

Imprecise probabilities and machine learning : a tradeoff between accuracy and epistemic uncertainty

M. Serrurier
IRIT, Toulouse, France

May 27, 2015

- Motivations
- Possibility distribution as a family of probability distributions
- Possibilistic loss function
- Possibilistic entropy
- Application to the learning of trees
- Conclusions

- Probability and density estimation widely used in machine learning.
- Limitations : parameters of the distributions are estimated from limited samples of data.
 - More data \rightarrow more complex models
- Model selection : best trade-off between accuracy and epistemic uncertainty
- Idea : use ability of possibility theory to represent epistemic uncertainty

- Possibility distribution π is a mapping from Ω to $[0, 1]$
- Possibility measure : $\forall A \subseteq \Omega, \Pi(A) = \sup_{x \in A} \pi(x)$
 - $\Pi(A \cup B) = \max(\Pi(A), \Pi(B))$
 - $\Pi(A \cap B) \leq \min(\Pi(A), \Pi(B))$
- Necessity measure : $\forall A \subseteq \Omega, N(A) = 1 - \Pi(\bar{A})$.
- States of knowledge :
 - complete knowledge: $\exists x \in \Omega$ such as $\pi(x) = 1$ and $\forall y \in \Omega, y \neq x, \pi(y) = 0$
 - total ignorance: $\forall x \in \Omega, \pi(x) = 1$.

Possibility distribution and upper bound of probability distribution

- Qualitative interpretation : description of imprecise concept (cheap, young, ...)
- Probabilistic interpretation : Upper bound of a family of probability distribution :

$$\mathcal{P}(\pi) = \{p \in \mathcal{P}, \forall A \in \Omega, N(A) \leq P(A) \leq \Pi(A)\}.$$

- States of knowledge :
 - complete knowledge : no uncertainty.
 - total ignorance : all probability distributions are possible.

- Specificity $\pi \preceq \pi'$, if and only if:

$$\pi \preceq \pi' \Leftrightarrow \forall x \in \Omega, \pi(x) \leq \pi'(x)$$

- σ -specificity (discrete case) : $\pi \preceq_{\sigma} \pi'$, if and only if there exists a permutation $\sigma \in \mathcal{S}_q$ such as:

$$\pi \preceq_{\sigma} \pi' \Leftrightarrow \forall x \in \Omega, \pi(x) \leq \pi'(\sigma(x))$$

- Specificity reflects the amount of information encoded by the distribution.

Probability-possibility transformation

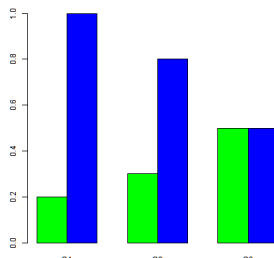
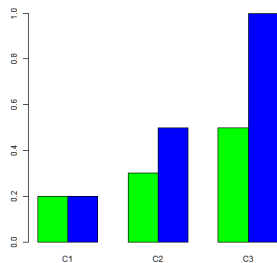
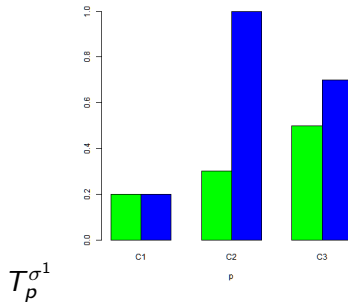
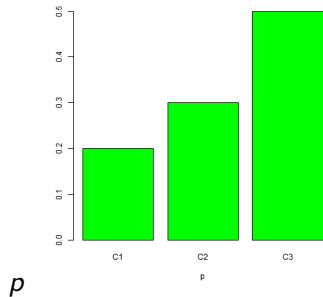
- $\forall \sigma \in S_q$ we have a cumulative distribution T_p^σ which encodes p :

$$\forall j \in \{1, \dots, q\}, T_p^\sigma(C_j) = \sum_{k, \sigma(k) \leq \sigma(j)} p(C_k).$$

$$\forall \sigma \in S_q, p \in \mathcal{P}(T_p^\sigma)$$

- Probability-possibility transformation (T_p^*) : the most σ specific possibility distribution which bounds the distribution
- T_p^* is a cumulative function of p
- $T_p^* = T_p^{\sigma^*}$ where $\sigma^* \in S_q$ follows the probability increasing order

