Audio-Visual
Automatic
Speech
Recognition

Helge Reikeras

Introduction

Acoustic
speech

Visual speech

Modeling

Experimental
results

Conclusion

# Audio-Visual Automatic Speech Recognition

Helge Reikeras

June 30, 2010
SciPy 2010: Python for Scientific Computing Conference

- What?
  - Integration of audio and visual speech modalities with the purpose of enhancing speech recognition performance.
- Why?
  - McGurk effect (e.g. visual /ga/ combined with an audio /ba/ is heard as /da/)
  - Performance increase in noisy environments
  - Progress in speech recognition seems to be stagnating

- Example: YouTube automatic captions

- Mel-frequency cepstrum coefficients (MFCCs).
- Cosine transform of the logarithm of the short-term energy spectrum of a signal, expressed on the mel-frequency scale.
- The result is a set of coefficients that approximates the way the human auditory system perceives sound.

- Visual speech information mainly contained in the motion of visible articulators such as lips, tongue and jaw.

$$\mathbf{s} = \mathbf{s}_0 + \sum_{i=1}^{N} p_i \mathbf{s}_i. \qquad \text{(PCA)}$$

# Active appearance models (appearance) (3/3)

$$A(\mathbf{x}) = A_0 + \sum_{i=1}^{M} \lambda_i A_i(\mathbf{x}), \qquad \mathbf{x} \in \mathbf{s}_0. \qquad \text{(PCA)}$$

# Facial feature tracking (1/2)

Audio-Visual
Automatic
Speech
Recognition

Helge Reikeras

Introduction

Acoustic
speech

Visual speech

Modeling

Experimental
results

Conclusion

- Minimize difference between AAM and input image (warped onto the base shape $\mathbf{s}_0$).

- Warp is a piecewise affine transformation (triangulated base shape).

- Nonlinear least squares problem

$$\operatorname*{argmin}_{\boldsymbol{\lambda},\mathbf{p}} \sum_{\mathbf{x}\in\mathbf{s}_0} \left[ A_0(\mathbf{x}) + \sum_{i=1}^{M} \lambda_i A_i(\mathbf{X}) - I(\mathbf{W}(\mathbf{x};\mathbf{p})) \right]^2$$

- Solve using non-linear numerical optimization methods.

- Gaussian Mixture Models (GMMs) provide a powerful method for modeling data distributions.
- Weighted linear combination of Gaussian distributions.

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Data: $\mathbf{x}$
- Model parameters:
    - Weights $\boldsymbol{\pi}$
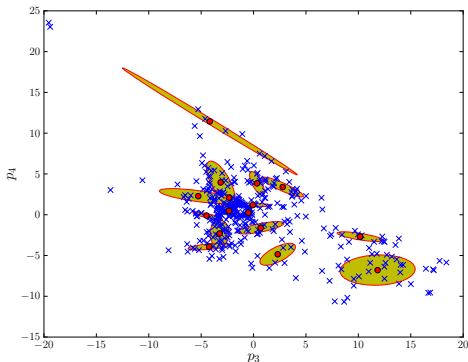    - Means $\boldsymbol{\mu}$
    - Covariances $\boldsymbol{\Sigma}$

- Log likelihood function gives the likelihood of the data $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n\}$ given GMM model parameters

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

- EM is an iterative algorithm for maximizing the log likelihood function w.r.t. GMM parameters.
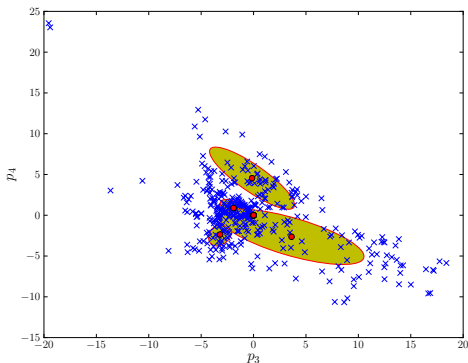
- Visual EM-GMM (16 mixture components)



(Note that in practice we use more than 2 dimensional feature vectors)

- How do we choose the number of Gaussian mixture components?
- VB differs from EM in that parameters are modeled as random variables.
- Suitable conjugate priors for GMM parameters are:
  - Weights; Dirichlet
  - Means: Gaussian
  - Covariances (precision): Wishart
- Avoids overfitting, singular solutions (when a Gaussian collapses onto a single data point) and leads to automatic model complexity selection.
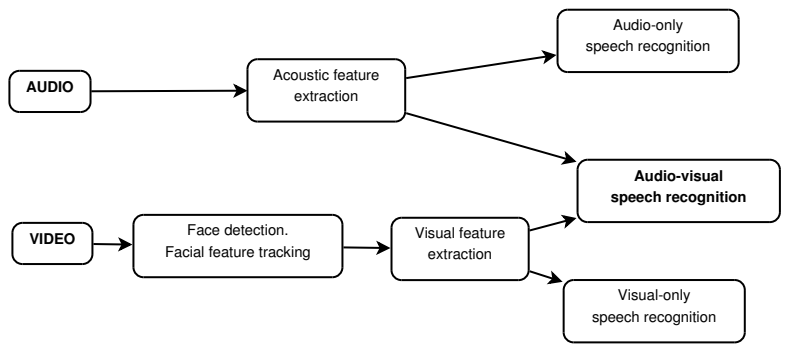
- Visual VB-GMM (16 mixture components)



- Remaining components have converged to their prior distributions and been assigned zero weights.

# Audio-visual fusion

- Acoustic GMM: $p(\mathbf{x}_A|c)$
- Visual GMM: $p(\mathbf{x}_V|c)$
- Classification (e.g. words or phonemes)
- Stream exponents $\lambda_A$, $\lambda_V$
- $\text{Score}(\mathbf{x}_{AV}|c) = p(\mathbf{x}_A|c)^{\lambda_A} p(\mathbf{x}_V|c)^{\lambda_V}$
- $0 \leq \lambda_A, \lambda_V \leq 1$
- $\lambda_A + \lambda_V = 1$
- Learn stream weights **discriminatively**.
- Minimize misclassification rate on development set.

# Python Implementation

Audio-Visual
Automatic
Speech
Recognition

Helge Reikeras

Introduction

Acoustic
speech

Visual speech

Modeling

Experimental
results

Conclusion

- Implemented in Python using SciPy (open source scientific computing Python library).

- Signal processing, computer vision and machine learning are active areas of development in the SciPy community.

- SciPy modules used:
  - scikits.talkbox.features.mfcc (MFCCs)
  - scikits.image (image processing)
  - scipy.optimize.fmin_ncg (facial feature tracking)
  - scipy.learn.em (EM)

- New modules developed as part of this research:
  - vb (VB inference)
  - aam (AAMs)

- Using the Clemson University audio-visual experiments (CUAVE) database.
- Contains video of 36 speakers, 19 male and 17 female, uttering isolated and connected digits in frontal, profile and while moving.

- Use separate training, development and test data sets (1/3, 1/3, 1/3).
- Add acoustic noise ranging from -5dB to 25 dB.
- Test audio-only, visual-only and audio-visual classifiers for different levels of acoustic noise.
- Evaluate performance based on misclassification rate.

# Conclusion

- Visual speech in itself does not contain sufficient information for speech recognition...

- ...but by combining visual and audio speech features we are able to achieve better performance than what is possible with audio-only ASR.

# Future work

- Speech features are not i.i.d. (hidden Markov models) (sprint)
- Audio and visual speech is asynchronous (dynamic Bayesian networks) (GrMPy)
- Adaptive stream weighting
- ...

# Thank you!
# Any questions?