} = 700 (\exp(f_{mel}/1125) - 1)$$

Traditional Approaches

- pre-selected features
- supervised training
- low-level features not appropriate for classification
- shallow structure of classifiers

Deep Learning Approaches

- learning representations from high-dim data
- extracting appropriate features without hand-crafting
- low-level features are used to build high-level features as network gets deeper
- frame-based classification

Foundations

Mel Frequency Cepstral Features

Emotion Recognition Approaches

Conditional Restricted Boltzmann Machine

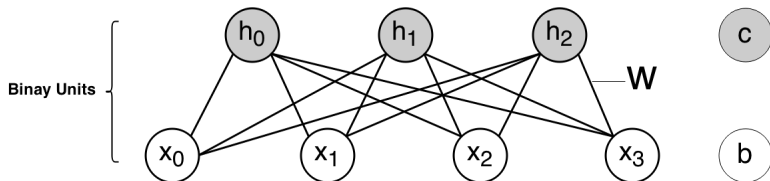
Restricted Boltzmann Machine

CRBM

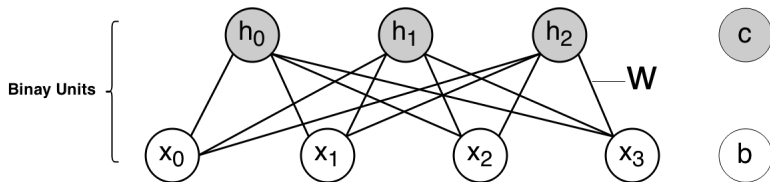
Conclusion and Outlook

- Generative graphical model, capture data distribution $P(\mathbf{x}|\theta)$
- Trained in unsupervised way, only use unlabeled input sequences \mathbf{x} for learning.
 - automatically extract useful features from data
 - Find hidden structure (distribution).
 - Learned features used for prediction or classification
- Successfully applied in motion capture (Graham W. Taylor, Geoffrey E. Hinton, 2006)
- Potential to be extend to capture temporal information

Structure



Structure



$$\text{Energy Function: } E_{\theta} = -\mathbf{x}^T \mathbf{W} \mathbf{h} - \mathbf{b}^T \mathbf{x} - \mathbf{c}^T \mathbf{h}$$

$$\text{Joint Distribution: } P^{RBM}(\mathbf{x}, \mathbf{h}) = \frac{1}{Z} e^{-E_{\theta}(\mathbf{x}, \mathbf{h})}$$

$$\text{Partition Function: } Z = \sum_{\mathbf{x}, \mathbf{h}} e^{-E_{\theta}(\mathbf{x}, \mathbf{h})}$$

$$\text{Free Energy: } \mathcal{F}(\mathbf{x}) = -\log \sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}$$

Inference

$$P(\mathbf{x}) = \sum_{\mathbf{h}} P(\mathbf{x}, \mathbf{h})$$

$$P(\mathbf{h}) = \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{h})$$

Inference

$$P(\mathbf{x}) = \sum_{\mathbf{h}} P(\mathbf{x}, \mathbf{h})$$

$$P(\mathbf{h}) = \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{h})$$

$$P(\mathbf{h}|\mathbf{x}) = \frac{P(\mathbf{x}, \mathbf{h})}{P(\mathbf{x})}$$

$$P(\mathbf{x}|\mathbf{h}) = \frac{P(\mathbf{x}, \mathbf{h})}{P(\mathbf{h})}$$

Inference

$$P(\mathbf{x}) = \sum_{\mathbf{h}} P(\mathbf{x}, \mathbf{h})$$

$$P(\mathbf{h}) = \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{h})$$

$$P(\mathbf{h}|\mathbf{x}) = \frac{P(\mathbf{x}, \mathbf{h})}{P(\mathbf{x})}$$

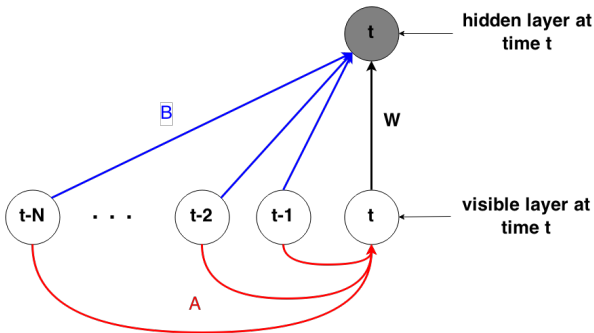
$$P(\mathbf{x}|\mathbf{h}) = \frac{P(\mathbf{x}, \mathbf{h})}{P(\mathbf{h})}$$