Probability-possibility transformation : example

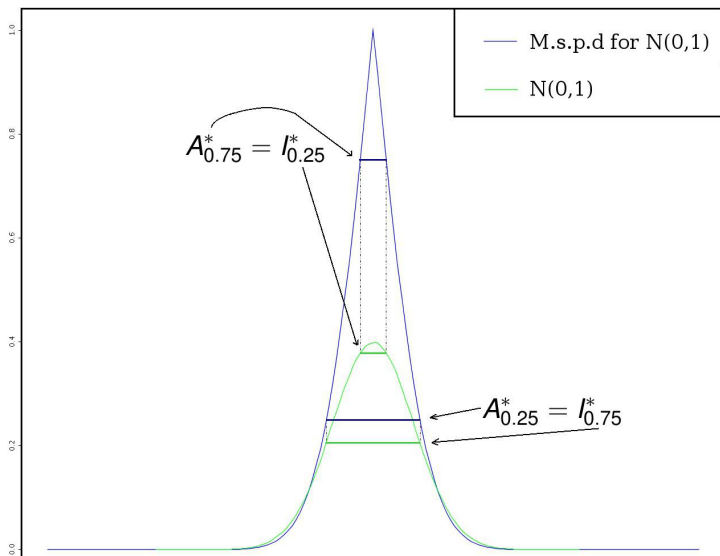


- α -cuts are subsets of Ω such that:

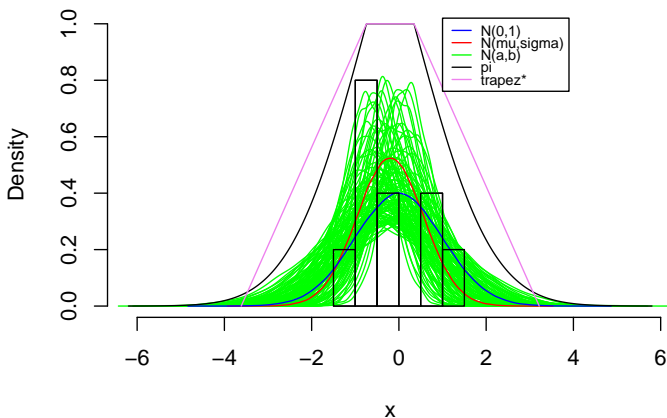
$$A_\alpha(\pi) = \{x \in \Omega, \pi(x) \geq \alpha\}.$$

- Specificity order based on inclusion of α -cuts
- Probability possibility transformation :
 - $T_p^*(x) = \max_{\alpha, x \in I_\alpha} (1 - \alpha)$
 - T_p^* corresponds to the cumulative distribution function of p according to the order the values of $p(x)$.
- α -cuts are upper bounds of $I_{(1-\alpha)}$

Probability-possibility transformation of a Gaussian distribution



Possibilistic distribution encoding uncertainty around Gaussian parameters



Why using possibility distribution ?

- Probability-possibility transformation \rightarrow loss of information but ..
- Possibility distribution can different state of knowledge from complete knowledge to total ignorance
- T_p^* + uncertainty around parameters of $p \rightarrow$ less specific possibility distribution
- σ -specificity respect the entropy order

$$T_p^* \preceq_{\sigma} T_{p'}^* \Rightarrow \mathcal{H}(p) \leq \mathcal{H}(p')$$

- Goal : describe loss and entropy functions that support the σ -specificity order and the probabilistic interpretation of possibility distribution

- Loss functions $\mathcal{L}(f, X)$ measure adequateness between data $X = \{x_1, \dots, x_n\}$ and a distribution f .
- Loss function $\mathcal{L}(f, X)$ is linear w.r.t. X : $\mathcal{L}(f, X) = \frac{\sum_{i=1}^n \mathcal{L}(f, x_i)}{n}$
- common loss functions :
 - Log loss : $\mathcal{L}_{\log}(p|X) = -\sum_{j=1}^q \alpha_j \log(p_j)$.
 - Squared loss : $\mathcal{L}_{sqr}(p|X) = \frac{1}{2} * \sum_{j=1}^q p_j^2 - (\sum_{j=1}^q \alpha_j * p_j)$
- Loss functions are minimal for the frequency distribution (i.e. $p_j = \alpha_j$)

