

Script generated by TTT

Title: groh: profile1 (25.06.2014)

Date: Wed Jun 25 08:16:07 CEST 2014

Duration: 92:37 min

Pages: 33



Person Detection from Audio

“Speaker Diarization / Segmentation“: given multi-party audio data (possibly with background noise):

→ who talks when?

- Typically 3 steps:
 - segmentation into speech / non-speech
 - detection of speaker transitions
 - clustering of speaker segments (+ classification of speaker)
- Segmentation into speech / non-speech:
 - Generate features:
 - ++ digital signal (pre-) processing (involving e.g. sub-division signal into overlapping samples of typically several ms, Fourier-transform etc.)
 - ++ MEL filters → MEL cepstrum coefficients
 - ++ Further Fourier- and other transformations
 - ++ additional features: zero-crossing rates, energy statistics etc.



Person Detection from Video

- First step: Face detection
 - Naive approach: simple pixel based binary classifier. Problem: too many possibilities for non-faces
 - Other approaches:
 - detect correct relatively positioned patches of skin, eyes or other face elements. Advantage; relatively robust against rotations
 - Approach [6]: Use special features instead of pixels (advantage: domain knowledge can be encoded into features), Intelligent feature selection / combination of simple binary classifiers that work on single features (AdaBoost)
- (optional second step: face recognition(e.g. via Eigenfaces (via PCA) [5])



Person Detection from Video

- First step: Face detection
 - Naive approach: simple pixel based binary classifier. Problem: too many possibilities for non-faces
 - Other approaches:
 - detect correct relatively positioned patches of skin, eyes or other face elements. Advantage; relatively robust against rotations
 - Approach [6]: Use special features instead of pixels (advantage: domain knowledge can be encoded into features), Intelligent feature selection / combination of simple binary classifiers that work on single features (AdaBoost)
- (optional second step: face recognition(e.g. via Eigenfaces (via PCA) [5])



Person Detection from Video

- **First step: Face detection**
 - Naive approach: simple pixel based binary classifier.
Problem: too many possibilities for non-faces
 - Other approaches:
 - detect correct relatively positioned patches of skin, eyes or other face elements. Advantage; relatively robust against rotations
 - Approach [6]: Use special features instead of pixels (advantage: domain knowledge can be encoded into features), Intelligent feature selection / combination of simple binary classifiers that work on single features (AdaBoost)
- (optional second step: face recognition(e.g. via Eigenfaces (via PCA) [5])



Person Detection from Video

- **Human figure detection:**
 - Main problem: too many options (clothes, accessoires)→ pixels as features won't work
 - **Approaches:**
 - features: histograms of directions of detected edges



Fig. 8. People detection. Examples of people detection in public spaces (pictures from [216]).

[1]



Person Detection from Video

- **First step: Face detection**
 - Naive approach: simple pixel based binary classifier.
Problem: too many possibilities for non-faces
 - Other approaches:
 - detect correct relatively positioned patches of skin, eyes or other face elements. Advantage; relatively robust against rotations
 - Approach [6]: Use special features instead of pixels (advantage: domain knowledge can be encoded into features), Intelligent feature selection / combination of simple binary classifiers that work on single features (AdaBoost)
- (optional second step: face recognition(e.g. via Eigenfaces (via PCA) [5])



Detecting Social Signals: Physical Appearance

- No specific studies in direct view of interpretation / effect as social signals ; however: general investigations:
 - Beauty assessment
 - via face symmetries as features + classifiers or
 - via distance to „ideal“ faces



morphed images taken from [7]





Detecting Social Signals: Physical Appearance

- No specific studies in direct view of interpretation / effect as social signals ; however: general investigations:
 - Beauty assessment
 - via face symmetries as features + classifiers or
 - via distance to „ideal“ faces



morphed images taken from [7]



