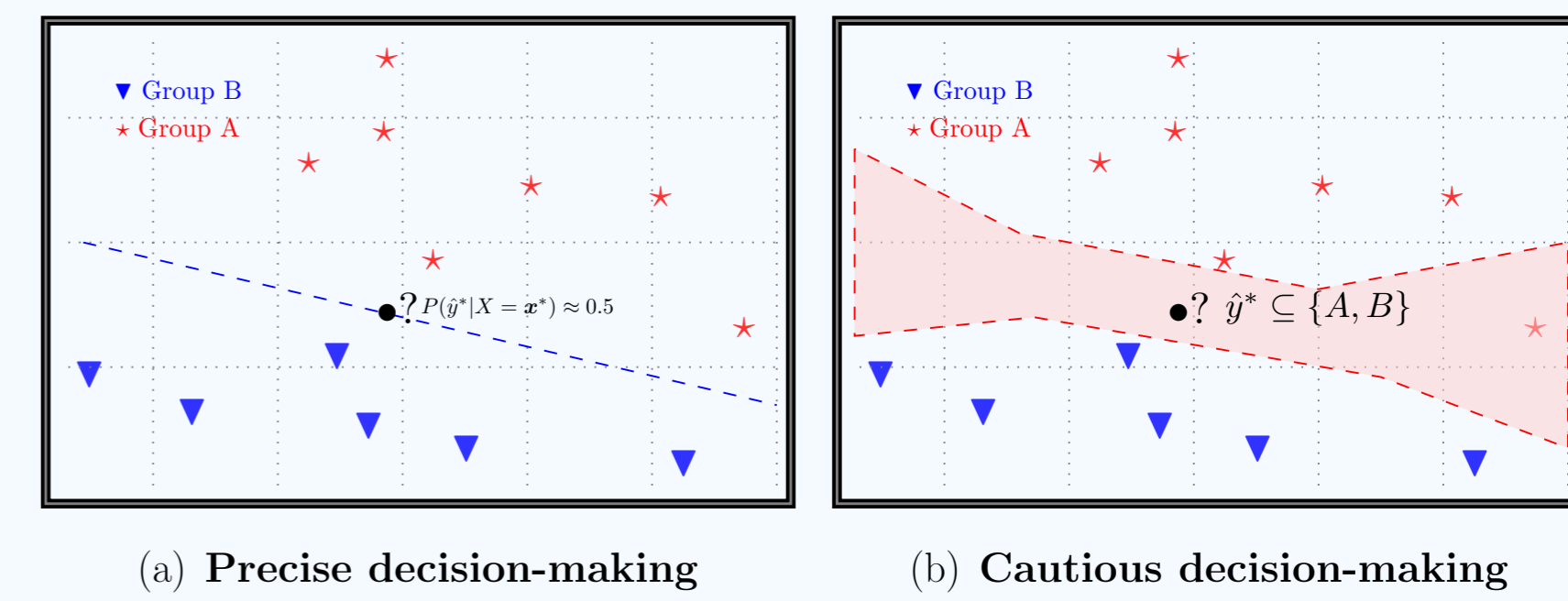


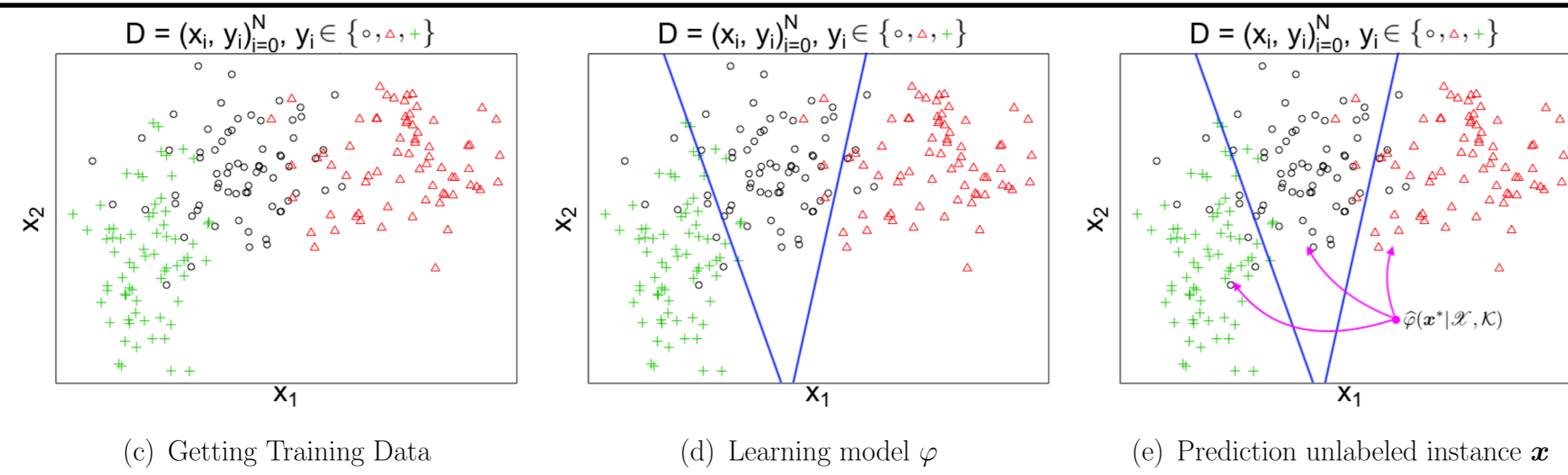
Problem statement

- **Setting:** Training data $D = \{x_i, y_i\}_{i=0}^N \subseteq \mathcal{X} \times \mathcal{K}$ where $\mathcal{X} = \mathbb{R}^p$ and $\mathcal{K} = \{m_1, \dots, m_K\}$
- **Motivation:** avoids mistakes performed by the precise models in hard-to-predict unlabeled instances by making cautious decisions (Figure 1(a) and 1(b))
- **Our proposal:**
 - **Cautious decision:** assigns to a new instance \mathbf{x} a set-valued predictions $\hat{Y} \subseteq \mathcal{K}$ in cases of high uncertainty.
 - **New classifier:** an extension of Gaussian Discriminant analysis aiming to quantify the lack of evidence of the component $\mathbb{P}_{X|Y=m_k}$.



(1) Classification problem

Objective
 Given a training data $D = \{x_i, y_i\}_{i=0}^N$
 Learning a classification rule:
 $\varphi : \mathcal{X} \rightarrow \mathcal{Y}$
 for predicting new observations $\hat{\varphi}(\mathbf{x})$



(2) Decision Making

Precise Decision

Definition 2 (Precise ordering [2, pp. 47])