Possibilistic log-loss function

- Principle given a possibility distribution π :
 - consider σ s.a. $\pi(C_{\sigma(1)}) \leq \dots \leq \pi(C_{\sigma(q)})$
 - consider $BC_j = \bigcup_{i=1}^j C_{\sigma(i)}$ and \overline{BC}_j as binary event space
 - $(\pi(C_{\sigma(j)}), 1 - \pi(C_{\sigma(j)}))$ is a probability distribution on $\Omega_j = \{BC_j, \overline{BC}_j\}$
 - apply re-scaled loss function to each Bernoulli distribution
- Poss-log loss :

$$\begin{aligned} \mathcal{L}_{\pi-l}(\pi|X) = & \\ & - \sum_{j=1}^q \left(\frac{cdf_j}{2} * \log\left(\frac{\pi_j}{2}\right) + \left(1 - \frac{cdf_j}{2}\right) * \log\left(1 - \frac{\pi_j}{2}\right) \right). \end{aligned}$$

- Poss-squared loss :

$$\mathcal{L}_{\pi-s}(\pi|X) = \frac{1}{2} \sum_{j=1}^q \pi_j^2 + \sum_{j=1}^q cdf_j - \sum_{j=1}^q \pi_j * cdf_j$$

Possibilistic log-loss function : properties

Linearity

\mathcal{L}_π is linear with respect to X

Optimality for probability possibility transformation

we have $\arg \min(\mathcal{L}_\pi(\pi|X)) = T_{p^\alpha}^*$ (where p^α is the frequency distribution).

Specificity order

$$\forall \sigma \in \mathcal{S}_q, T_{p^\alpha}^\sigma \preceq \pi_1 \preceq \pi_2 \Rightarrow$$

$$\mathcal{L}_\pi(T_{p^\alpha}^\sigma|X) \leq \mathcal{L}_\pi(\pi_1|X) \leq \mathcal{L}_\pi(\pi_2|X)$$

- Direct extension of the discrete case based on α -cuts
- Poss-log loss :

$$\begin{aligned}\mathcal{L}_{\pi-I}(\pi|x) &= - \int_{\mathbb{R}} \log(1 - \pi(x)/2) dx \\ &\quad - 0.5 * \int_{A_{\pi_x}} \log(\pi(x)/2) - \log(1 - \pi(x)/2) dx\end{aligned}$$

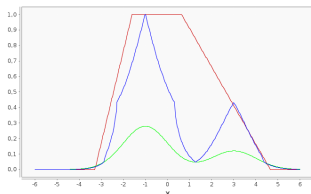
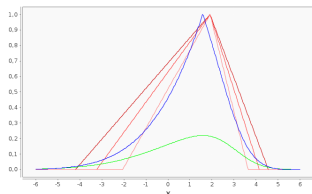
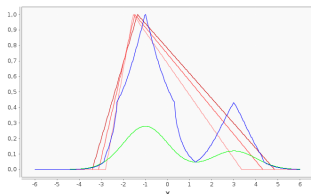
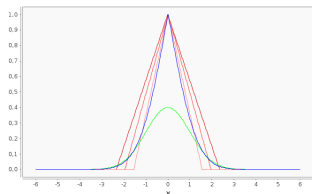
- Poss-squared loss :

$$\mathcal{L}_{\pi-I}(\pi|x) = \int_{A_{\pi(x)}} \pi(t) dt - |A_{\pi(x)}| - \frac{1}{2} \int_{\mathbb{R}} \pi(t)^2 dt$$

- Same properties than in the discrete case

- Gaussian distribution
 - Confidence intervals obtained by mood regions
 - Analytic formulas for the possibility distribution that encodes such family
 - Poss-log loss has to be approximated
- Triangular and trapezoidal possibility distributions
 - Triangular distribution upper bound any unimodal distribution at a given confidence threshold
 - Poss-squared loss easy to compute.