Detecting Social Signals: Gaze and Face

- **AU**: smallest discernable temporal feature sequence: sequence of geometry or appearance features (modeled e.g. via Dynamic Bayesian Networks (DBN))
- Detection: example: basic integrative methods based on **optical flow** on detected faces:
 - optical flow: motion pattern of picture elements (e.g. pixels): represented by vector field of velocity $V(x,y,t)$ of intensity:

$$I(x + dx, y + dy, t + dt) = I(x, y, t) + \frac{\partial I}{\partial x} dx + \frac{\partial I}{\partial y} dy + \frac{\partial I}{\partial t} dt + O(d^2)$$

$$\rightarrow \frac{\partial I}{\partial x} V_x + \frac{\partial I}{\partial y} V_y + \frac{\partial I}{\partial t} = 0 \quad (\text{optical flow equation})$$

use numeric methods to compute solutions



Detecting Social Signals: Gaze and Face

- **AU**: smallest discernable temporal feature sequence: sequence of geometry or appearance features (modeled e.g. via Dynamic Bayesian Networks (DBN))
- Detection: example: basic integrative methods based on **optical flow** on detected faces:
 - optical flow: motion pattern of picture elements (e.g. pixels): represented by vector field of velocity $V(x,y,t)$ of intensity:

$$I(x + dx, y + dy, t + dt) = I(x, y, t) + \frac{\partial I}{\partial x} dx + \frac{\partial I}{\partial y} dy + \frac{\partial I}{\partial t} dt + O(d^2)$$

$$\rightarrow \frac{\partial I}{\partial x} V_x + \frac{\partial I}{\partial y} V_y + \frac{\partial I}{\partial t} = 0 \quad (\text{optical flow equation})$$

use numeric methods to compute solutions



Detecting Social Signals: Gaze and Face

- **AU**: smallest discernable temporal feature sequence: sequence of geometry or appearance features (modeled e.g. via Dynamic Bayesian Networks (DBN))
- Detection: example: basic integrative methods based on **optical flow** on detected faces:
 - optical flow: motion pattern of picture elements (e.g. pixels): represented by vector field of velocity $V(x,y,t)$ of intensity:

$$I(x + dx, y + dy, t + dt) = I(x, y, t) + \frac{\partial I}{\partial x} dx + \frac{\partial I}{\partial y} dy + \frac{\partial I}{\partial t} dt + O(d^2)$$

$$\rightarrow \frac{\partial I}{\partial x} V_x + \frac{\partial I}{\partial y} V_y + \frac{\partial I}{\partial t} = 0 \quad (\text{optical flow equation})$$

use numeric methods to compute solutions





Detecting Social Signals: Gaze and Face

- **AU**: smallest discernable temporal feature sequence: sequence of geometry or appearance features (modeled e.g. via Dynamic Bayesian Networks (DBN))
- Detection: example: basic integrative methods based on **optical flow** on detected faces:
 - optical flow: motion pattern of picture elements (e.g. pixels): represented by vector field of velocity $V(x,y,t)$ of intensity:

$$I(x + dx, y + dy, t + dt) = I(x, y, t) + \frac{\partial I}{\partial x} dx + \frac{\partial I}{\partial y} dy + \frac{\partial I}{\partial t} dt + O(d^2)$$

$$\rightarrow \frac{\partial I}{\partial x} V_x + \frac{\partial I}{\partial y} V_y + \frac{\partial I}{\partial t} = 0 \quad (\text{optical flow equation})$$

use numeric methods to compute solutions



Detecting Social Signals: From Audio

- Vocal features: up to now: mostly investigated for speech detection
- **Prosody**: **pitch, tempo, energy**
 - pitch: first fundamental frequency (1st maximum in Fourier transform (e.g. 30ms frames))
 - tempo: vowels / sec. ; vowel: phonetically relevant unit
 - energy E of signal $s(t)$: $E = \sum_i s(t_i)^2$
- Few efforts so far in analysis of **non-linguistic vocalizations**
 - example: laughter detection (e.g. via SVMs) and **linguistic vocalizations**
- **silence detection**: e.g. via energy as feature (often as by-product of speaker diarization)