 Given a general loss function $\mathcal{L}(\cdot, \cdot)$, a conditional probability distribution $\mathbb{P}_{Y|\mathbf{x}}$ and a new unlabeled instance \mathbf{x} , m_a is preferred to m_b , denoted by

$$m_a \succ m_b \iff \mathbb{E}_{\mathbb{P}_{Y|\mathbf{x}}}[\mathcal{L}(\cdot, m_a)] < \mathbb{E}_{\mathbb{P}_{Y|\mathbf{x}}}[\mathcal{L}(\cdot, m_b)]$$

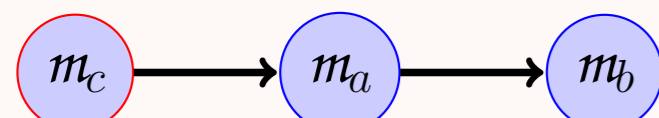
if $\mathcal{L}(\cdot, \cdot)$ is 0/1 loss function, then

$$m_a \succ m_b \iff P(Y = m_a | X = \mathbf{x}) > P(Y = m_b | X = \mathbf{x})$$

Example 1 Given a set of labels $\mathcal{K} = \{m_a, m_b, m_c\}$, a new unlabeled instance \mathbf{x} , and the probability estimates of the conditional distribution $\mathbb{P}_{Y|\mathbf{x}}$:

$$\begin{aligned} \hat{P}(Y = m_a | X = \mathbf{x}) &= 0.3, \\ \hat{P}(Y = m_b | X = \mathbf{x}) &= 0.1, \\ \hat{P}(Y = m_c | X = \mathbf{x}) &= 0.6, \end{aligned}$$

the complete preorder over labels w.r.t. estimated probabilities is $m_c \succ m_a \succ m_b$ where m_c is the maximal predicted label dominating all others.