Approximation with Poss-squared loss



- Probabilistic case :
 - The entropy is the loss function value of the frequency distribution (i.e. $\mathcal{H}(p^\alpha) = \mathcal{L}(p^\alpha|X)$)
 - Entropy is maximal for the uniform distribution
 - Entropy is minimal when all the data pertain to the same class
- Expected properties in the possibilistic case
 - The possibilistic entropy applied to probability possibility transformations respects the specificity order.
 - The possibilistic entropy increases when uncertainty around the considered probability distribution increases.

- The possibility cumulative entropy is the entropy of a possibility distribution π with respect to a probability distribution p

$$\begin{aligned} \mathcal{H}_{\pi-I}(p, \pi) &= \\ &= - \sum_{j=1}^q \frac{\frac{T_p^*(C_j)}{2} * \log\left(\frac{\pi(C_j)}{2}\right) + \left(1 - \frac{T_p^*(C_j)}{2}\right) * \log\left(1 - \frac{\pi(C_j)}{2}\right)}{q * \log(q)}. \end{aligned} \quad (1)$$

- Given X and his associated frequency distribution p^α we have
$$\mathcal{H}_{\pi-1}(p^\alpha, \pi) = \frac{\mathcal{L}_{\pi-1}(\pi|X)}{q^* \log(q)}$$

Specificity order

$$T_p^* \preceq T_{p'}^* \Rightarrow \mathcal{H}_{\pi-1}(p, T_p^*) \leq \mathcal{H}_{\pi-1}(p', T_{p'}^*)$$

Increase with uncertainty

$$T_p^* \preceq \pi \preceq \pi' \Rightarrow \mathcal{H}_{\pi-1}(p, T_p^*) \leq \mathcal{H}_{\pi-1}(p, \pi) \leq \mathcal{H}_{\pi-1}(p, \pi')$$

- Given $p(c)$ estimated from n pieces of data, we compute the upper bound $p_{\gamma,n}^*$ of the $(1 - \gamma)\%$ confidence interval with Agresti-Coull method.
- Possibility distribution as an upper bound of the frequency distribution

$$\pi_{p,n}^{\gamma}(C_j) = P_{\gamma,n}^* \left(\bigcup_{i=1}^j C_{\sigma(i)} \right)$$

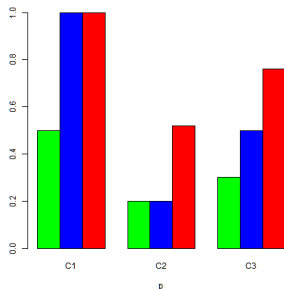
where $\sigma \in S_q$ follows the probability order.

- Possibilistic entropy of a frequency distribution estimated from n pieces of data :

$$\mathcal{H}_{\pi-l}^*(p, n, \gamma) = \mathcal{H}_{\pi-l}(p, \pi_{p,n}^{\gamma})$$

Example

- $p(C_1) = 0.5$, $p(C_2) = 0.2$ and $p(C_3) = 0.3$.
- $n = 10$ ($\gamma = 0.05$)
- $\pi_{p,10}^{0.05}(C_1) = P_{0.05,10}^*(C_1 \cup C_2 \cup C_3) = 1$
- $\pi_{p,10}^{0.05}(C_2) = p_{0.05,10}^*(C_2) = 0.52$
- $\pi_{p,10}^{0.05}(C_3) = P_{0.05,10}^*(C_2 \cup C_3) = 0.76$
- $\mathcal{H}_{\pi-I}^*(p, 50, 0.05) = 0.38$.



Possibilistic cumulative entropy for a limited set of data : properties

Encoding

$$p \in \mathcal{P}(\pi_{p,n}^\gamma)$$

$$\forall n > 0, \pi_p^* \preceq \pi_{p,n}^\gamma \text{ and } \lim_{n \rightarrow \infty} \pi_{p,n}^\gamma = \pi_p^*$$

