

Statistics and Machine Learning using belief functions

Lecture 4 – Estimation from Uncertain Data

Thierry Denœux

Université de Technologie de Compiègne
HEUDIASYC (UMR CNRS 6599)

<https://www.hds.utc.fr/~tdenoeux>

Beijing University of Technology
May 2017

Outline

- 1 Introduction
 - Motivations
 - Examples
- 2 Evidential EM algorithm
 - Evidential Likelihood
 - E²M algorithm
- 3 Partially supervised classification
 - Linear discriminant analysis
 - Logistic regression
 - Results

Outline

- 1 Introduction
 - Motivations
 - Examples
- 2 Evidential EM algorithm
 - Evidential Likelihood
 - E²M algorithm
- 3 Partially supervised classification
 - Linear discriminant analysis
 - Logistic regression
 - Results

Uncertain data

- Uncertain data arise in many applications (but it is usually neglected).
- Uncertainty may be due to:
 - **Limitations of the underlying measuring equipment** (unreliable sensors, indirect measurements), e.g.: biological sensor for toxicity measurement in water.
 - Use of **imputation, interpolation or extrapolation techniques**, e.g.: clustering of moving objects whose position is measured asynchronously by a sensor network,
 - **Partial or uncertain responses in surveys or subjective data annotation**, e.g.: sensory analysis experiments, data labeling by experts, etc.
- How to carry out statistical analysis of uncertain data?

Introductory example

- Let us consider a population in which some disease is present in proportion θ .
- n patients have been selected **at random** from that population. Let $x_i = 1$ if patient i has the disease, $x_i = 0$ otherwise. Each x_i is a realization of $X_i \sim \mathcal{B}(\theta)$.
- We assume that the x_i 's are **not observed directly**. For each patient i , a physician gives a **degree of plausibility** $pl_i(1)$ that patient i has the disease and a **degree of plausibility** $pl_i(0)$ that patient i does not have the disease.
- The observations are **uncertain data** of the form pl_1, \dots, pl_n .
- How to estimate θ ?

Aleatory vs. epistemic uncertainty

- In the previous example, uncertainty has **two distinct origins**:
 - 1 **Before** a patient has been drawn at random from the population, uncertainty is due to the **variability** of the variable of interest in the population. This is **aleatory uncertainty**.
 - 2 **After** the random experiment has been performed, uncertainty is due to **lack of knowledge** of the state of each particular patient. This is **epistemic uncertainty**.
- Epistemic uncertainty can be reduced by carrying out further investigations. Aleatory uncertainty cannot.

Approach

- In this lecture, we will consider statistical estimation problems in which **both kinds of uncertainty are present**: it will be assumed that each data item x
 - has been generated at random from a population (aleatory uncertainty), but
 - it is ill-known because of imperfect measurement or perception (epistemic uncertainty).
- The proposed model treats these two kinds of uncertainty separately:
 - **Aleatory uncertainty** will be represented by a **parametric statistical model**;
 - **Epistemic uncertainty** will be represented using **belief functions**.
- Application: partially supervised learning.

Outline

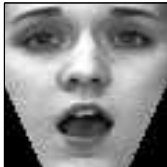
- 1 Introduction
 - Motivations
 - **Examples**
- 2 Evidential EM algorithm
 - Evidential Likelihood
 - E²M algorithm
- 3 Partially supervised classification
 - Linear discriminant analysis
 - Logistic regression
 - Results

Facial expressions

joy



surprise



sadness



disgust



anger



fear

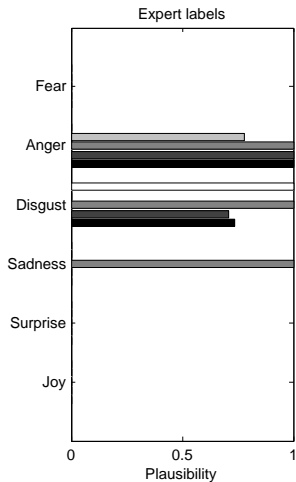


Recognition of facial expressions

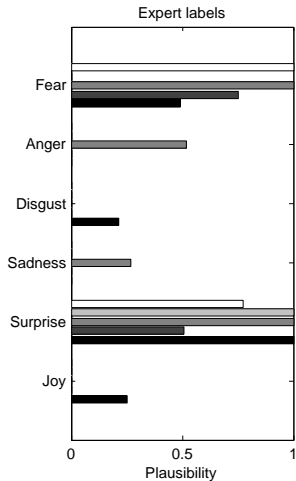
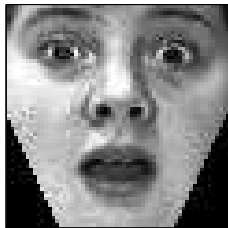
Experiment

- To achieve good performances in such tasks (object classification in images or videos), we need a large number of labeled images.
- However, **ground truth is usually not available** or difficult to determine with high precision and reliability: it is necessary to have the images subjectively annotated (labeled) by humans.
- How to **account for uncertainty** in such subjective annotations?
- Experiment:
 - Images were labeled by 5 subjects;
 - For each image, subjects were asked to give a **degree of plausibility** for each of the 6 basic expressions.

