

Some results on Imprecise discriminant analysis

11th Workshop on Principles and Methods of Statistical Inference with Interval Probability

CARRANZA-ALARCON Yonatan-Carlos
Ph.D. Candidate in Computer Science

DESTERCKE Sébastien
Ph.D Director



30 July 2018 to 01 August 2018

Overview

Imprecise Discriminant Analysis Classification

- Classification
 - Decision Making
 - Discriminant Analysis
- Imprecise Classification
 - Imprecise Decision
 - Imprecise Linear discriminant analysis
- Future work
- Conclusions

Overview

- Classification
 - Decision Making
 - Discriminant Analysis
- Imprecise Classification
 - Imprecise Decision
 - Imprecise Linear discriminant analysis
- Future work
- Conclusions

Classification - Setting

A classic classification problem is composed of :

- Data training $D = \{x_i, y_i\}_{i=0}^N$ such as :
 - (Input) $x_i \in \mathcal{X}$ are regressors or features (often $x_i \in \mathbb{R}^P$).
 - (Output) $y_i \in \mathcal{K}$ is a response category variable, with $\mathcal{K} = \{m_1, \dots, m_K\}$

Classification - Setting

A classic classification problem is composed of :

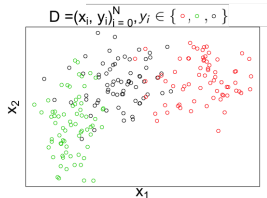
- Data training $D = \{x_i, y_i\}_{i=0}^N$ such as :
 - (Input) $x_i \in \mathcal{X}$ are regressors or features (often $x_i \in \mathbb{R}^p$).
 - (Output) $y_i \in \mathcal{K}$ is a response category variable, with $\mathcal{K} = \{m_1, \dots, m_K\}$

Objective

Given training data $D = \{x_i, y_i\}_{i=0}^N$, we need to learn a classification rule : $\phi : \mathcal{X} \rightarrow \mathcal{Y}$ in order to predict a new observation $\phi(\mathbf{x}^*)$

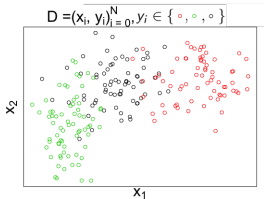
Classification - Outline (Example)

Getting Training Data



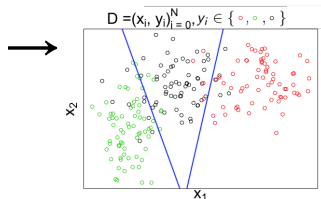
Classification - Outline (Example)

Getting Training
Data



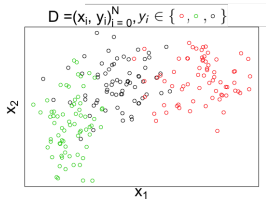
Learning a
classification rule :

$$\phi : \mathcal{X} \rightarrow \mathcal{Y}$$



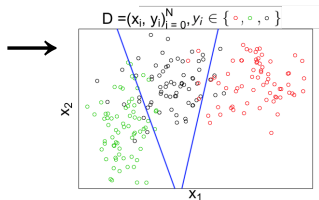
Classification - Outline (Example)

Getting Training Data



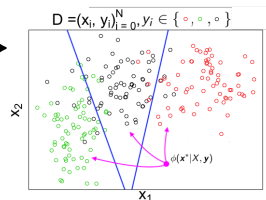
Learning a classification rule :

$$\phi : \mathcal{X} \rightarrow \mathcal{Y}$$



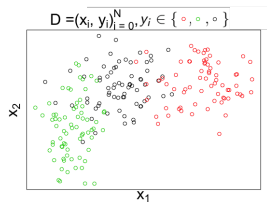
Predict class for new instances :

$$\hat{y}^* := \phi(\mathbf{x}^* | X, \mathbf{y})$$



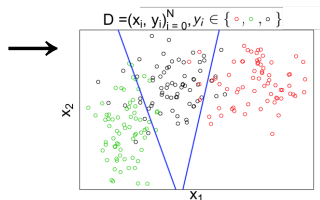
Classification - Outline (Example)

Getting Training Data



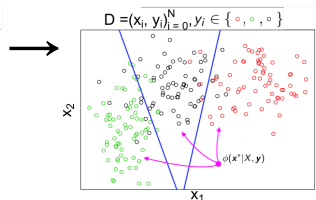
Learning a classification rule :

$$\phi : \mathcal{X} \rightarrow \mathcal{Y}$$



Predict class for new instances :

$$\hat{y}^* := \phi(\mathbf{x}^* | X, \mathbf{y})$$



But :

- How can we learn the “classification rule” (model) from training data ?