Increases when uncertainty increases

given $n' \leq n$ we have

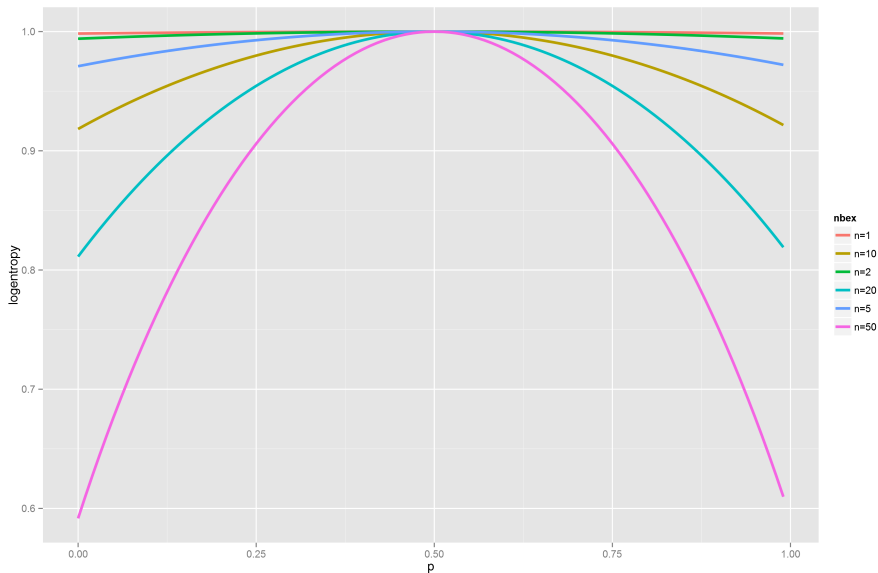
$$\forall \gamma \in]0, 1[, \mathcal{H}_{\pi-l}^*(p, n, \gamma) \leq \mathcal{H}_{\pi-l}^*(p, n', \gamma)$$

Stability

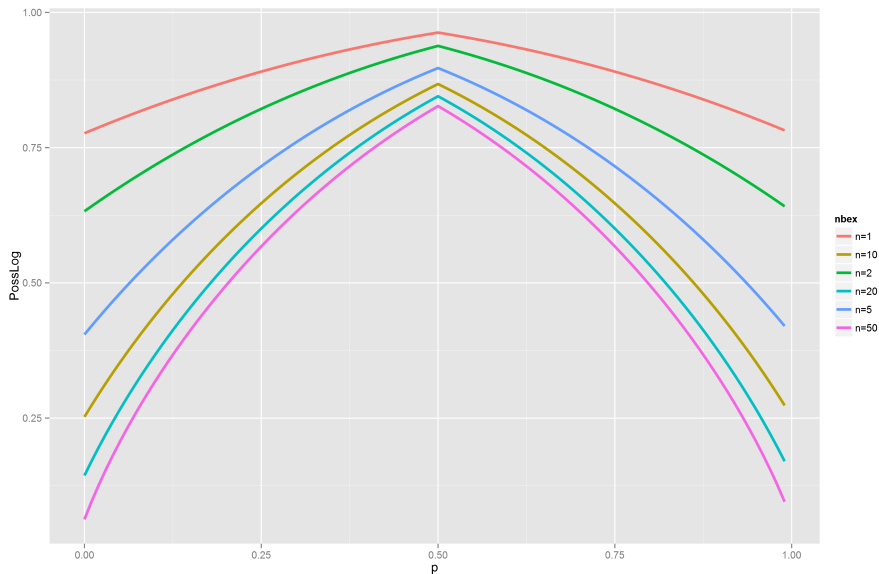
given p and p' we have

$$\forall \gamma \in]0, 1[, T_p^* \preceq T_{p'}^* \Rightarrow \mathcal{H}_{\pi-l}^*(p, n, \gamma) \leq \mathcal{H}_{\pi-l}^*(p', n, \gamma)$$