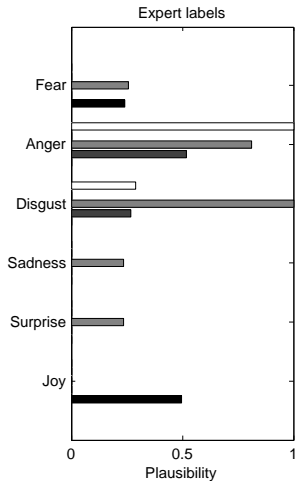
Example 1



Example 2



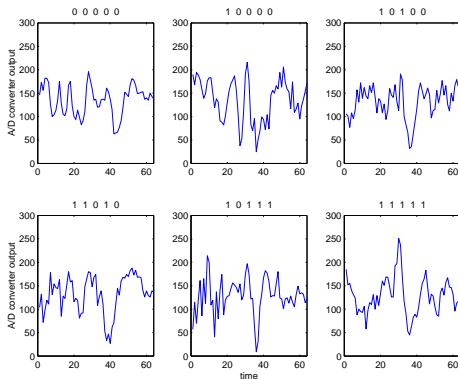
Example 3



Detection of K-complexes in EEG signals

- K-complexes: EEG waveforms that occur during stage 2 of non-rapid eye movement sleep . May aid sleep-based memory consolidation.
- Goal: build a system that automatically detects K-complexes in EEG data.
- We need a learning set of EEG signals labeled as positive and negative instances.
- Problem: no ground truth!
- Solution: label data by experts.

K-complex dataset



Each learning instance is composed of a feature vector and an uncertain class label.

Partially supervised learning

- **Complete data:** $\mathbf{x} = \{(\mathbf{w}_i, z_i)\}_{i=1}^n$ with
 - \mathbf{w}_i : feature vector for image i (pixel gray levels)
 - z_i : class of image i (one the six expressions).
- The feature vectors \mathbf{w}_i are perfectly observed but class labels are only **partially known** through subjective evaluations.
- **Observed data:**

$$\mathcal{L}_{ps} = \{(\mathbf{w}_i, m_i)\}_{i=1}^n,$$

where m_i is a mass function representing **partial information** about the class of object i .

- How to **learn a decision rule** from such data?

General approach

- 1 Postulate a parametric statistical model $p_{\mathbf{x}}(\mathbf{x}; \theta)$ for the complete data;
- 2 Represent epistemic data uncertainty using **belief functions** (observed data);
- 3 Estimate θ by **minimizing the conflict** between the model and the observed data using an extension of the **EM algorithm**: the evidential EM (E^2M) algorithm.

Outline

- 1 Introduction
 - Motivations
 - Examples
- 2 **Evidential EM algorithm**
 - Evidential Likelihood
 - E²M algorithm
- 3 Partially supervised classification
 - Linear discriminant analysis
 - Logistic regression
 - Results

Model

- Let \mathbf{X} be a (discrete) random vector taking values in $\Omega_{\mathbf{X}}$, with probability mass function $p_{\mathbf{X}}(\cdot; \theta)$ depending on an **unknown parameter** $\theta \in \Theta$.
- Let \mathbf{x} be a realization of \mathbf{X} (**complete data**).
- We assume that \mathbf{x} is only **partially observed**, and partial knowledge of \mathbf{x} is described by a **mass function** m on $\Omega_{\mathbf{X}}$ (“observed” data).
- Problem: estimate θ .

Outline

- 1 Introduction
 - Motivations
 - Examples
- 2 Evidential EM algorithm
 - Evidential Likelihood
 - E²M algorithm
- 3 Partially supervised classification
 - Linear discriminant analysis
 - Logistic regression
 - Results

Likelihood function (reminder)

- Given a parametric model $p_{\mathbf{X}}(\cdot; \theta)$ and an observation \mathbf{x} , the **likelihood function** is the mapping from Θ to $[0, 1]$ defined as

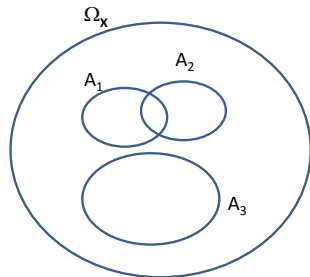
$$\theta \rightarrow L(\theta; \mathbf{x}) = p_{\mathbf{X}}(\mathbf{x}; \theta).$$

- It measures the “likelihood” or plausibility of each possible value of the parameter, after the data has been observed.
- If we observe that $\mathbf{x} \in A$, then the likelihood function is:

$$L(\theta; A) = \mathbb{P}_{\mathbf{X}}(A; \theta) = \sum_{\mathbf{x} \in A} p_{\mathbf{X}}(\mathbf{x}; \theta).$$

Evidential Likelihood function

Definition



- Assume that m has focal sets A_1, \dots, A_r .
- If we knew that $\mathbf{x} \in A_i$, the likelihood would be

$$L(\theta; A_i) = \mathbb{P}_{\mathbf{x}}(A_i; \theta) = \sum_{\mathbf{x} \in A_i} p_{\mathbf{x}}(\mathbf{x}; \theta).$$

- Taking the expectation with respect to m :

$$L(\theta; m) = \sum_{i=1}^r m(A_i) L(\theta; A_i)$$

Interpretation

- We have

$$\begin{aligned}
 L(\theta; m) &= \sum_{i=1}^r m(A_i) \sum_{\mathbf{x} \in A_i} p_{\mathbf{X}}(\mathbf{x}; \theta) \\
 &= \sum_{\mathbf{x} \in \Omega_{\mathbf{X}}} p_{\mathbf{X}}(\mathbf{x}; \theta) \sum_{A_i \ni \mathbf{x}} m(A_i) \\
 &= \sum_{\mathbf{x} \in \Omega_{\mathbf{X}}} p_{\mathbf{X}}(\mathbf{x}; \theta) pl(\mathbf{x}) = 1 - \kappa,
 \end{aligned}$$

where κ is the **degree of conflict** between $p_{\mathbf{X}}(\cdot; \theta)$ and m .