Cautious Decision

Definition 3 (Partial Ordering by Maximality Criterion [4, §3.2])
 Let $\mathcal{L}(\cdot, \cdot)$ be a general loss function, \mathbf{x} an observed instance and $\mathcal{S}_{Y|\mathbf{x}}$ a set of conditional probability distributions. m_a is preferred to m_b according to the maximality criterion if the cost of exchanging m_a with m_b has a positive lower expectation:

$$m_a \succ_M m_b \iff \inf_{\mathbb{P}_{Y|\mathbf{x}} \in \mathcal{S}_{Y|\mathbf{x}}} \mathbb{E}_{\mathbb{P}_{Y|\mathbf{x}}}[\mathcal{L}(\cdot, m_b) - \mathcal{L}(\cdot, m_a)] > 0 \quad (6)$$

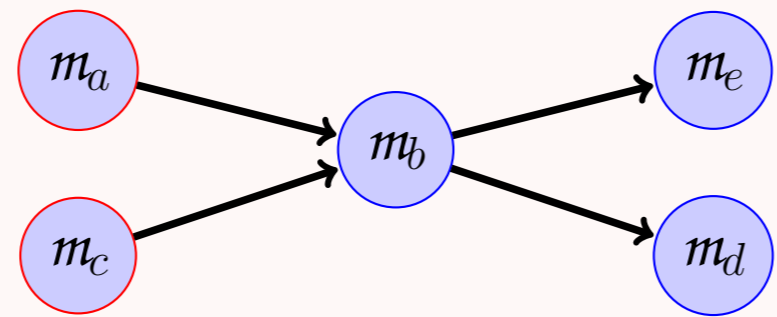
if $\mathcal{L}(\cdot, \cdot)$ is 0/1 loss function, $m_a \succ_M m_b$ if and only if:

$$\inf_{\mathbb{P}_{Y|\mathbf{x}} \in \mathcal{S}_{Y|\mathbf{x}}} P(Y = m_a | \mathbf{x}) - P(Y = m_b | \mathbf{x}) > 0 \quad (7)$$

Example 2 Given a set of labels $\mathcal{K} = \{m_a, m_b, m_c, m_d, m_e\}$ and a possible partial ordering could be the following:

$$\mathcal{B} = \{m_a \succ_M m_b, m_c \succ_M m_b, m_a \succ_M m_c, m_b \succ_M m_d, m_b \succ_M m_e, m_d \succ_M m_e\}$$

where $\hat{Y}_M = \{m_a, m_c\}$ is the predicted set obtained from the set \mathcal{B} of comparisons derived by the criterion of maximality.



(4) Gaussian Discriminant Classification(GDC)

Precise GDC

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

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

Assuming normality on $P_{X|Y=m_k}$:

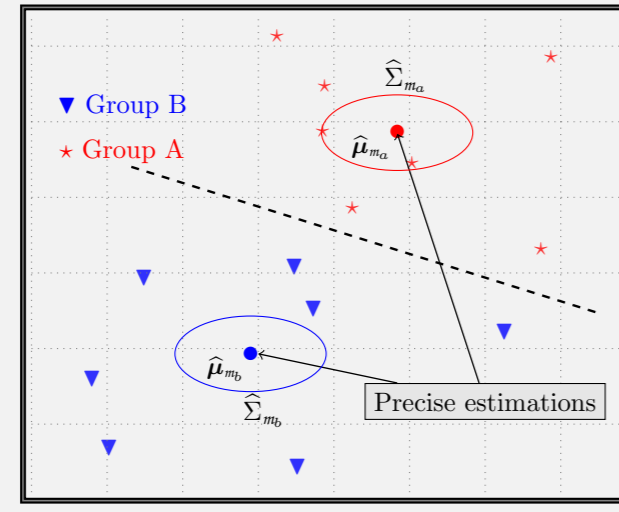
$$\mathcal{G}_{m_k} := \mathbb{P}_{X|Y=m_k} \sim \mathcal{N}(\mu_{m_k}, \Sigma_{m_k}) \quad (8)$$