Decision Making in Statistic

- In statistic : classification rule often seen as a decision-making problem under risk of getting missclassification.

$$\mathcal{R}(y, \varphi(X)) = \arg \min_{\varphi(X) \in \mathcal{K}} \mathbb{E}_{\mathcal{X} \times \mathcal{Y}} [\mathcal{L}(y, \varphi(X))] \quad (1)$$

- Under 1/0 loss function \mathcal{L} , minimizing \mathcal{R} equivalent to :

$$\phi(\mathbf{x}^* | X, \mathbf{y}) := \arg \max_{m_k \in \mathcal{K}} P(y = m_k | X = \mathbf{x}^*) \quad (2)$$

- Where :

- The predicted class $\hat{y}^* = \phi(\mathbf{x}^* | X, \mathbf{y})$ is the most probable (equation (2)).
- This last equation (2) is also known as Bayes classifier [1, pp. 21].

Decision Making in Statistic

- In statistic : classification rule often seen as a decision-making problem under risk of getting missclassification.

$$\mathcal{R}(y, \varphi(X)) = \arg \min_{\varphi(X) \in \mathcal{K}} \mathbb{E}_{\mathcal{X} \times \mathcal{Y}} [\mathcal{L}(y, \varphi(X))] \quad (1)$$

- Under 1/0 loss function \mathcal{L} , minimizing \mathcal{R} equivalent to :

$$\phi(\mathbf{x}^* | X, \mathbf{y}) := \arg \max_{m_k \in \mathcal{K}} P(y = m_k | X = \mathbf{x}^*) \quad (2)$$

- Where :
 - The predicted class $\hat{y}^* = \phi(\mathbf{x}^* | X, \mathbf{y})$ is the most probable (equation (2)).
 - This last equation (2) is also known as Bayes classifier [1, pp. 21].

Decision Making in Statistic

Definition (Preference ordering [5, pp. 47])

With general loss $\mathcal{L}(\cdot, \cdot)$, m_a is preferred to m_b , denoted by $m_a > m_b$, if and only if :

$$\mathbb{E}_P[\mathcal{L}(\cdot, m_a) | \mathbf{x}^*] < \mathbb{E}_P[\mathcal{L}(\cdot, m_b) | \mathbf{x}^*]$$

In the particular case where $\mathcal{L}(\cdot, \cdot)$ is the 0/1 loss function we get :

$$m_a \succeq m_b \iff \frac{P(y = m_a | X = \mathbf{x}^*)}{P(y = m_b | X = \mathbf{x}^*)} > 1$$

where $P(Y = m_a | X = \mathbf{x}^*)$ is the class probability. We then take the **maximal element** of the complete order \succeq , i.e.

$$m_{i_K} \succeq m_{i_{K-1}} \succeq \dots \succeq m_{i_1} \iff P(y = m_{i_K} | \mathbf{x}^*) \geq \dots \geq P(y = m_{i_1} | \mathbf{x}^*)$$

(Precise) Discriminant Analysis

Applying Baye's rules to $P(Y = m_a | X = \mathbf{x}^*)$:

$$P(y = m_k | X = \mathbf{x}^*) = \frac{P(X = \mathbf{x}^* | y = m_k)P(y = m_k)}{\sum_{m_l \in \mathcal{K}} P(X = \mathbf{x}^* | y = m_l)P(y = m_l)}$$

where $\pi_k := \mathbb{P}_{Y=y_k}$ such as $\sum_j \pi_j = 1$ and $\mathcal{G}_k := \mathbb{P}_{X|Y=m_k} \sim \mathcal{N}(\mu_k, \Sigma_k)$

A frequentist point estimation :

$$\hat{\pi}_k = \frac{n_k}{N}$$

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} x_{i,k}$$

$$\hat{\Sigma}_k = \frac{1}{N - n_k} \sum_{i=1}^{n_k} (x_{i,k} - \bar{\mathbf{x}}_k)(x_{i,k} - \bar{\mathbf{x}}_k)^t$$