- Consequently, maximizing $L(\theta; m)$ with respect to θ amounts to **minimizing the conflict** between the parametric model and the uncertain observations

Case of fuzzy data

- We can also write $L(\theta; m)$ as:

$$L(\theta; m) = \sum_{\mathbf{x} \in \Omega_X} p_X(\mathbf{x}; \theta) p_I(\mathbf{x}) = \mathbb{E}_\theta [p_I(\mathbf{X})]$$

- If m is **consonant**, p_I may be interpreted as the membership function of a fuzzy subset of Ω_X : it can be seen as **fuzzy data**.
- $L(\theta; m)$ is then the **probability of the fuzzy data**, according to the definition given by Zadeh (1968).

Independence assumptions

- Let us assume that $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{np}$, where each \mathbf{x}_i is a realization from a p -dimensional random vector \mathbf{X}_i .
- Independence assumptions:
 - Stochastic independence** of $\mathbf{X}_1, \dots, \mathbf{X}_n$:

$$p_{\mathbf{X}}(\mathbf{x}; \theta) = \prod_{i=1}^n p_{\mathbf{X}_i}(\mathbf{x}_i; \theta), \quad \forall \mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \Omega_{\mathbf{X}}$$

- Cognitive independence** of $\mathbf{x}_1, \dots, \mathbf{x}_n$ with respect to m :

$$pl(\mathbf{x}) = \prod_{i=1}^n pl_i(\mathbf{x}_i), \quad \forall \mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \Omega_{\mathbf{X}}.$$

- Under these assumptions:

$$\log L(\theta; m) = \sum_{i=1}^n \log \mathbb{E}_{\theta} [pl_i(\mathbf{X}_i)].$$

Outline

- 1 Introduction
 - Motivations
 - Examples
- 2 Evidential EM algorithm
 - Evidential Likelihood
 - E²M algorithm
- 3 Partially supervised classification
 - Linear discriminant analysis
 - Logistic regression
 - Results

Evidential EM algorithm

- The evidential log-likelihood function $\log L(\theta; m)$ can be maximized using an **iterative algorithm** composed of two steps:

E-step: Compute the expectation of $\log L(\theta; \mathbf{X})$ with respect to $m \oplus p_{\mathbf{X}}(\cdot; \theta^{(q)})$:

$$Q(\theta, \theta^{(q)}) = \frac{\sum_{\mathbf{x} \in \Omega_{\mathbf{X}}} \log(L(\theta; \mathbf{x})) p_{\mathbf{X}}(\mathbf{x}; \theta^{(q)}) p_l(\mathbf{x})}{\sum_{\mathbf{x} \in \Omega_{\mathbf{X}}} p_{\mathbf{X}}(\mathbf{x}; \theta^{(q)}) p_l(\mathbf{x})}.$$

M-step: Maximize $Q(\theta, \theta^{(q)})$ with respect to θ .

- E- and M-steps are iterated until the increase of $\log L(\theta; m)$ becomes smaller than some threshold.

Properties

- 1 When m is categorical: $m(A) = 1$ for some $A \subseteq \Omega$, then the previous algorithm reduces to the EM algorithm \rightarrow **evidential EM (E²M) algorithm**.
- 2 Monotonicity: any sequence $L(\theta^{(q)}; m)$ for $q = 0, 1, 2, \dots$ of evidential likelihood values obtained using the E²M algorithm is non decreasing, i.e., it verifies

$$L(\theta^{(q+1)}; m) \geq L(\theta^{(q)}; m), \quad \forall q.$$

- 3 The algorithm **only uses the contour function p_l** , which drastically reduces the complexity of calculations.

Example: uncertain Bernoulli sample

Model and data

- Let us assume that the complete data $\mathbf{x} = (x_1, \dots, x_n)$ is a realization from an i.i.d. sample X_1, \dots, X_n from $\mathcal{B}(\theta)$ with $\theta \in [0, 1]$.
- We only have **partial information** about the x_i 's in the form: pI_1, \dots, pI_n , where $pI_i(x)$ is the plausibility that $x_i = x$, $x \in \{0, 1\}$.
- Under the cognitive independence assumption:

$$\begin{aligned}\log L(\theta; pI_1, \dots, pI_n) &= \sum_{i=1}^n \log \mathbb{E}_{\theta} [pI_i(X_i)] \\ &= \sum_{i=1}^n \log [(1 - \theta)pI_i(0) + \theta pI_i(1)]\end{aligned}$$

E- and M-steps

Complete data log-likelihood:

$$\log L(\theta, \mathbf{x}) = n \log(1 - \theta) + \log \left(\frac{\theta}{1 - \theta} \right) \sum_{i=1}^n x_i.$$

E-step: compute

$$Q(\theta, \theta^{(q)}) = n \log(1 - \theta) + \log \left(\frac{\theta}{1 - \theta} \right) \sum_{i=1}^n \xi_i^{(q)}, \text{ with}$$