Defining the marginal distribution as $\pi_{m_k} := P(Y = m_k)$, so under 0/1 loss function, the optimal prediction becomes:

$$\arg \max_{m_k \in \mathcal{K}} \log \pi_{m_k} - \log |\Sigma_{m_k}|^{\frac{1}{2}} - \frac{1}{2}(\mathbf{x}^T - \mu_{m_k})^T \Sigma_{m_k}^{-1} (\mathbf{x}^T - \mu_{m_k}) \quad (9)$$

Estimating parameters by MLE on a subset $\mathcal{D}_{m_k} = \{(x_{i,k}, y_{i,k} = m_k)\}_{i=1, \dots, n_{m_k}} \subseteq \mathcal{D}$:

- $\hat{\pi}_{m_k} = n_{m_k}/N$ (frequency of m_k)
- $\hat{\mu}_{m_k} = \bar{\mathbf{x}}_{m_k}$ (sample mean of \mathcal{D}_{m_k})
- If we assume:
 - (heteroscedasticity) $\rightarrow \hat{\Sigma}_{m_k} = \hat{\Sigma}_{m_k}$ (sample covariance matrix of \mathcal{D}_{m_k})
 - (homoscedasticity) $\rightarrow \hat{\Sigma}_{m_k} = \hat{\Sigma}$ (within-class covariance matrix \mathcal{D})



Imprecise Gaussian Discriminant Classification(IGDC)

Using the maximality criterion and applying Baye's rule, to know whether $m_a \succ_M m_b$, we need to solve

$$\inf_{\mathbb{P}_Y \in \mathcal{P}_Y} \inf_{\mathbb{P}_{X|m_a} \in \mathcal{P}_{X|m_a}} \inf_{\mathbb{P}_{X|m_b} \in \mathcal{P}_{X|m_b}} [P(X = \mathbf{x} | Y = m_a)P(Y = m_a) - P(X = \mathbf{x} | Y = m_b)P(Y = m_b)] > 0 \quad (10)$$

Assuming (A2), an precise estimation for marginal: $P(Y = m_k) := \hat{\pi}_{m_k} > 0$:

$$\inf_{\mathbb{P}_{X|m_a} \in \mathcal{P}_{X|m_a}} \inf_{\mathbb{P}_{X|m_b} \in \mathcal{P}_{X|m_b}} [P(X = \mathbf{x} | Y = m_a) \hat{\pi}_{m_a} - P(X = \mathbf{x} | Y = m_b) \hat{\pi}_{m_b}] > 0 \quad (11)$$

As conditional distributions sets $\mathcal{P}_{X|Y=m_k}$ are independent of each others, then

$$\hat{\pi}_{m_a} \inf_{\mathbb{P}_{X|m_a} \in \mathcal{P}_{X|m_a}} P(X = \mathbf{x} | Y = m_a) - \hat{\pi}_{m_b} \sup_{\mathbb{P}_{X|m_b} \in \mathcal{P}_{X|m_b}} P(X = \mathbf{x} | Y = m_b) > 0 \quad (12)$$