Overview

- Classification
 - Decision Making
 - Discriminant Analysis
- Imprecise Classification
 - Imprecise Decision
 - Imprecise Linear discriminant analysis
- Future work
- Conclusions

Decision Making in Imprecise Probabilities

Definition (Partial Ordering by Maximality Criterion)

Let \mathcal{P} a set of probabilities, then m_a is preferred to m_b if the cost of exchanging m_a with m_b have a positive lower expectation :

$$m_a \succ_M m_b \iff \inf_{P \in \mathcal{P}} \mathbb{E}_P[\mathcal{L}(\cdot, m_b) - \mathcal{L}(\cdot, m_a) | \mathbf{x}^*] > 0$$

if $\mathcal{L}(\cdot, \cdot)$ is 1/0 loss function, so :

$$m_a \succ_M m_b \iff \inf_{P \in \mathcal{P}} \frac{P(y = m_a | X = \mathbf{x}^*)}{P(y = m_b | X = \mathbf{x}^*)} > 1$$

Decision Making in Imprecise Probabilities

By applying Bayes theorem on $P(y = m_a | X = \mathbf{x}^*)$, so :

$$m_a \succ_M m_b \iff \inf_{P_{X|y} \in \mathcal{P}_1, P_y \in \mathcal{P}_2} \frac{P(\mathbf{x}^* | y = m_a) P(y = m_a)}{P(\mathbf{x}^* | y = m_b) P(y = m_b)} > 1$$

The resulting set of cautions decisions is :

$$Y_M = \{m_a \in \mathcal{K} \mid \nexists m_b : m_a \succ_M m_b\}$$

For instance, if $\mathcal{K} = \{m_a, m_b, m_c\}$, we can have :

$$\hat{Y}_M = \{m_a \succ_M m_b, m_c \succ_M m_b, m_a \succ_{<M} m_c\} = \{m_a, m_c\}$$

Imprecise Linear Discriminant Analysis (ILDA)

Objective :

Making imprecise the parameter mean μ_k of each Gaussian distribution family $\mathcal{G}_k := \mathbb{P}_{X|Y=m_k} \sim \mathcal{N}(\mu_k, \hat{\Sigma})$

Assumptions :

- Covariances precisely estimated and Homoscedasticity, i.e. $\Sigma_k = \Sigma$:

$$\hat{\Sigma} = \frac{1}{(N-K)} \sum_{k=1}^K \sum_{i=1}^{n_k} (x_{i,k} - \bar{\mathbf{x}}_k)(x_{i,k} - \bar{\mathbf{x}}_k)^t$$

- Prior probabilities precisely estimated : $\hat{\pi}_k = \frac{n_k}{N}$

Decision Making in ILDA

We take the previously maximality criterion and assumptions, so :

$$\begin{aligned}
 m_a \succ_M m_b &\iff \inf_{P_{X|Y} \in \mathcal{D}_1, P_Y \in \mathcal{D}_2} \frac{P(\mathbf{x}^* | y = m_a) P(y = m_a)}{P(\mathbf{x}^* | y = m_b) P(y = m_b)} > 1 \quad (3) \\
 &\iff \inf_{P_{X|Y} \in \mathcal{D}_1} \frac{P(\mathbf{x}^* | y = y_a) \hat{\pi}_a}{P(\mathbf{x}^* | y = y_b) \hat{\pi}_b} > 1
 \end{aligned}$$

Given $\mathcal{G}_k := \mathbb{P}_{X|Y=m_k} \sim \mathcal{N}(\mu_k, \hat{\Sigma})$ are independent :

$$\iff \frac{\inf_{P \in \mathcal{G}_a} P(\mathbf{x}^* | y = y_a) \hat{\pi}_a}{\sup_{P \in \mathcal{G}_b} P(\mathbf{x}^* | y = y_b) \hat{\pi}_b} > 1$$

Decision Making in ILDA (cont...)