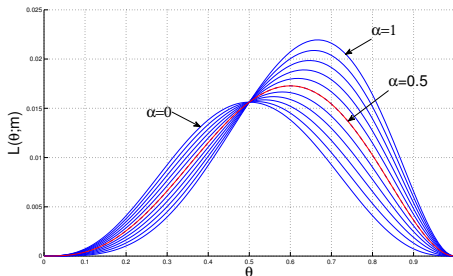
$$\xi_i^{(q)} = \mathbb{E}_{\theta^{(q)}} [X_i | p_i] = \frac{\theta^{(q)} p_i(1)}{(1 - \theta^{(q)}) p_i(0) + \theta^{(q)} p_i(1)}.$$

M-step:

$$\theta^{(q+1)} = \frac{1}{n} \sum_{i=1}^n \xi_i^{(q)}.$$

Numerical example

i	1	2	3	4	5	6
$pl_i(0)$	1	1	1	α	0	0
$pl_i(1)$	0	0	0	$1 - \alpha$	1	1



$$\alpha = 0.5$$

q	$\theta^{(q)}$	$L(\theta^{(q)}; pl)$
0	0.3000	6.6150
1	0.5500	16.8455
2	0.5917	17.2676
3	0.5986	17.2797
4	0.5998	17.2800
5	0.6000	17.2800

$$\hat{\theta} = 0.6$$

Outline

- 1 Introduction
 - Motivations
 - Examples
- 2 Evidential EM algorithm
 - Evidential Likelihood
 - E²M algorithm
- 3 Partially supervised classification**
 - Linear discriminant analysis
 - Logistic regression
 - Results

Classification

- We consider a population of **objects** partitioned in **g classes**.
- Each object is described by **d continuous features** $\mathbf{W} = (W^1, \dots, W^d)$ and a class variable Z .
- The goal of **classification** is to learn a **decision rule** that classifies any object from its feature vector, based on a learning set.

Partially supervised learning

- Classically, different learning tasks are considered:

Supervised learning: $\mathcal{L}_s = \{(\mathbf{w}_i, z_i)\}_{i=1}^n$;

Unsupervised learning: $\mathcal{L}_{ns} = \{\mathbf{w}_i\}_{i=1}^n$;

Semi-supervised learning: $\mathcal{L}_{ss} = \{(\mathbf{w}_i, z_i)\}_{i=1}^{n_s} \cup \{\mathbf{w}_i\}_{i=n_s+1}^n$

- Here, we consider **partially supervised learning**:

$$\mathcal{L}_{ps} = \{(\mathbf{w}_i, m_i)\}_{i=1}^n,$$

where m_i is a mass function representing **partial information** about the class of object i .

- This problem can be solved using the E²M algorithm using a suitable parametric model.
- In this lecture, I will present two models:
 - Linear discriminant analysis
 - Logistic regression

Outline

- 1 Introduction
 - Motivations
 - Examples
- 2 Evidential EM algorithm
 - Evidential Likelihood
 - E²M algorithm
- 3 **Partially supervised classification**
 - **Linear discriminant analysis**
 - Logistic regression
 - Results

Model

- Generative model:
 - Complete data: $\mathbf{x} = \{(\mathbf{w}_i, Z_i)\}_{i=1}^n$, assumed to be a realization of an **iid random sample** $\mathbf{X} = \{(\mathbf{W}_i, Z_i)\}_{i=1}^n$;
 - Given $Z_i = k$, \mathbf{W}_i is **multivariate normal** with mean $\boldsymbol{\mu}_k$ and **common variance matrix** Σ .
 - The proportion of class k in the population is π_k .
 - Parameter vector: $\boldsymbol{\theta} = (\{\pi_k\}_{k=1}^g, \{\boldsymbol{\mu}_k\}_{k=1}^g, \Sigma)$.
- The **Bayes rule** is approximated by assigning each object to the class k^* that maximizes the estimated posterior probability

$$p(Z = k | \mathbf{w}; \hat{\boldsymbol{\theta}}) = \frac{\phi(\mathbf{w}; \hat{\boldsymbol{\mu}}_k, \hat{\Sigma}) \hat{\pi}_k}{\sum_{\ell} \phi(\mathbf{w}; \hat{\boldsymbol{\mu}}_{\ell}, \hat{\Sigma}) \hat{\pi}_{\ell}},$$

where $\hat{\boldsymbol{\theta}}$ is the MLE of $\boldsymbol{\theta}$.

Complete-data likelihood

The complete-data likelihood is

$$L_c(\theta) = \prod_{i=1}^n p(\mathbf{w}_i | Z_i = z_i) p(z_i) \quad (1a)$$

$$= \prod_{i=1}^n \prod_{k=1}^g \phi(\mathbf{w}_i; \mu_k, \Sigma)^{z_{ik}} \pi_k^{z_{ik}}, \quad (1b)$$

where $\phi(\cdot; \mu_k, \Sigma)$ is the multivariate normal density,

$$\phi(\mathbf{w}; \mu_k, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{w} - \mu)^T \Sigma^{-1} (\mathbf{w} - \mu) \right\},$$

and z_{ik} is a binary class indicator variable, such that $z_{ik} = 1$ if $z_i = k$ and $z_{ik} = 0$ otherwise.