We then have two optimization problems with constraint convex space:

$$\sup_{\mathbb{P}_{X|m_b} \in \mathcal{P}_{X|m_b}} P(X = \mathbf{x} | Y = m_b) \iff \bar{\mu}_{m_b} = \arg \max_{\mu_{m_b} \in \mathcal{G}_{m_b}} -\frac{1}{2}(\mathbf{x} - \mu_{m_b})^T \hat{\Sigma}_{m_b}^{-1} (\mathbf{x} - \mu_{m_b}) \quad (\text{BQP})$$

$$\inf_{\mathbb{P}_{X|m_a} \in \mathcal{P}_{X|m_a}} P(X = \mathbf{x} | Y = m_a) \iff \underline{\mu}_{m_a} = \arg \min_{\mu_{m_a} \in \mathcal{G}_{m_a}} -\frac{1}{2}(\mathbf{x} - \mu_{m_a})^T \hat{\Sigma}_{m_a}^{-1} (\mathbf{x} - \mu_{m_a}) \quad (\text{NBQP})$$

- First problem box-constrained quadratic problem (BQP).
- Second problem non-convex BQP \rightarrow solved through Branch and Bound method.

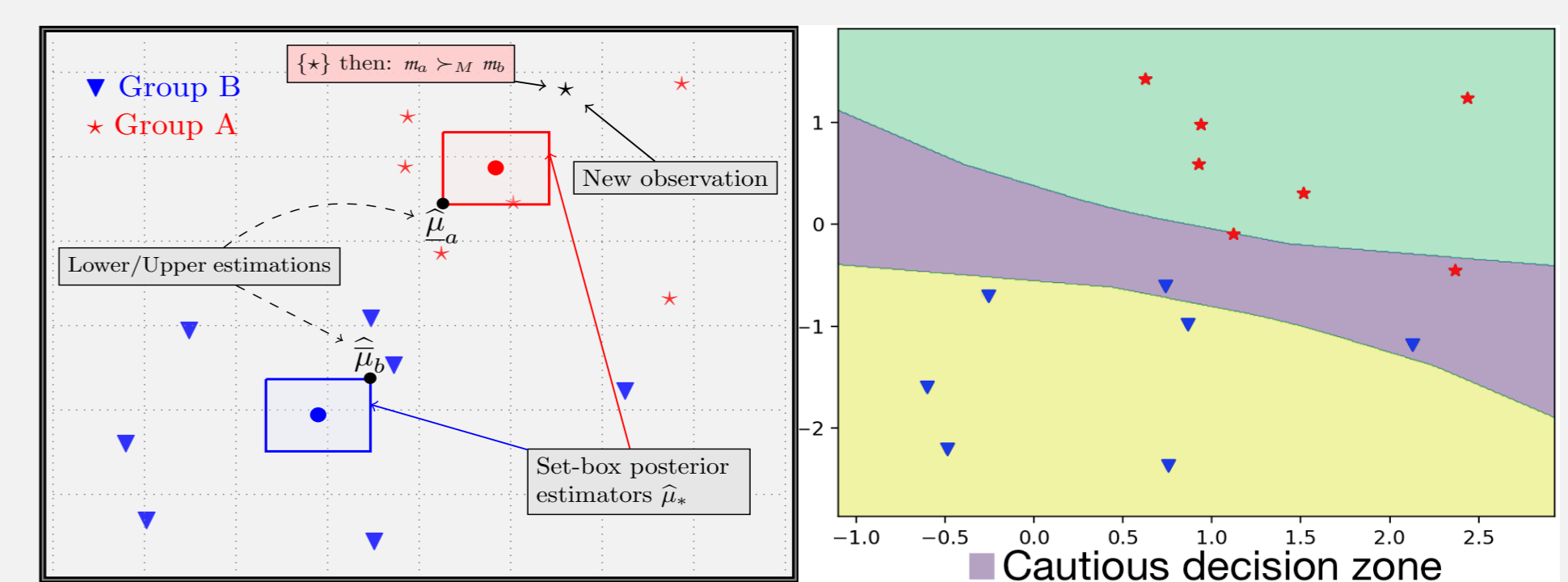


Figure 1: Cautious decision zone for a binary classification