Then, the problem reduces to two optimisation problems :

$$\underline{P}(\mathbf{x}^* | y = y_a) = \inf_{P \in \mathcal{G}_a} P(\mathbf{x}^* | y = y_a) \quad (4)$$

$$\overline{P}(\mathbf{x}^* | y = y_b) = \sup_{P \in \mathcal{G}_b} P(\mathbf{x}^* | y = y_b) \quad (5)$$

As $\mathbb{P}_{X|Y=m_k} \sim \mathcal{N}(\mu_k, \hat{\Sigma})$ and $\Sigma_b = \hat{\Sigma}$, so :

$$\underline{P}(\mathbf{x}^* | y = y_a) \iff \underline{\mu}_a = \inf_{P \in \mathcal{G}_a} -\frac{1}{2}(\mathbf{x}^* - \mu_a)^T \hat{\Sigma}^{-1}(\mathbf{x}^* - \mu_a) \quad (6)$$

$$\overline{P}(\mathbf{x}^* | y = y_b) \iff \overline{\mu}_b = \sup_{P \in \mathcal{G}_b} -\frac{1}{2}(\mathbf{x}^* - \mu_b)^T \hat{\Sigma}^{-1}(\mathbf{x}^* - \mu_b) \quad (7)$$

Imprecise Linear Discriminant Analysis

Now, the question is : How could we make imprecise the unknown mean parameter μ_k ?

- Confidence intervals.
- Neighbors around μ_k .
- P-Box
- Robust Bayesian
-

*We would use **robust Bayesian with conjugate distributions for exponential families***

Imprecise Linear Discriminant Analysis

Bayesian inference context

In classic Bayesian inference is based on two components :

- The distribution of the observed data conditional on its unknown parameters (or Likelihood).
- A belief information of expert (or prior distribution).

In order to build procedures of posterior inference on the unknown parameter, in this case μ_k .

$$p(\mu_k | X, \mathbf{y} = m_k) \propto p(X | \mu_k, \mathbf{y} = m_k)p(\mu_k) \quad (8)$$

Where $p(\mu_k) \in \mathcal{P}_{\mu_k}$ could belong a set of prior distributions \mathcal{P}_{μ_k}

Imprecise Linear Discriminant Analysis

We propose to use a set of prior distributions based on near-ignorance approach of [6, eq. 16] :

$$\mathcal{M}_0^\mu = \left\{ \mu \in \mathbb{R}^d \mid p(\mu|m) \propto \exp(\ell^T \mu), m = [\ell_1, \dots, \ell_d]^T \in \mathbb{L} \right\} \quad (9)$$

where m is a hyper-parameter which belong to convex space \mathbb{L} :

$$\mathbb{L} = \left\{ \ell \in \mathbb{R}^d : \ell_i \in [-c_i, c_i], c_i > 0, i = \{1, \dots, d\} \right\}$$

[6] Alessio BENAVALI et Marco ZAFFALON. "Prior near ignorance for inferences in the k-parameter exponential family". In : *Statistics* 49.5 (2015), p. 1104-1140

Remark

\mathcal{M}_0^μ satisfy the four minimal properties that model of prior ignorance require : invariance, near-ignorance, learning and convergence (more details [6]).

Imprecise Linear Discriminant Analysis

By applying Baye's rule (8) (or equation [6, eq 17]), we get a set of posterior distribution :

$$\mathcal{M}_{n_k}^{\mu_k} = \left\{ \mu_k | \bar{\mathbf{x}}_{n_k}, m \propto \mathcal{N} \left(\frac{\ell + n_k \bar{\mathbf{x}}_{n_k}}{n_k}, \frac{1}{n_k} \hat{\Sigma} \right), \right\} \quad (10)$$

where $\bar{\mathbf{x}}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} \mathbf{x}_{i,k}$ and $\ell \in \mathbb{L}$, and :

$$\inf_{\mathcal{M}_{n_k}^{\mu_k}} \mathbb{E}[\boldsymbol{\mu}_k | \bar{\mathbf{x}}_{n_k}, \ell] = \underline{\mathbb{E}}[\boldsymbol{\mu}_k | \bar{\mathbf{x}}_{n_k}, m] = \frac{-\ell + n_k \bar{\mathbf{x}}_{n_k}}{n} \quad (11)$$

$$\sup_{\mathcal{M}_{n_k}^{\mu_k}} \mathbb{E}[\boldsymbol{\mu}_k | \bar{\mathbf{x}}_{n_k}, \ell] = \bar{\mathbb{E}}[\boldsymbol{\mu}_k | \bar{\mathbf{x}}_{n_k}, m] = \frac{\ell + n_k \bar{\mathbf{x}}_{n_k}}{n_k} \quad (12)$$