Observed-data likelihood

- Under the assumption of cognitive independence, the contour function on Ω_X is $pl(x) = \prod_{i=1}^n pl_i(x_i)$, with

$$pl_i(x_i) = \begin{cases} pl_{ik} & \text{if } x_i = (w_i, k) \text{ for some } k = 1, \dots, g \\ 0 & \text{otherwise.} \end{cases}$$

- The evidential likelihood is thus

$$L(\theta) = \prod_{i=1}^n \sum_{k=1}^g pl_{ik} \phi(w_i; \mu_k, \Sigma) \pi_k. \quad (2)$$

Special cases

- When there is no uncertainty, i.e., when $p_{iik} = z_{iik}$ for all (i, k) , we have

$$\sum_{k=1}^g p_{iik} \phi(\mathbf{w}_i; \mu_k, \Sigma) \pi_k = \prod_{k=1}^g \phi(\mathbf{w}_i; \mu_k, \Sigma)^{p_{iik}} \pi_k^{p_{iik}},$$

and the evidential likelihood (2) becomes identical to the complete-data likelihood (1b).

- When uncertainty is maximal, i.e., class labels are completely unknown, then $p_{iik} = 1$ for all (i, k) , and the evidential likelihood (2) becomes

$$L(\theta) = \prod_{i=1}^n \sum_{k=1}^g \phi(\mathbf{w}_i; \mu_k, \Sigma) \pi_k,$$

which is the likelihood function corresponding to the unsupervised case.

E²M algorithm: E-step

In the E-step of the E²M algorithm for this model, we compute the expectation of the complete-data log-likelihood

$$\ell_c(\theta) = \sum_{i=1}^n \sum_{k=1}^K z_{ik} [\log \phi(\mathbf{w}_i; \mu_k, \Sigma) + \log \pi_k]$$

with respect to the combined probability mass function

$$p_X(x|p_l; \theta^{(q)}) = \prod_{i=1}^n p(x_i|p_l; \theta^{(q)}),$$

with

$$p(x_i|p_l; \theta^{(q)}) = \begin{cases} \frac{p_l \pi_k^{(q)} \phi(\mathbf{w}_i; \mu_k^{(q)}, \Sigma^{(q)})}{\sum_{\ell} p_l \pi_{\ell}^{(q)} \phi(\mathbf{w}_i; \mu_{\ell}^{(q)}, \Sigma^{(q)})} & \text{if } x_i = (\mathbf{w}_i, k) \text{ for some } k \\ 0 & \text{otherwise.} \end{cases}$$

E²M algorithm: E-step (continued)

We get

$$Q(\theta, \theta^{(q)}) = \sum_{i=1}^n \sum_{k=1}^g t_{ik}^{(q)} [\log \phi(\mathbf{w}_i; \mu_k, \Sigma) \pi_k + \log \pi_k], \quad (3)$$

with

$$t_{ik}^{(q)} = \mathbb{E}(Z_{ik} | \mathbf{p}l; \theta^{(q)}) = \frac{\mathbf{p}l_{ik} \pi_k^{(q)} \phi(\mathbf{w}_i; \mu_k^{(q)}, \Sigma^{(q)})}{\sum_{\ell} \mathbf{p}l_{i\ell} \pi_{\ell}^{(q)} \phi(\mathbf{w}_i; \mu_{\ell}^{(q)}, \Sigma^{(q)})}. \quad (4)$$

E²M algorithm: M-step

- The parameter values maximizing $Q(\theta, \theta^{(q)})$ can be readily obtained as

$$\pi_k^{(q+1)} = \frac{1}{n} \sum_{i=1}^n t_{ik}^{(q)}, \quad \boldsymbol{\mu}_k^{(q+1)} = \frac{\sum_{i=1}^n t_{ik}^{(q)} \mathbf{w}_i}{\sum_{i=1}^n t_{ik}^{(q)}}.$$

$$\Sigma^{(q+1)} = \frac{1}{n} \sum_{i,k} t_{ik}^{(q)} (\mathbf{w}_i - \boldsymbol{\mu}_k^{(q+1)}) (\mathbf{w}_i - \boldsymbol{\mu}_k^{(q+1)})^T$$

- The complexity is the same as that of the EM algorithm with unsupervised data and precise attributes.

Outline

- 1 Introduction
 - Motivations
 - Examples
- 2 Evidential EM algorithm
 - Evidential Likelihood
 - E^2M algorithm
- 3 **Partially supervised classification**
 - Linear discriminant analysis
 - **Logistic regression**
 - Results

Model

- In contrast with LDA, LR starts with a model of the conditional distribution of Z given $W = w$. The conditional probabilities are

$$p_k(w; \theta) = \frac{\exp(\beta_k^T \tilde{w})}{1 + \sum_{\ell=1}^{g-1} \exp(\beta_\ell^T \tilde{w})}, \quad k = 1, \dots, g-1 \quad (5a)$$

$$p_g(w; \theta) = \frac{1}{1 + \sum_{\ell=1}^{g-1} \exp(\beta_\ell^T \tilde{w})}, \quad (5b)$$

where $p_k(w; \theta) = \mathbb{P}(Z = k | W = w; \theta)$, β_k is a $p + 1$ -dimensional vector of coefficients, $\theta = (\beta_1^T, \dots, \beta_{g-1}^T)^T$ is the vector of all parameters in the model, and $\tilde{w} = (1, w^T)^T$ is an extended input vector.

- Logistic regression maximizes the conditional likelihood

$$L_c(\theta) = \prod_{i=1}^n \mathbb{P}(Z_i = z_i | w_i; \theta) = \prod_{i=1}^n \prod_{k=1}^g p_k(w; \theta)^{z_{ik}}, \quad (6)$$

with $p_k(w; \theta)$ equal to (5).

