## Why 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

linguistic channel

Speaker → Listener

paralinguistic channel

## Deep Network Applications

- Handwriting Digit Recognition

- Image Recognition

## Why 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.

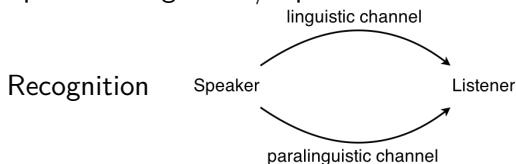- Speech Recognition / Speaker Identification / Emotion

Recognition

linguistic channel

Speaker → Listener

paralinguistic channel

## Deep Network Applications

- Handwriting Digit Recognition
- Image Recognition

# Table of Contents

## Table of Contents

# Mel Frequency Cepstral Features

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



Plots of Mel scale versus Hertz scale

$$f_{mel} = 1125 \ \ln \left(1 + f_{Hz}/700\right) \quad (1)$$
$$f_{Hz} = 700 \left(\exp(f_{mel}/1125) - 1\right) \quad (2)$$

# Emotion Recognition Approaches

## Traditional Approaches

- pre-selected features
- supervised training
- low-level features not appropriate for claasification
- 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 classfication

# Table of Contents

## Character

- Generative model
- Undirected graphical model, good at modeling high-dimensioanl data (speech emotion)
- Trained in unsupervised way, only use unlabeled input sequencex for learning.
  - □ automatically extract useful features from data
  - □ Find hidden structure (distribution).
  - □ Learned features used for prediction or classification
- Potential to be extend to capture temporal information.

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

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

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

## Inference

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

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

# Table of Contents

Computing net-activation

$$\underline{z}_k^{(l+1)} = \mathbf{W}^{(l)}\underline{a}_k^{(l)} + \underline{b}^{(l)}$$

$$\underline{a}_k^{(l+1)} = \underline{\Phi}\left(\underline{z}_k^{(l+1)}\right)$$

$$\hat{\underline{y}}_k = \underline{a}_k^{(ol)}$$

- Arbitrary non-linear mapping from $\underline{x}_k$ to $\hat{\underline{y}}_k$ possible
- Relation $N \Leftrightarrow$ Complexity
- Deep Architectures ($l \uparrow$) more efficient than shallow ones ($l \downarrow, N_l \uparrow$)

## Determining the parameters

### Training objective

$$J(\mathbf{W}, \underline{b}) = \sum_{\forall k} \frac{1}{2} ||\underline{y}_k - \hat{\underline{y}}_k||^2 + \frac{\lambda}{2} \sum_{\forall l} ||\mathbf{W}^{(l)}||_F^2 \qquad (3)$$

$$\mathbf{W}, \underline{b} = \arg\min_{\mathbf{W}, \underline{b}} J(\mathbf{W}, \underline{b}) \qquad (4)$$

### Numerical minimization

- Gradient calculation with Backpropagation
- Stochastic gradient descent
- **L**imited memory **B**royden-**F**letcher-**G**oldfarb-**S**hanno (L-BFGS)

## Problems

- Optimization problem non-convex
  $\Rightarrow$ getting stuck in poor local minima
- Diffusion of gradients
- Large p small n problem $\Rightarrow$ overfitting

- Layerwise Pre-training
- Layerwise Pre-training

# Table of Contents

# Concepts of RNN

- modelling sequential data, emotion in speech .

- Same Structure as MLP but differs from feed-forward network, enabling nonlinear mapping.

- Feedback connection between previous hidden units and current hidden units, enabling memory past hidden state.

- Potentially to model arbitrary dynamic system.

- Trained with backpropagation through time (BPTT)

## Recurrent Neural Network

### Concepts of RNN

- modelling sequential data, emotion in speech .
- Same Structure as MLP but differs from feed-forward network, enabling nonlinear mapping.
- Feedback connection between previous hidden units and current hidden units, enabling memory past hidden state.
- Potentially to model arbitary dynamic system.
- Trained with **b**ack**p**ropagation **t**hrough **t**ime (BPTT)

## Recurrent Neural Network

### Concepts of RNN

- modelling sequential data, emotion in speech .
- Same Structure as MLP but differs from feed-forward network, enabling nonlinear mapping.

$$h_t = \mathcal{H}(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$$
$$y_t = W_{hy}h_t + b_y$$