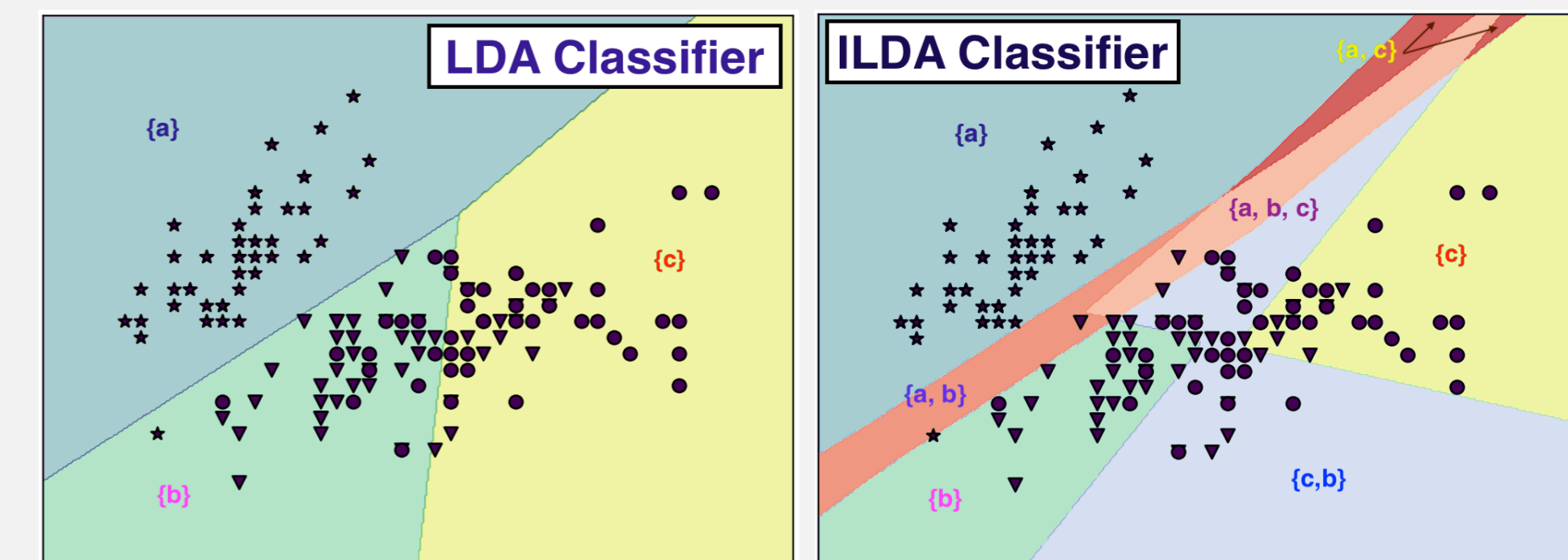


Figure 2: Cautious decision zone with three class $\{a, b, c\}$

(5) Experimental results

- 9 data sets issued from UCI repository.
- 10x10-fold cross-validation procedure.
- Utility-discounted accuracy measure proposed to Zaffalon et al on [5].

$$u(y, \hat{Y}_M) = \begin{cases} 0 & \text{if } y \notin \hat{Y}_M \\ \frac{\alpha}{|\hat{Y}_M|} - \frac{1-\alpha}{|\hat{Y}_M|^2} & \text{else} \end{cases} \implies \text{Reward cautiousness to some degree} \begin{cases} \alpha = 1 : & \text{cautiousness} = \text{randomness} \\ \alpha \rightarrow \infty : & \text{best classifier vacuous} \end{cases}$$

where the usual measures $u_{65}(\cdot, \cdot)$ with $\alpha = 1.6$ and $u_{80}(\cdot, \cdot)$ with $\alpha = 2.2$ have been used in this work.