Evidential likelihood

- Under the cognitive independence assumption, the evidential likelihood is

$$L(\theta) = \prod_{i=1}^n \sum_{k=1}^g p l_{ik} p_k(w_i; \theta). \quad (7)$$

- We can easily check that $L(\theta) = L_c(\theta)$ whenever $p l_{ik} = z_{ik}$ for all (i, k) , i.e., when there is no label uncertainty.
- On the other hand, in case of maximal uncertainty, i.e., when $p l_{ik} = 1$ for all (i, k) , we have $L(\theta) = 1$ for all θ , and the model parameters can no longer be estimated.

E²M algorithm: E-step

- In the E-step, we compute the expectation of the complete-data log-likelihood with respect to the combined probability mass function

$$p_Z(z|p_l; \theta^{(q)}) = \prod_{i=1}^n p_{Z_i}(z_i|p_l; \theta^{(q)}),$$

with

$$p_{Z_i}(k|p_l; \theta^{(q)}) = \frac{p_{l_{ik}} p_k(\mathbf{w}_i; \theta^{(q)})}{\sum_{\ell} p_{l_{i\ell}} p_{\ell}(\mathbf{w}_i; \theta^{(q)})}, \quad k = 1, \dots, g.$$

- We get

$$Q(\theta, \theta^{(q)}) = \sum_{i=1}^n \left\{ \sum_{k=1}^{g-1} t_{ik}^{(q)} \beta_k^T \tilde{\mathbf{w}}_i - \log \left(1 + \sum_{k=1}^{g-1} \beta_k^T \tilde{\mathbf{w}}_i \right) \right\}, \quad (8)$$

with

$$t_{ik}^{(q)} = \mathbb{E}(Z_{ik}|p_l; \theta^{(q)}) = \frac{p_{l_{ik}} p_k(\mathbf{w}_i; \theta^{(q)})}{\sum_{\ell} p_{l_{i\ell}} p_{\ell}(\mathbf{w}_i; \theta^{(q)})}. \quad (9)$$

E²M algorithm: M-step

- The maximization of (8) cannot be performed in one step and requires an iterative optimization procedure, such as the Newton-Raphson algorithm.
- It is actually not necessary to maximize function $Q(\theta, \theta^{(q)})$: we may simply make a step uphill, i.e., find some new estimate $\theta^{(q+1)}$ such that $Q(\theta^{(q+1)}, \theta^{(q)}) > Q(\theta^{(q)}, \theta^{(q)})$. Such a procedure is classically called a *Generalized EM algorithm*.
- An uphill step starting from the previous estimate $\theta^{(q)}$ can be made by carrying out one iteration of the Newton-Raphson algorithm with line search, i.e., by using the following update rule,

$$\theta^{(q+1)} = \theta^{(q)} - \eta \left[\frac{\partial^2 Q(\theta, \theta^{(q)})}{\partial \theta \partial \theta^T} \right]_{\theta=\theta^{(q)}}^{-1} \left. \frac{\partial Q(\theta, \theta^{(q)})}{\partial \theta} \right|_{\theta=\theta^{(q)}},$$

where η is the step size.

Outline

- 1 Introduction
 - Motivations
 - Examples
- 2 Evidential EM algorithm
 - Evidential Likelihood
 - E^2M algorithm
- 3 **Partially supervised classification**
 - Linear discriminant analysis
 - Logistic regression
 - **Results**

Sleep data

- 1178 EEG signals encoded as 64-dimensional patterns.
- Each example (positive or negative) was then assigned a soft label consisting of a Bayesian mass function m_i such that

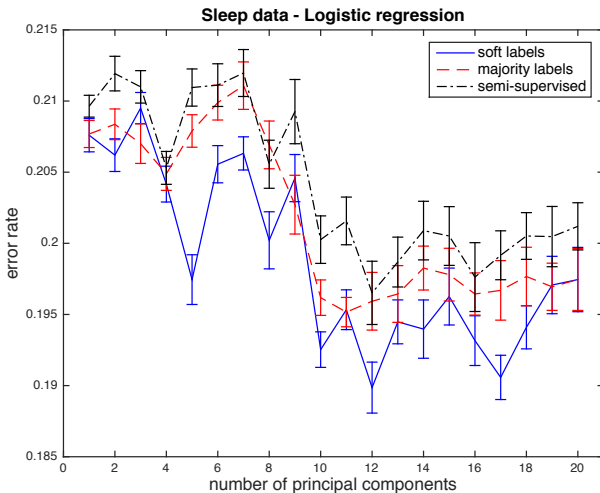
$$m_i(\{1\}) = k_i/5, \quad m_i(\{0\}) = 1 - k_i/5, \quad (10)$$

where 1 and 0 represent, respectively, the positive (*K*-complex) and negative (delta wave) class, and k_i denotes the number of experts who classified the pattern as positive.