Detecting Social Signals: Gaze and Face

- **AU**: smallest discernable temporal feature sequence: sequence of geometry or appearance features (modeled e.g. via Dynamic Bayesian Networks (DBN))
- Detection: example: basic integrative methods based on **optical flow** on detected faces:
 - optical flow: motion pattern of picture elements (e.g. pixels): represented by vector field of velocity $V(x,y,t)$ of intensity:

$$I(x + dx, y + dy, t + dt) = I(x, y, t) + \frac{\partial I}{\partial x} dx + \frac{\partial I}{\partial y} dy + \frac{\partial I}{\partial t} dt + O(d^2)$$

$$\rightarrow \frac{\partial I}{\partial x} V_x + \frac{\partial I}{\partial y} V_y + \frac{\partial I}{\partial t} = 0 \quad (\text{optical flow equation})$$

use numeric methods to compute solutions



Detecting Social Signals: From Audio

- Vocal features: up to now: mostly investigated for speech detection
- **Prosody**: **pitch, tempo, energy**
 - pitch: first fundamental frequency (1st maximum in Fourier transform (e.g. 30ms frames))
 - tempo: vowels / sec. ; vowel: phonetically relevant unit
 - energy E of signal $s(t)$: $E = \sum_i s(t_i)^2$
- Few efforts so far in analysis of **non-linguistic vocalizations**
 - example: laughter detection (e.g. via SVMs) and **linguistic vocalizations**
- **silence detection**: e.g. via energy as feature (often as by-product of speaker diarization)



- Important issue: behavioral cues can have different meaning if happening in different **outer contexts**
- Example: **temporal dynamics** of behavioral cues / social signals (e.g. relative person-person timing, person-environment timing etc.)
- Other important issue: **multi-modal combination / fusion** of social signals (e.g. audio and interaction geometry)

Social Situation Models as Models of Social Context

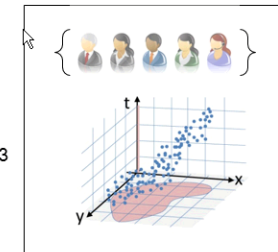
Social Situation:

Co-located social interaction with full mutual awareness



Simplified Social Situation Model:

- Participating **persons**: P: set of IDs
- **Spatio-temporal** reference: X: sub-set of $\mathbb{R} \times \mathbb{R}^3$
- $\rightarrow S = (P, X)$



Social Situation Models as Models of Social Context

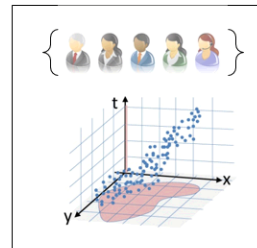
Social Situation:

Co-located social interaction with full mutual awareness



Simplified Social Situation Model:

- Participating **persons**: P: set of IDs
- **Spatio-temporal** reference: X: sub-set of $\mathbb{R} \times \mathbb{R}^3$
- $\rightarrow S = (P, X)$



Social Situation Models as Models of Social Context

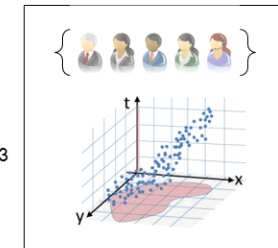
Social Situation:

Co-located social interaction with full mutual awareness



Simplified Social Situation Model:

- Participating **persons**: P: set of IDs
- **Spatio-temporal** reference: X: sub-set of $\mathbb{R} \times \mathbb{R}^3$
- $\rightarrow S = (P, X)$



Detecting Social Situations: Mobile Social Signal Processing

- **Social Situations: detection and understanding:** Social Signal Processing (see work by A. Vinciarelli et al.)
- **Social signal processing:** deriving higher level social models from low level social signal-sources (audio, video, etc.)
- Use **sensors in mobile devices** (see work by A. Pentland et al.)