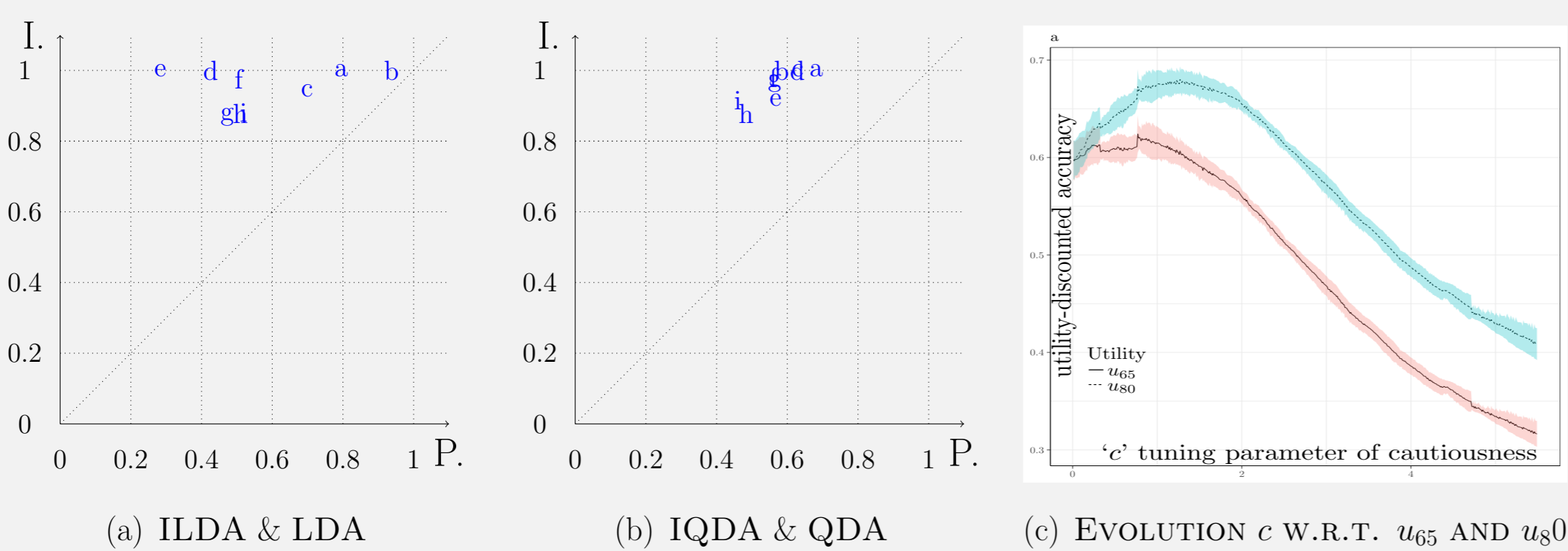


Figure 3: (a) Correctness of the Imprecise LDA in the case of abstention versus accuracy of the Precise LDA. (b) Correctness of the ImpreciseQLDA in the case of abstention versus accuracy of the Precise QDA (graphs are given for the u_{80} accuracies), and the Figure (c) Prediction performance of ILDA model w.r.t. utility-discount accuracy and c tuning parameter on vowel dataset.

#	name	# instances	# features	# labels
a	iris	150	4	3
b	wine	178	13	3
c	forest	198	27	4
d	seeds	210	7	3
e	dermatology	385	34	6
f	vehicle	846	18	4
g	vowel	990	10	11
h	wine-quality	1599	11	6
i	wall-following	5456	24	4

Table 1: Data sets used in the experiments

#	LDA	ILDA	QDA	IQDA	Avg. time (sec.)		
a	97.96	98.38	97.16	97.29	98.08	97.13	0.56
b	98.85	98.99	98.95	99.03	99.39	99.09	1.49
c	94.61	94.56	94.05	89.43	91.77	88.90	12.14
d	96.35	96.59	96.51	94.64	95.20	94.72	1.50
e	96.58	97.06	96.94	82.47	84.24	84.05	19.24
f	77.96	81.98	79.59	85.07	87.96	86.13	3.10
g	60.10	67.45	62.41	87.83	89.96	88.40	4.95
h	59.25	65.83	60.31	55.62	65.85	60.36	34.85
i	67.96	71.34	66.65	65.87	71.79	69.75	10.77
avg.	83.68	86.05	84.03	80.34	87.16	85.33	10.1