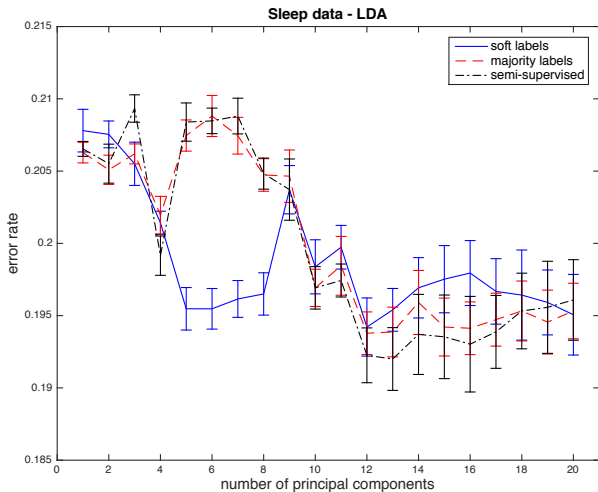
Methodology

- Both LR and LDA were applied to these data. To reduce the input dimension, Principal Component Analysis (PCA) was first used as a preprocessing step, and the number of components was varied between 1 and 20.
- The LR and LDA classifiers were trained using three different sets of labels:
 - 1 Soft labels (10), taking into account the proportion of experts in favor of each class;
 - 2 Crisp labels, corresponding to the majority decision;
 - 3 “Semi-supervised labels”: instances classified as positive by two or three experts were considered as ambiguous and were labeled by the vacuous mass function $m_?$; the other instances were labeled unambiguously according to the majority class.
- We used 10-fold cross-validation, repeated 10 times with different random partitions. The mean cross-validation error rates with corresponding 95% confidence intervals are represented as functions of the number of principal components in the following figures.

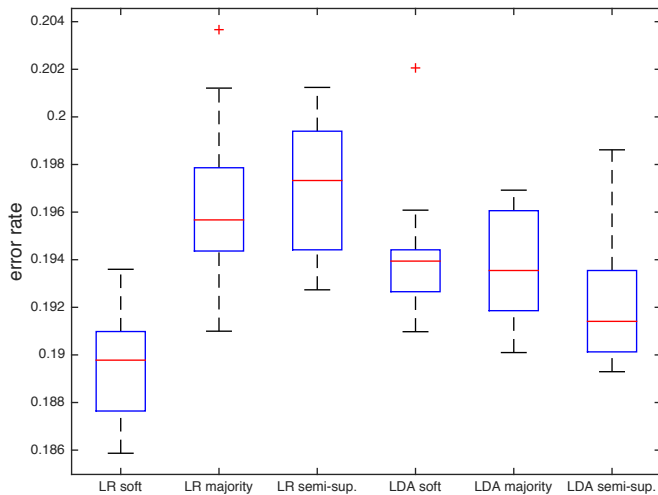
Results: logistic regression



Results: LDA



Results: comparison



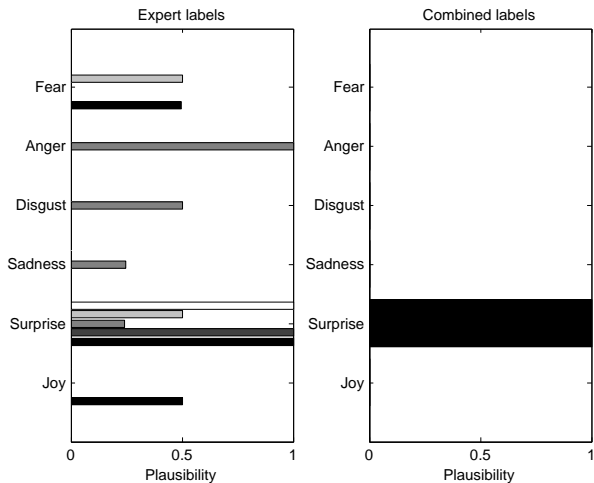
Expression recognition problem

Experimental settings

- 216 images of 60×70 pixels, 36 in each class.
- One half for training, the rest for testing.
- A reduced number of features was extracted using Principal component analysis (PCA).
- Each training image was labeled by 5 subjects who gave **degrees of plausibility** for each image and each class.
- The plausibilities were combined using **Dempster's rule** (after some discounting to avoid total conflict).