Mobile Social Signal Processing

Detecting Social Situations: Mobile Social Signal Processing

Social Situation **detection**

- **Example:** microphone → audio-signals → speaker diarization → set of interacting persons
- **Example:** gyroscope, accelerometer, ultrasound-s. → relative body distance & orientation → set of interacting persons

Social Situation **understanding**

- **Example:** microphone → audio-signals → analysis of prosody → emotion detection → model of state of mind of person(s)

Detecting Social Situations: Mobile Social Signal Processing

Social Situation **detection**

- **Example:** microphone → audio-signals → speaker diarization → set of interacting persons
- **Example:** gyroscope, accelerometer, ultrasound-s. → relative body distance & orientation → set of interacting persons

Social Situation **understanding**

- **Example:** microphone → audio-signals → analysis of prosody → emotion detection → model of state of mind of person(s)

Geometry of Social Interaction

Interpersonal **distances**

- Hall: „general **quality**“ of social relation → 4 personal **zones**
- Other **influences** (?):
social context:
 architectural environment (socio-petal, socio-fugal forces (Watson)), density, gender, etc.
individual context:
 culture, age, self-esteem, disabilities,



Body **angles**

- Kendon: F-Formations [8]

Geometry of Social Interaction

Interpersonal distances

- Hall: „general quality“ of social relation
→ 4 personal zones
- Other influences (?):
social context:
architectural environment (socio-petal,
socio-fugal forces (Watson)), density,
gender, etc.
individual context:
culture, age, self-esteem, disabilities,



Body angles

- Kendon: F-Formations [8]

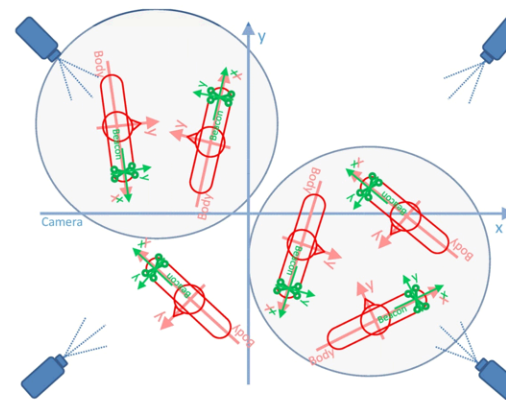
6 / 18

Model for Human Social Int. Geometry: Function $p(\delta\theta, \delta d)$

- Idea: Reduce n-ary social interaction to binary;
infer n-ary by graph clustering
- Binary: $p(\delta\theta, \delta d)$
- $\delta d(t) = \pm |P_{x,y} s_1(t) - P_{x,y} s_2(t)|$
- $\delta\theta(t) = \theta_z(R_{12}(t)) = \theta_z((R_1(t) (R_2(t))^T)$
- Optional: $p(\delta\theta, \delta d, \overline{\delta d})$

9 / 18

Experiment



8 / 18

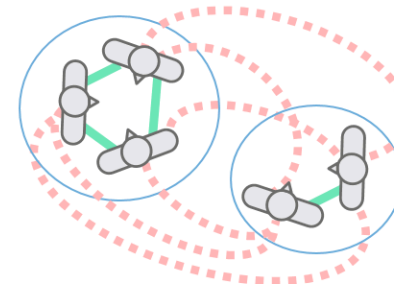
Results

Experiment data: Manual annotation

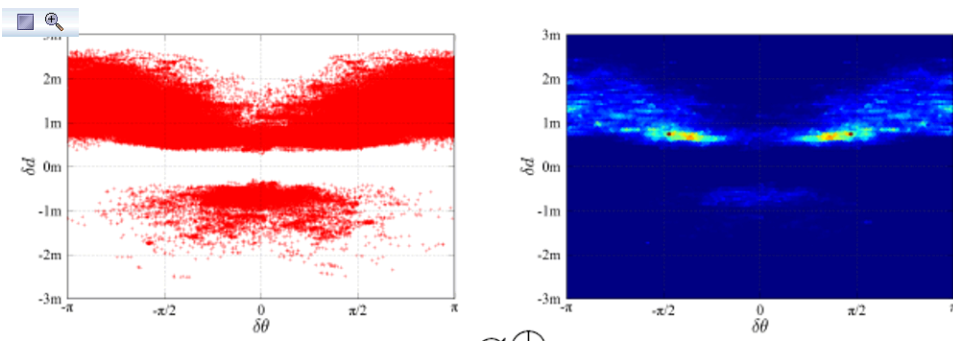
- $|S^{\oplus}| = 321307 (\delta\theta, \delta d)$ pairs corresponding to „in a social situation“
- $|S^{\ominus}| = 398335 (\delta\theta, \delta d)$ pairs corresponding to „not in a social situation“