Imprecise Linear Discriminant Analysis

The two last estimations describe a convex set around μ :

$$\mathbb{G}_k = \left\{ \hat{\mu}_k \in \mathbb{R}^d \left| \begin{array}{l} \hat{\mu}_{i,k} \in \left[\frac{-c_i + n_k \bar{\mathbf{x}}_{i,n_k}}{n_k}, \frac{c_i + n_k \bar{\mathbf{x}}_{i,n_k}}{n_k} \right], \\ \forall i = \{1, \dots, d\} \end{array} \right. \right\}$$

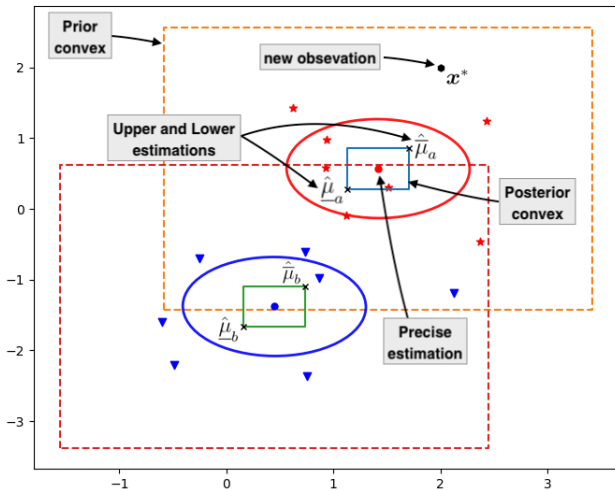
That we use as constraint in on our two optimisation problems.

$$\underline{P}(\mathbf{x}^* | y = m_a) \iff \underline{\hat{\mu}}_a = \arg \max_{\hat{\mu}_a \in \mathbb{G}_a} \frac{1}{2} \hat{\mu}_a^T \hat{\Sigma}^{-1} \hat{\mu}_a + \mathbf{x}^{*T} \hat{\Sigma}^{-1} \hat{\mu}_a \quad (\text{NPQB})$$

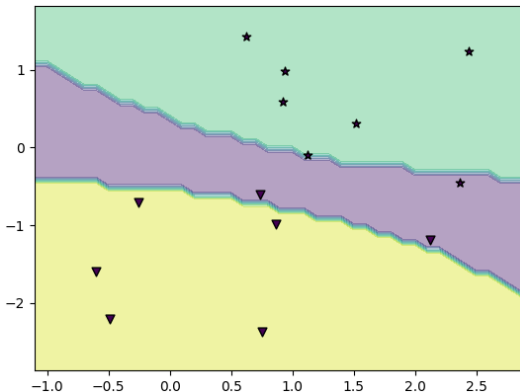
$$\overline{P}(\mathbf{x}^* | y = m_b) \iff \overline{\hat{\mu}}_b = \arg \min_{\hat{\mu}_b \in \mathbb{G}_b} \frac{1}{2} \hat{\mu}_b^T \hat{\Sigma}^{-1} \hat{\mu}_b + \mathbf{x}^{*T} \hat{\Sigma}^{-1} \hat{\mu}_b \quad (\text{PQB})$$

First problem non-convex \rightarrow solved through B&B method.

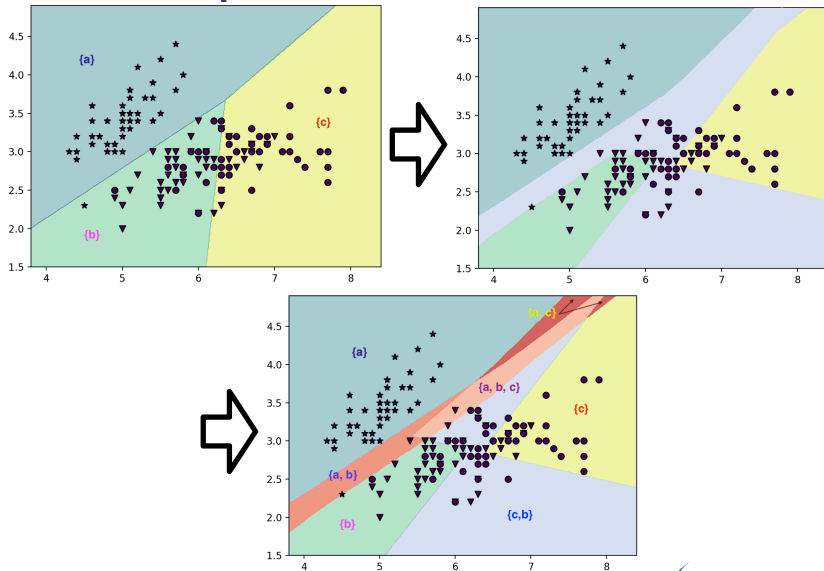
Example



Example (cont..)