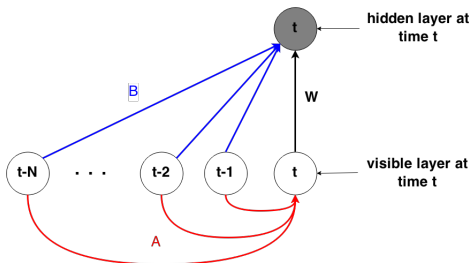
$$P(h_j = 1 \mid \mathbf{x}) = \text{sigmoid}(\sum_i x_i W_{ij} + c_j)$$

$$P(x_i = 1 \mid \mathbf{h}) = \text{sigmoid}(\sum_j W_{ij} h_j + b_i)$$

- Consider visible units from previous time step as additional bias for current visible and hidden layer
- A and B are weight parameter of visible (history) - visible and visible (history) - hidden connections
- Visible layer is linear units with independent Gaussian noise to model real-valued data, e.g. spectral features

- Consider visible units from previous time step as additional bias for current visible and hidden layer
- A and B are weight parameter of visible (history) - visible and visible (history) - hidden connections
- Visible layer is linear units with independent Gaussian noise to model real-valued data, e.g. spectral features





$$\text{Energy Function: } E_{\theta}^{CRBM}(\mathbf{x}, \mathbf{h}) = \left\| \frac{\mathbf{x} - \tilde{\mathbf{b}}}{2} \right\|^2 - \tilde{\mathbf{c}}^T \mathbf{h} - \mathbf{x}^T \mathbf{W} \mathbf{h}$$

$$\tilde{\mathbf{b}} = \mathbf{b} + \mathbf{A} \cdot \mathbf{x}_{<t}$$

$$\tilde{\mathbf{c}} = \mathbf{c} + \mathbf{B} \cdot \mathbf{x}_{<t}$$

$$\theta = \{\mathbf{W}, \mathbf{A}, \mathbf{B}, \mathbf{b}, \mathbf{c}\}$$

$$\text{Free Energy: } \mathcal{F}(\mathbf{x}) = \left\| \mathbf{x} - \tilde{\mathbf{b}} \right\|^2 - \log(1 + e^{\tilde{\mathbf{c}} + \mathbf{x} \cdot \mathbf{W}})$$

$$\text{Energy Function: } E_{\theta}^{CRBM}(\mathbf{x}, \mathbf{h}) = \left\| \frac{\mathbf{x} - \tilde{\mathbf{b}}}{2} \right\|^2 - \tilde{\mathbf{c}}^T \mathbf{h} - \mathbf{x}^T \mathbf{W} \mathbf{h}$$

$$\text{Free Energy: } \mathcal{F}(\mathbf{x}) = \left\| \mathbf{x} - \tilde{\mathbf{b}} \right\|^2 - \log(1 + e^{\tilde{\mathbf{c}} + \mathbf{x} \cdot \mathbf{W}})$$

$$\tilde{\mathbf{b}} = \mathbf{b} + \mathbf{A} \cdot \mathbf{x}_{<t}$$

$$\tilde{\mathbf{c}} = \mathbf{c} + \mathbf{B} \cdot \mathbf{x}_{<t}$$

$$\theta = \{\mathbf{W}, \mathbf{A}, \mathbf{B}, \mathbf{b}, \mathbf{c}\}$$

Maximum Likelihood Estimation $P(\mathbf{x}|\boldsymbol{\theta})$

Note that KL is non-negative

Maximum Likelihood Estimation $P(\mathbf{x}|\boldsymbol{\theta})$

Kullback-Leibler Divergence:

$$\begin{aligned} KL(Q(\mathbf{x})||P(\mathbf{x}|\boldsymbol{\theta})) &= \int_{-\infty}^{\infty} Q(\mathbf{x}) \cdot \log \frac{Q(\mathbf{x})}{P(\mathbf{x}|\boldsymbol{\theta})} d\mathbf{x} \\ &= \int_{-\infty}^{\infty} Q(\mathbf{x}) \cdot \log Q(\mathbf{x}) d\mathbf{x} - \int_{-\infty}^{\infty} Q(\mathbf{x}) \cdot \log P(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} \\ &= \langle \log Q(\mathbf{x}) \rangle_{Q(\mathbf{x})} - \langle \log P(\mathbf{x}|\boldsymbol{\theta}) \rangle_{Q(\mathbf{x})} \end{aligned}$$