Example:

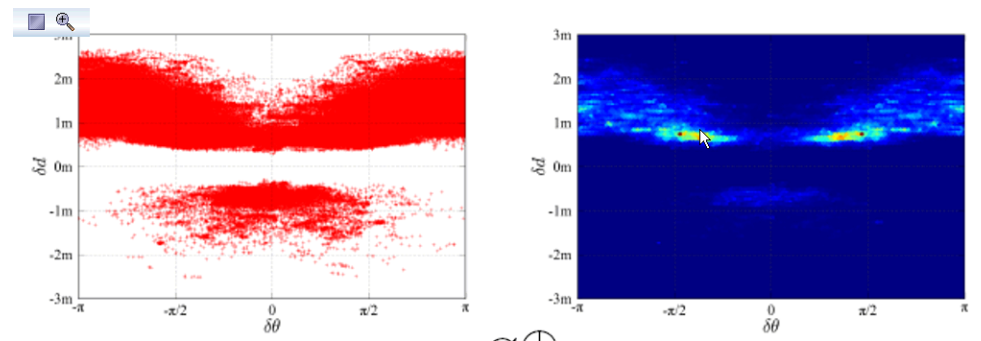
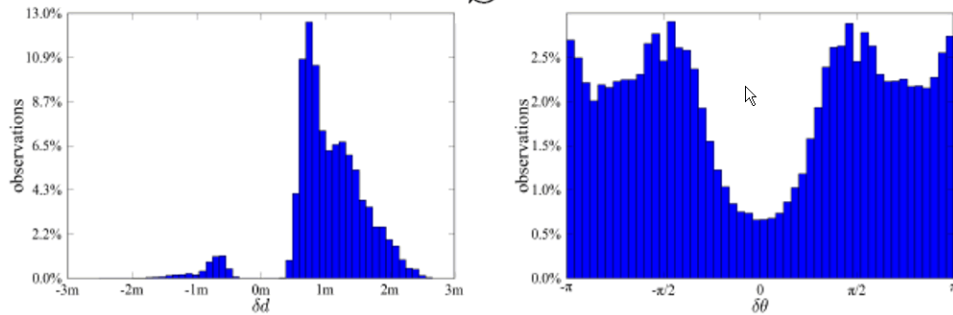
- $|S^{\oplus}| = 4$
- $|S^{\ominus}| = 6$