Another Example with 3 class



Experiments

Average utility-discounted accuracy measure of [4]

$$u(y, \hat{Y}_M) = \begin{cases} 0 & \text{if } y \notin \hat{Y}_M \\ \frac{\alpha}{|\hat{Y}_M|} - \frac{\beta}{|\hat{Y}_M|} & \text{else} \end{cases}$$

Where u_{65} with $(\alpha, \beta) = (1.6, 0.6)$ and u_{80} with $(\alpha, \beta) = (2.2, 1.2)$.

#	Name	# Obs.	# Regr.	# Classes	#	DLA	IDLA		Inference time
							u_{65}	u_{80}	
a	iris	150	4	3	a	0.961	0.969	0.975	0.56 sec.
b	seeds	210	7	3	b	0.959	0.959	0.962	1.50 sec.
c	glass	214	9	6	c	0.594	0.589	0.642	8.66 sec

Overview

- Classification
 - Decision Making
 - Discriminant Analysis
- Imprecise Classification
 - Imprecise Decision
 - Imprecise Linear discriminant analysis
- Future work
- Conclusions

Imprecise Quadratic discriminant analysis

(1) Release homoscedasticity assumption, i.e. $\Sigma_k \neq \Sigma$

$$\begin{aligned}
 \hat{\underline{\mu}}_a &= \arg \max \frac{1}{2} \hat{\underline{\mu}}_a^T \hat{\Sigma}_k^{-1} \hat{\underline{\mu}}_a - \mathbf{x}^{*T} \hat{\Sigma}_k^{-1} \hat{\underline{\mu}}_a \\
 \text{s.t. } \frac{-c_j + n\bar{x}_{j,n}}{n} &\leq \hat{\underline{\mu}}_{j,a} \leq \frac{c_j + n\bar{x}_{j,n}}{n} \\
 \forall j &= \{1, \dots, d\}
 \end{aligned}
 \tag{PQB}$$

- Making imprecise $P(y = m_a) = [\underline{P}(y = m_a), \bar{P}(y = m_a)]$ and to solve :

$$\inf_{P_{X|y} \in \mathcal{P}_1, P_y \in \mathcal{P}_2} \frac{P(\mathbf{x}^* | y = m_a) P(y = m_a)}{P(\mathbf{x}^* | y = m_b) P(y = m_b)} > 1$$

Imprecise Quadratic discriminant analysis

Space Convex Matrices \mathbf{S}_+

(2) Make imprecise the covariance matrix (i.e. Σ_k or Σ) by using a prior Wishart distribution :

$$\underline{\Sigma}_k = \inf_{\Omega \in \mathbf{S}_+^n} \mathbb{E}[\Sigma_k | X, y = m_k, \tau_0, \Omega] \quad (13)$$

$$\underline{\Sigma}_k = \inf_{\Omega \in \mathbf{S}_+^n} \frac{\Omega + (n-1)\hat{\Sigma}_k^{\text{MLE}}}{n + \tau_0} \quad (14)$$

where $\hat{\Sigma}_k^{\text{MLE}}$ is the maximal likelihood estimator of covariance matrix Σ_k and \mathbf{S}_+^n is a convex space of families of positive semi-definite positive matrices.

Imprecise Quadratic discriminant analysis

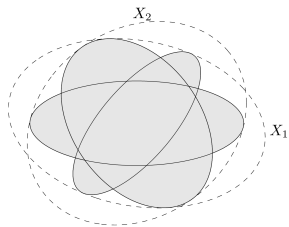
Space Convex Matrices \mathbf{S}_+

In [2], we can find a good intuitions for minimize the last optimization problem, where Φ_ϵ is a perturbation in the neighbourhood of Ω_0 prior parameter value, and $\|\cdot\|_F$ is Frobenius norm.

$$\arg \min_{\Omega_0 \in \mathbf{S}_+^n} \underline{\Sigma} = \frac{\Omega_0 + (n-1)\widehat{\Sigma}_e}{n + \tau_0}$$

$$\text{s.t. } \underline{\Sigma} \succeq X_i, \quad \forall X_i \in \mathbf{S}_+^n, i = \{1, \dots, m\}$$

$$\mathbf{S}_+^n = \{\Omega_0 \mid \|\Omega_0 - \Phi_\epsilon\|_F \preceq \Omega_0 \preceq \|\Omega_0 + \Phi_\epsilon\|_F\}$$



Imprecise Quadratic discriminant analysis

Space Convex of eigenvalues or eigenvectors

(3) Imprecise eigenvalues and eigenvectors of Σ_k .