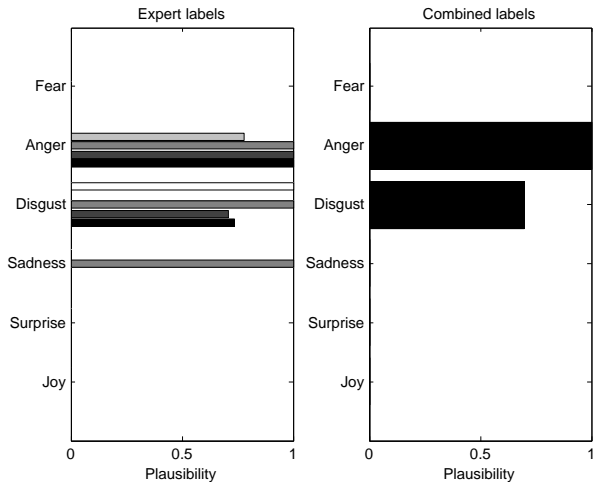
Combined labels

Example 1



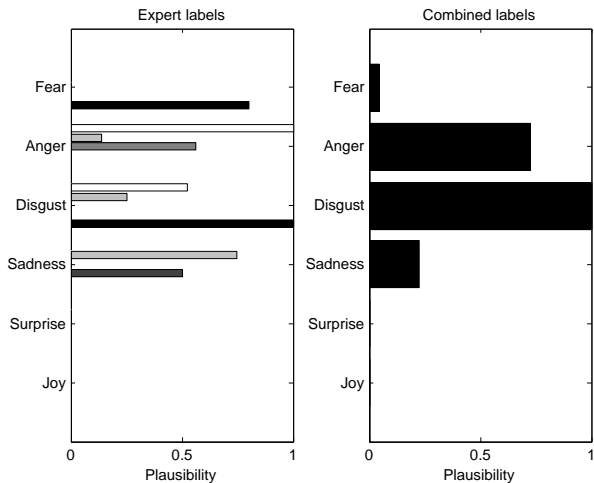
Combined labels

Example 2

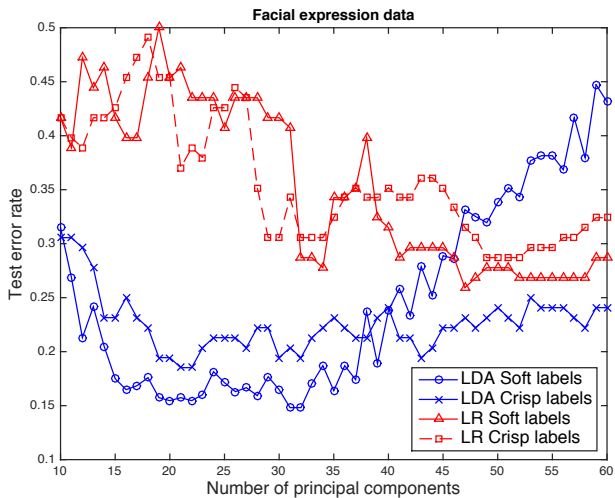


Combined labels

Example 3

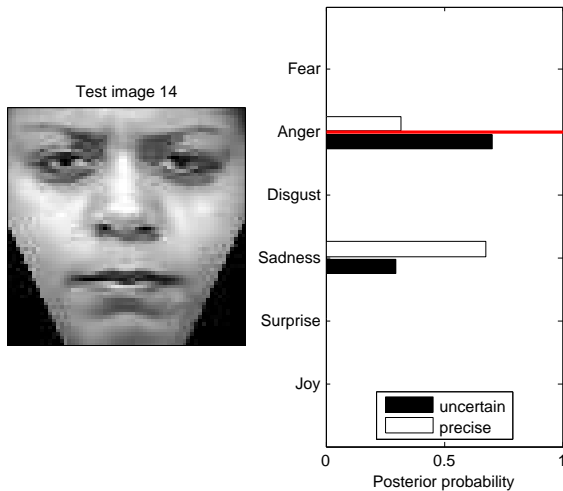


Results



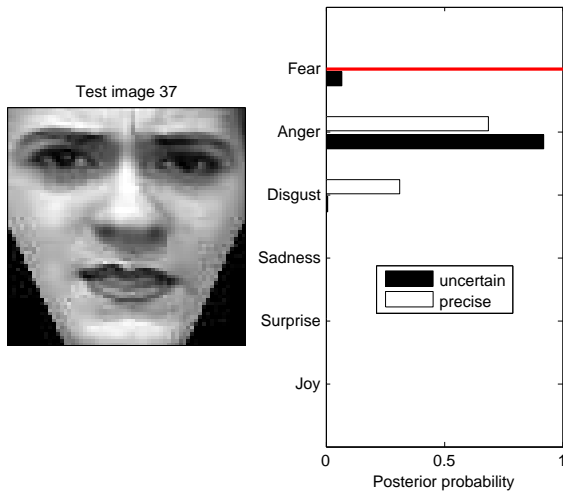
Results

Example 1



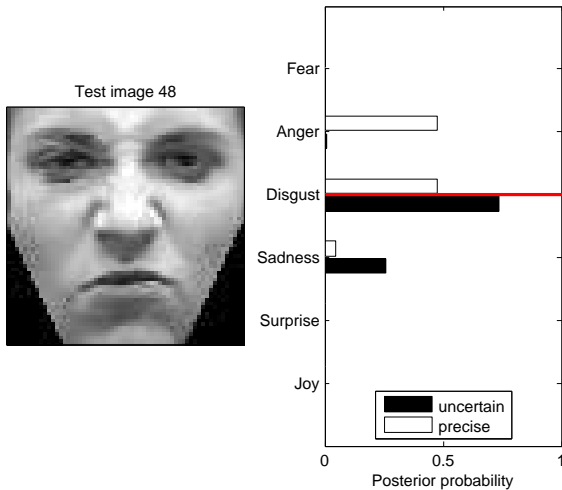
Results

Example 2



Results

Example 3



References I

cf. <https://www.hds.utc.fr/~tdenoeux>



T. Denœux.

Maximum likelihood estimation from fuzzy data using the EM algorithm.
Fuzzy Sets and Systems, 183:72-91, 2011.



T. Denœux.

Maximum likelihood estimation from Uncertain Data in the Belief Function Framework.
IEEE Trans. Knowledge and Data Engineering, Vol. 25, Issue 1, pages 119-130, 2013.



Z. L. Cherfi, L. Oukhellou, E. Côme, T. Denœux and P. Aknin.

Partially supervised Independent Factor Analysis using soft labels elicited from multiple experts: Application to railway track circuit diagnosis.
Soft Computing, Vol. 16, Number 5, pages 741-754, 2012.

References II

cf. <https://www.hds.utc.fr/~tdenoeux>



E. Ramasso and T. Denoeux.

Making use of partial knowledge about hidden states in HMMs: an approach based on belief functions.

IEEE Transactions on Fuzzy Systems, Vol. 22, Issue 2, pages 395-405, 2014.



B. Quost and T. Denoeux.

Clustering and classification of fuzzy data using the fuzzy EM algorithm.

Fuzzy Sets and Systems, Vol. 286, pages 134-156, 2016.



B. Quost, T. Denoeux and S. Li.

Parametric classification with soft labels using the Evidential EM algorithm.

Advances in Data Analysis and Classification, submitted, 2017.