Bayesian revision

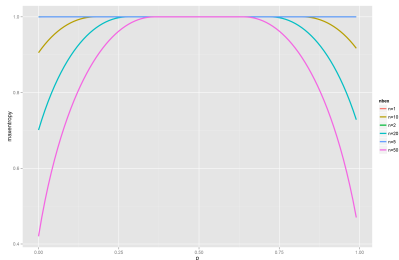


Possibilistic cumulative entropy

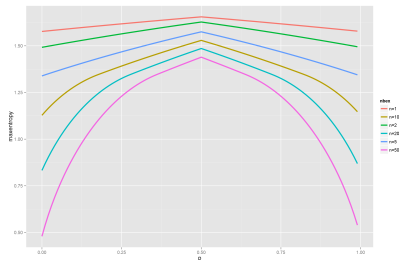


Upper entropy on credal sets

Without specificity term



With specificity term



- Learning of decision trees is based on entropy of frequency distributions
- When we go deeper downward the tree, the examples become rarer and the faithfulness of entropy decreases
- Log entropy-based gain : splitting a node always decreases the weighted entropy of the leaves obtained

- Principal : recursively choose the attribute that maximize the gain function
- Log gain function :

$$G(Z, A) = \mathcal{H}(p_Z) - \sum_{k=1}^r \frac{|v_k|}{n} \mathcal{H}(p_{v_k})$$

- Possibilistic gain function :

$$G_{\gamma}^{\pi}(Z, A) = \mathcal{H}_{\pi-1}^*(p_Z, n, \gamma) - \sum_{k=1}^r \frac{|v_k|}{n} \mathcal{H}_{\pi-1}^*(p_{v_k}, |v_k|, DS(\gamma, r)).$$

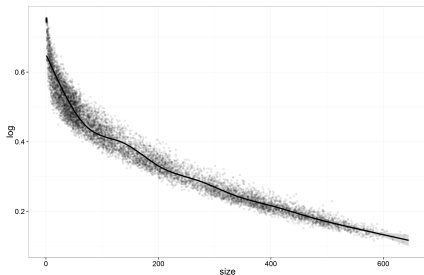
Advantages of the approach

- Significant choices of split
- Statistically relevant stopping criterion
- Reasonable estimator of the performances of a decision tree
- Provide well sized and well balanced trees

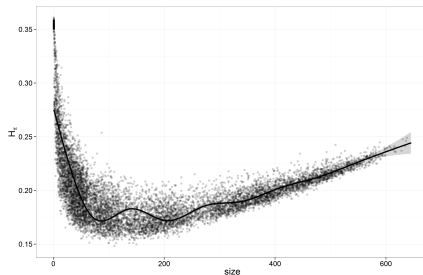
- Algorithm :
 - ① browse recursively the tree to the corresponding leaf
 - ② add x to the set of examples
 - ③ search the attribute with the best G_{γ}^{π}
 - ④ if the gain is positive, create a new node with the corresponding attribute, else do nothing.
- Advantages :
 - Incremental algorithm
 - Built new leaves only when the split gives a statistical significant gain
 - Only consider the leaf concerned by the new example