- Feedback connection between previous hidden units and current hidden units, enabling memory past hidden state.
- Potentially to model arbitary dynamic system.
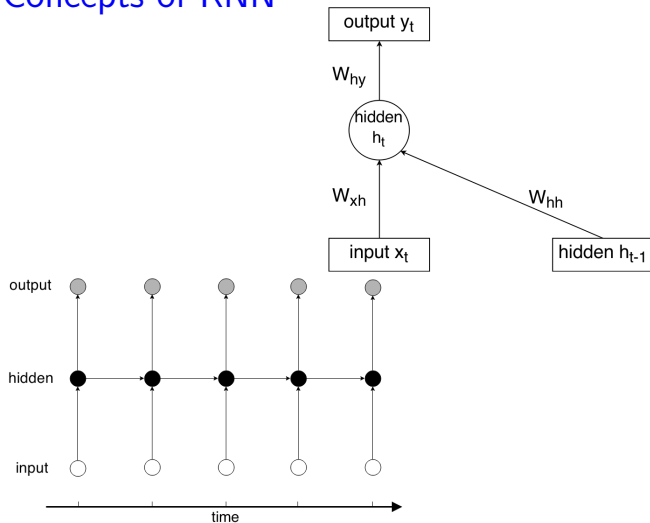- Trained with **b**ack**p**ropagation **t**hrough **t**ime (BPTT)

## Concepts of RNN

- modelling sequential data, emotion in speech .
- Same Structure as MLP but differs from feed-forward network, enabling nonlinear mapping.
- Feedback connection between previous hidden units and current hidden units, enabling memory past hidden state.
- Potentially to model arbitary dynamic system.
- Trained with **b**ack**p**ropagation **t**hrough **t**ime (BPTT)

## Recurrent Neural Network

### Concepts of RNN

- modelling sequential data, emotion in speech .
- Same Structure as MLP but differs from feed-forward network, enabling nonlinear mapping.
- Feedback connection between previous hidden units and current hidden units, enabling memory past hidden state.
- Potentially to model arbitary dynamic system.
- Trained with backpropagation through time (BPTT)

# Recurrent Neural Network

## Concepts of RNN

- modelling sequential data, emotion in speech .
- Same Structure as MLP but differs from feed-forward network, enabling nonlinear mapping.
- Feedback connection between previous hidden units and current hidden units, enabling memory past hidden state.
- Potentially to model arbitary dynamic system.
- Trained with **b**ack**p**ropagation **t**hrough **t**ime (BPTT)
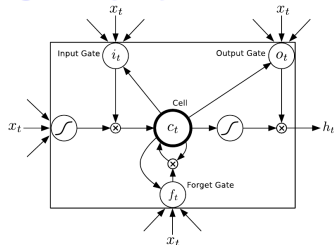
## Concepts of RNN

## Problems with RNN

- gradient vanishing during backpropagation as time steps increases ($>100$)
- difficult to capture long-time dependency (which is required in emotion recognition)

## Solutions

-

S. Hochreiter and J. Schmidhuber, Lovol. 9, pp. 1735-1780, 1997.

LSTM unit



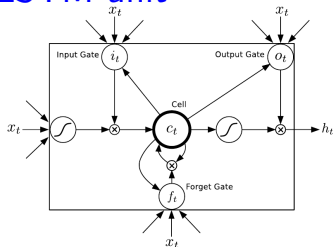$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i)$$
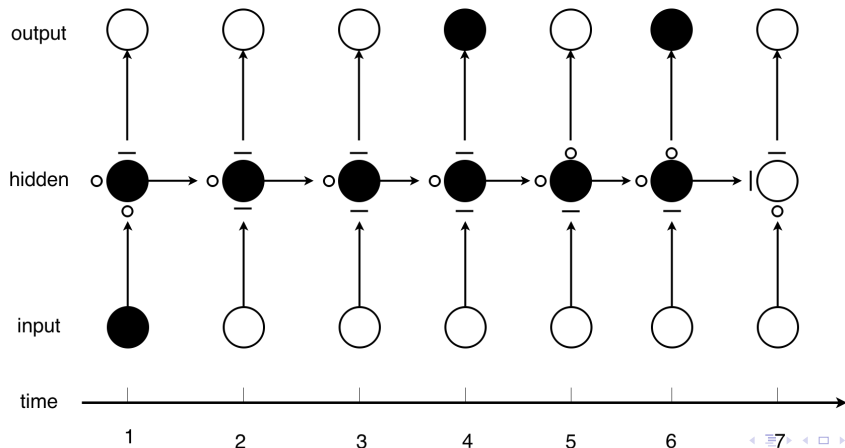$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f)$$
$$c_t = f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b)$$
$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o)$$
$$h_t = o_t \tanh(c_t)$$

## LSTM unit



$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i)$$

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f)$$

$$c_t = f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o)$$

$$h_t = o_t \tanh(c_t)$$

# Long short term memory

## Features in LSTM

- gates are trained to learn when it shoud be open/closed.
- Constant Error Carousel
- preserve long-time dependency by maintaining gradient over time.

# Table of Contents

# Table of Contents