$Q(\mathbf{x})$, true data distribution

$P(\mathbf{x}|\boldsymbol{\theta})$, model distribution

$\langle \cdot \rangle_{Q(\mathbf{x})}$, expectation w.r.t. $Q(\mathbf{x})$

Note that KL is non-negative

$$-\log P(\mathbf{x}|\boldsymbol{\theta}) = \mathcal{F}(\mathbf{x}) + \log \sum_{\mathbf{x}} \sum_{\mathbf{h}} e^{-E_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{h})} \quad \text{Free Energy}$$

$$-\log P(\mathbf{x}|\boldsymbol{\theta}) = \mathcal{F}(\mathbf{x}) + \log \sum_{\mathbf{x}} \sum_{\mathbf{h}} e^{-E_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{h})} \quad \text{Free Energy}$$

$$-\frac{\partial \log P(\mathbf{x})}{\partial \boldsymbol{\theta}} = \frac{\partial \mathcal{F}(\mathbf{x})}{\partial \boldsymbol{\theta}} - \sum_{\tilde{\mathbf{x}}} P(\tilde{\mathbf{x}}) \frac{\partial \mathcal{F}(\tilde{\mathbf{x}})}{\partial \boldsymbol{\theta}} \quad \leftarrow \text{intractable!}$$

$$-\log P(\mathbf{x}|\boldsymbol{\theta}) = \mathcal{F}(\mathbf{x}) + \log \sum_{\mathbf{x}} \sum_{\mathbf{h}} e^{-E_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{h})}$$

Free Energy

$$-\frac{\partial \log P(\mathbf{x})}{\partial \boldsymbol{\theta}} = \frac{\partial \mathcal{F}(\mathbf{x})}{\partial \boldsymbol{\theta}} - \sum_{\tilde{\mathbf{x}}} P(\tilde{\mathbf{x}}) \frac{\partial \mathcal{F}(\tilde{\mathbf{x}})}{\partial \boldsymbol{\theta}}$$

← intractable!

$$-\frac{\partial \log P(\mathbf{x})}{\partial \boldsymbol{\theta}} = \frac{\partial \mathcal{F}(\mathbf{x})}{\partial \boldsymbol{\theta}} - \frac{1}{|\mathcal{N}|} \sum_{\tilde{\mathbf{x}} \in \mathcal{N}} P(\tilde{\mathbf{x}}) \frac{\partial \mathcal{F}(\tilde{\mathbf{x}})}{\partial \boldsymbol{\theta}}$$

sampling

t steps-Gibbs sampling

$$\mathbf{x}_1 \sim \hat{P}(\mathbf{x})$$

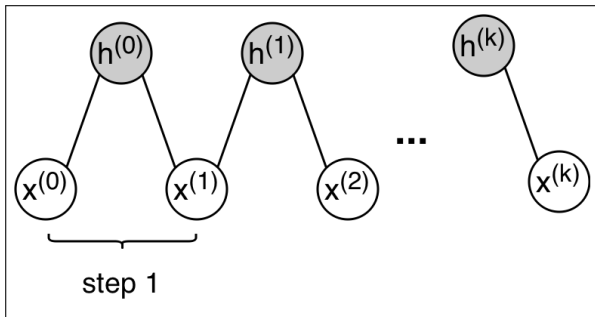
$$\mathbf{h}_1 \sim \hat{P}(\mathbf{h}|\mathbf{x}_1)$$