We'll propose to use $\hat{\Omega}$ estimation of [3, §3], i.e $\hat{\Omega} = \frac{\text{tr}(\Sigma_k^{\text{MLE}})}{d}$, and then applying it the spectral decomposition :

$$\frac{\hat{\Omega} + (n-1)\hat{\Sigma}_k^{\text{MLE}}}{n + \tau_0} \iff \frac{\text{tr}(\sum_{j=1}^d \lambda_j u_j u_j^t)}{d(n + \tau_0)} \mathbb{1} + \frac{n-1}{n + \tau_0} \sum_{j=1}^d \lambda_j u_j u_j^t \quad (15)$$

$$\sum_{j=1}^d \lambda_j \left[\frac{\text{tr}(u_j u_j^t)}{d} \mathbb{1} + (n-1)u_j u_j^t \right] \quad (16)$$

Imprecise Quadratic discriminant analysis

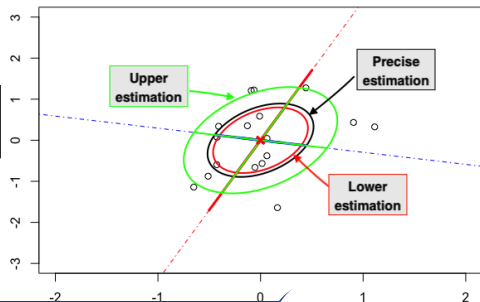
Space Convex of eigenvalues or eigenvectors

In [3], it has been proven that eigenvalues have estimations either biased high (overestimated) or biased low (underestimated) for small and noisy samples.

Then, we could assume the variability of direction is “correctly” estimated (i.e eigenvectors)

$$\bar{\lambda} = \arg \max_{\lambda \in \mathbf{S}_+^n} \sum_{j=1}^d \lambda_j \left[\frac{\text{tr}(u_j u_j^t)}{d} \mathbb{1} + (n-1) u_j u_j^t \right]$$

$$\text{s.t. } \mathbf{S}_+^n = \left\{ \hat{\Sigma} = \sum_{j=1}^d \lambda_j v v^t \mid \underline{\Sigma} \leq \hat{\Sigma} \leq \bar{\Sigma} \right\}$$



Overview

- Classification
 - Decision Making
 - Discriminant Analysis
- Imprecise Classification
 - Imprecise Decision
 - Imprecise Linear discriminant analysis
- Future work
- Conclusions

Conclusions

Imprecise Analysis Discriminant Classification

- Increasing in imprecision on the estimators has allowed us to be more cautious in doubt and to improve the prediction of classification [7].
- More experiments with all imprecise components.
- Creation of new imprecise statistic models for a sensibility analysis and a more (cautious) robust prediction.

References



Jerome FRIEDMAN, Trevor HASTIE et Robert TIBSHIRANI. *The elements of statistical learning*. T. 1. Springer series in statistics New York, 2001.



Stephen BOYD et Lieven VANDENBERGHE. *Convex optimization*. Cambridge university press, 2004.



Santosh SRIVASTAVA. *Bayesian Minimum Expected Risk Estimation of Distributions for Statistical Learning*. University of Washington, 2007.



Marco ZAFFALON, Giorgio CORANI et Denis MAUÁ. "Evaluating credal classifiers by utility-discounted predictive accuracy". In : *International Journal of Approximate Reasoning* 53.8 (2012), p. 1282-1301.



James O BERGER. *Statistical decision theory and Bayesian analysis*. Springer Science & Business Media, 2013.



Alessio BENAOLI et Marco ZAFFALON. "Prior near ignorance for inferences in the k-parameter exponential family". In : *Statistics* 49.5 (2015), p. 1104-1140.



Yonatan-Carlos CARRANZA-ALARCON et Sébastien DESTERCKE. "Analyse Discriminante Imprécise basé sur l'inference Bayésienne robuste". In : *27 emes rencontres francophones sur la logique floue et ses applications* (2018).