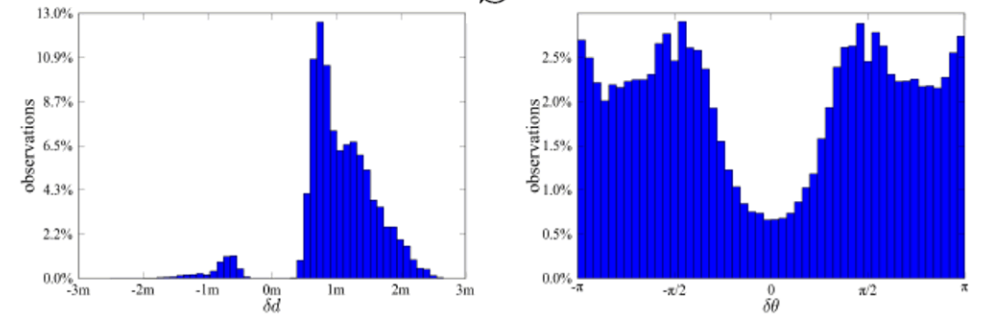
10 / 18



$S \oplus$



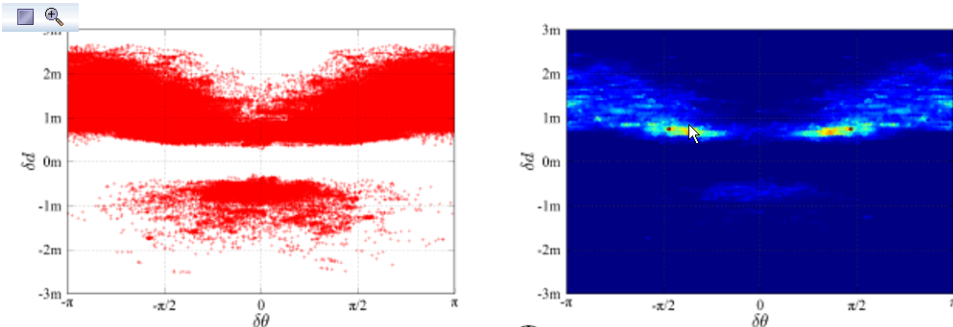
$S \oplus$



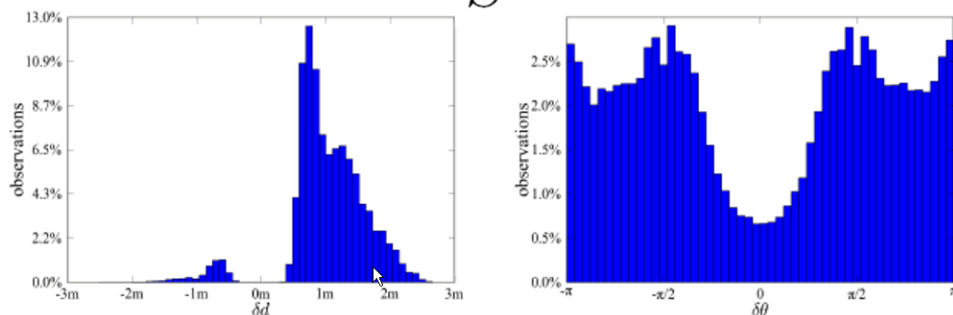
Classification Results

Classifier	Accuracy*
Gaussian Mixture Model (3 Gaussians)	74,34 %
Gaussian Mixture Model (5 Gaussians)	74,67 %
Gaussian Mixture Model (7 Gaussians)	74,59 %
Naive Bayes	65,45 %
Support Vector Machine (Polyn. Kernel)	77,81 %

(*) w. 10-fold cross validation



$S \oplus$



Reconstructing Social Situations

- For each t : complete weighted Graph $G(V,E,w,t)$ with V =set of persons,

$$w((s_1,s_2)) = \frac{p^{\oplus}(\delta\theta_{s_1s_2}, \delta d_{s_1s_2})}{p^{\oplus}(\delta\theta_{s_1s_2}, \delta d_{s_1s_2}) + p^{\ominus}(\delta\theta_{s_1s_2}, \delta d_{s_1s_2})}$$

- Average Link Clustering of $G(V,E,w,t)$ + Maximum Modularity Dendrogram Cut \rightarrow Partition X of V
- Compare X with annotation X' via $RAND(X,X')$ \rightarrow Accuracy of Social Situation Detection for each t
- Average over all t : RAND \sim 0.76 Adj.Rand \sim 0.529

Reconstructing Social Situations

- For each t : complete weighted Graph $G(V,E,w,t)$ with V =set of persons,

$$w((s_1,s_2)) = \frac{p^{\oplus}(\delta\theta_{s_1s_2}, \delta d_{s_1s_2})}{p^{\oplus}(\delta\theta_{s_1s_2}, \delta d_{s_1s_2}) + p^{\ominus}(\delta\theta_{s_1s_2}, \delta d_{s_1s_2})}$$

- Average Link Clustering of $G(V,E,w,t)$ + Maximum Modularity Dendrogram Cut \rightarrow Partition X of V
- Compare X with annotation X' via $RAND(X,X')$ \rightarrow Accuracy of Social Situation Detection for each t
- Average over all t : RAND \sim 0.76 Adj.Rand \sim 0.529