$$\mathbf{x}_2 \sim \hat{P}(\mathbf{x}|\mathbf{h}_1)$$

$$\mathbf{h}_2 \sim \hat{P}(\mathbf{h}|\mathbf{x}_2)$$

⋮

$$\mathbf{x}_{t+1} \sim \hat{P}(\mathbf{x}|\mathbf{h}_t)$$



t steps-Gibbs sampling

$$\mathbf{x}_1 \sim \hat{P}(\mathbf{x})$$

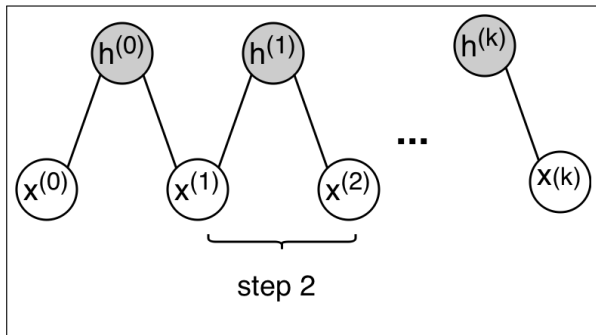
$$\mathbf{h}_1 \sim \hat{P}(\mathbf{h}|\mathbf{x}_1)$$

$$\mathbf{x}_2 \sim \hat{P}(\mathbf{x}|\mathbf{h}_1)$$

$$\mathbf{h}_2 \sim \hat{P}(\mathbf{h}|\mathbf{x}_2)$$

⋮

$$\mathbf{x}_{t+1} \sim \hat{P}(\mathbf{x}|\mathbf{h}_t)$$



t steps-Gibbs sampling

$$\mathbf{x}_1 \sim \hat{P}(\mathbf{x})$$

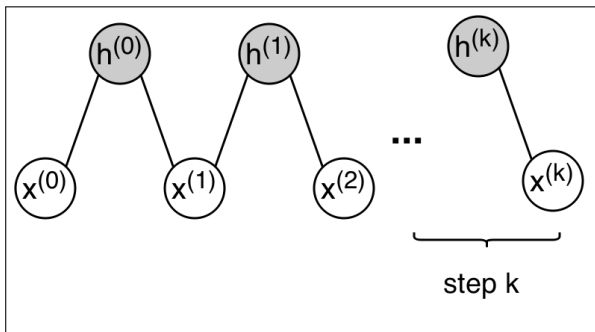
$$\mathbf{h}_1 \sim \hat{P}(\mathbf{h}|\mathbf{x}_1)$$

$$\mathbf{x}_2 \sim \hat{P}(\mathbf{x}|\mathbf{h}_1)$$

$$\mathbf{h}_2 \sim \hat{P}(\mathbf{h}|\mathbf{x}_2)$$

\vdots

$$\mathbf{x}_{t+1} \sim \hat{P}(\mathbf{x}|\mathbf{h}_t)$$



- Performing k-Gibbs steps to generate $P_k(\mathbf{x}|\boldsymbol{\theta})$, approximation of model distribution
- Difference between approximation and true model distribution:

$$KL(P_k(\mathbf{x}|\boldsymbol{\theta})||P(\mathbf{x}|\boldsymbol{\theta}))$$

- Contrastive Divergence (CD):

$$KL(Q(\mathbf{x})||P(\mathbf{x}|\boldsymbol{\theta})) - KL(P_k(\mathbf{x}|\boldsymbol{\theta})||P(\mathbf{x}|\boldsymbol{\theta}))$$

- With enough steps the Markov chain converges to stationary distribution:

$$P_{k \rightarrow \infty}(\mathbf{x}|\boldsymbol{\theta}) = P(\mathbf{x}|\boldsymbol{\theta})$$

- CD-1 performs well in practice

Parameter Update

$$\Delta\theta \sim KL(P(\mathbf{x}|\theta) || P_k(\mathbf{x}|\theta))$$

Foundations

Mel Frequency Cepstral Features

Emotion Recognition Approaches

Conditional Restricted Boltzmann Machine

Restricted Boltzmann Machine

CRBM

Conclusion and Outlook

- Model with long-term dependencies shall be used for speech emotion
- CRBM is appropriate for short-term modelling, but not for long-term variation
- LSTM is good at modelling long time dependency
- Frame-based classification can also reach good result
 - CRBM-LSTM 71.98%
 - LSTM 81.59%
 - LSTM with rectifier layers 83.43%

- Stacking CRBM to form deeper structure
- Train CRBM with more/larger database
- Second order optimization to speed up learning process
- Bi-directional LSTM, capturing future dependencies

Thank You!