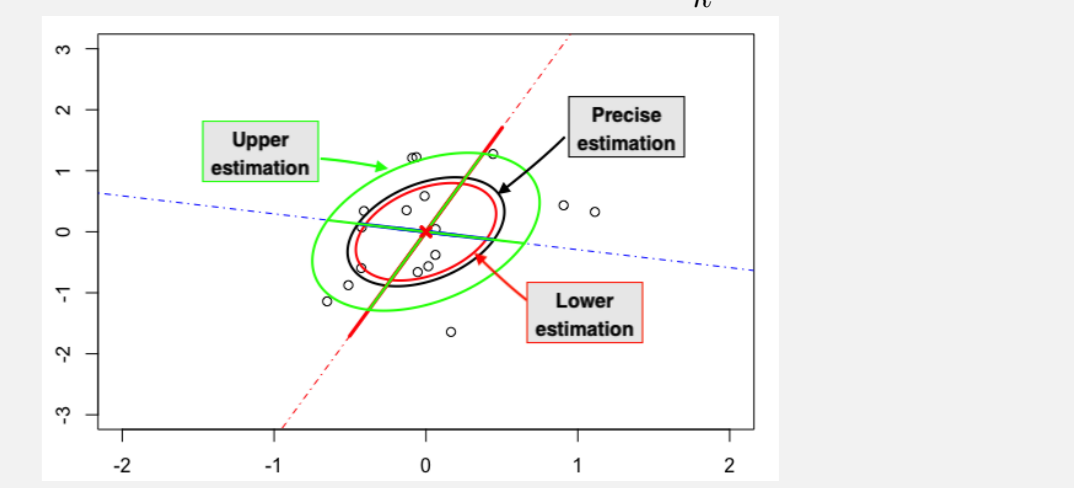
Table 2: Average utility-discounted accuracies (%)

(6) Conclusion and Perspectives

- Increasing imprecision on the estimators has allowed us to be more cautious and to improve the prediction of classification.
- Works done since submission of ISIPTA paper:
 - Considering a diagonal structure of the covariance matrix, i.e. $\Sigma_{m_k} = \sigma_{m_k}^T \mathbb{I}$.
 - Considering a set of marginals distribution \mathcal{P}_Y instead of \mathbb{P}_Y (i.e. release Assumption (A2))
 - Considering the use of a generic loss function instead of zero-one loss function $\mathcal{L}_{0/1}$.

What remains to do:

- Make imprecise the covariance matrix Σ_{m_k} by using the following set of prior distributions
- $$\mathcal{M} \propto \{|\Lambda|^{\frac{\nu_0}{2}} \exp\{-\frac{1}{2}tr(\Lambda \ell^T)\}, \ell \in \mathbb{L}, \ell_i \in [-c_i, c_i], \nu_0 > p\}$$
- Making imprecise the components eigenvalues and eigenvectors of covariance matrix Σ_{m_k} .



Acknowledgments

This work was carried out in the framework of the Labex MS2T, funded by the French Government, through the National Agency for Research (Reference ANR-11-IDEX-0004-02).

References

- [1] Alessio Benavoli and Marco Zaffalon. Prior near ignorance for inferences in the k-parameter exponential family. *Statistics*, 49(5):1104–1140, 2014.
- [2] James O Berger. *Statistical decision theory and Bayesian analysis; 2nd ed.* Springer Series in Statistics. Springer, New York, 1985.
- [3] José M Bernardo and Adrian FM Smith. *Bayesian Theory*. John Wiley & Sons Ltd., 2000.
- [4] Matthias CM Troffaes. Decision making under uncertainty using imprecise probabilities. *International Journal of Approximate Reasoning*, 45(1):17–29, 2007.
- [5] Marco Zaffalon, Giorgio Corani, and Denis Mauá. Evaluating credal classifiers by utility-discounted predictive accuracy. *International Journal of Approximate Reasoning*, 53(8):1282–1301, 2012.