Experiments : Entropy vs size

size vs log loss

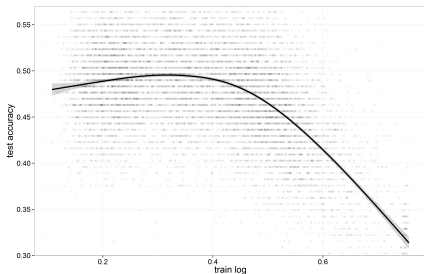


size vs possibilistic log-loss

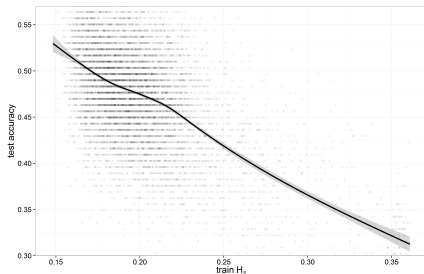


Experiments : entropy vs accuracy

Accuracy vs log-loss



Accuracy vs possibilistic log-loss

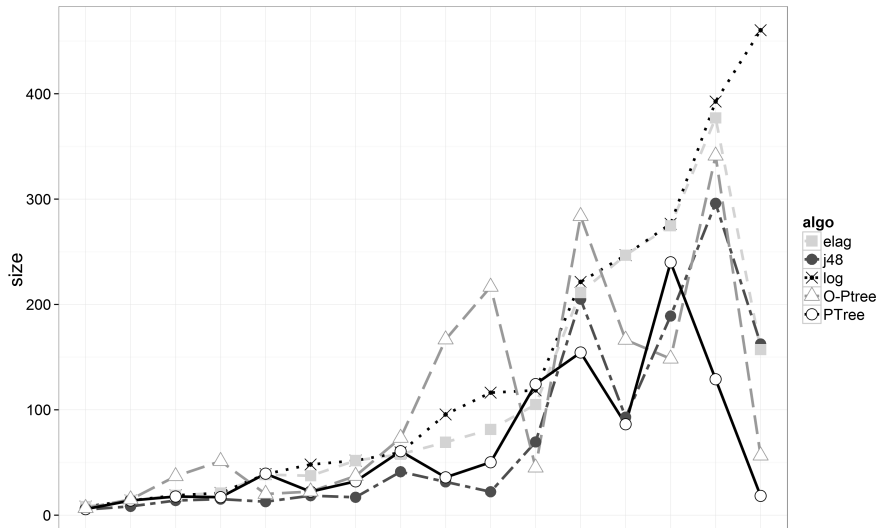


Accuracy comparison

Data set	Log Tree	PrunTree	Π Tree	O- Π Tree	J48
soybean	89.4 \pm 5.0	89.4 \pm 5.0	<u>94.0\pm2.8</u>	89.0 \pm 3.8	91.7 \pm 3.1
lymph	72.9 \pm 11.8	72.9 \pm 11.8	<u>78.3\pm7.9</u>	<u>78.3\pm8.2</u>	75.8 \pm 11.0
zoo2	97.0\pm4.8	97.0\pm4.8	97.0\pm4.8	96.0 \pm 5.1	92.6 \pm 7.3
ilpd	67.9 \pm 5.5	67.4 \pm 5.6	69.9\pm5.3	66.8 \pm 4.7	68.1 \pm 5.6
yeast	52.0 \pm 4.1	57.0 \pm 3.3	57.1\pm3.4	56.7 \pm 3.6	56.6 \pm 3.7
waveform	75.2 \pm 1.5	75.3 \pm 1.5	<u>77.4\pm1.5</u>	72.6 \pm 1.8	75.2 \pm 1.9
diabetes	68.7 \pm 5.7	70.4 \pm 4.7	<u>74.3\pm4.4</u>	70.4 \pm 3.4	74.4\pm5.2
banknote	98.3 \pm 1.1	98.3 \pm 1.1	98.3 \pm 1.0	97.4 \pm 2.1	98.5\pm1.0
ecoli	78.9 \pm 7.7	80.4 \pm 7.4	82.4 \pm 7.9	<u>83.6\pm7.2</u>	82.8 \pm 5.7
vehicle	71.6 \pm 4.7	71.6 \pm 4.0	<u>74.1\pm4.1</u>	69.1 \pm 3.1	72.2 \pm 4.3
ionosphere	90.3 \pm 4.7	90.3 \pm 4.7	91.1\pm3.6	87.7 \pm 4.0	89.7 \pm 4.3
segment	96.8 \pm 0.6	96.7 \pm 0.7	96.9\pm1.2	94.7 \pm 1.4	96.7 \pm 1.2
pendigits	96.5\pm0.5	96.4 \pm 0.5	96.4 \pm 0.2	93.2 \pm 1.0	96.5 \pm 0.6
spambase	91.8 \pm 1.2	91.7 \pm 1.3	<u>94.0\pm1.3</u>	90.5 \pm 1.2	92.8 \pm 1.0
breast-wv2	92.9 \pm 2.4	92.9 \pm 2.4	93.9 \pm 3.1	<u>94.7\pm1.6</u>	94.1 \pm 2.5
wine2	92.5 \pm 8.7	92.5 \pm 8.7	93.7 \pm 7.3	<u>94.3\pm8.3</u>	93.2 \pm 5.9

Size comparison

Number leaves of LogTree, PrunTree, Π Tree and O- Π Tree, J4.8 comparison for different databases.



- Same approach than for decision trees with poss-log loss function
- Use possibilistic cumulative entropy of the possibility distribution that encodes the family of Gaussian distribution that has parameters inside the confidence interval based on mood confidence region.
- Online algorithm also works
- Promising results

- Possibility loss functions and entropies
 - Agrees with the probabilistic view of possibility theory
 - Reflects both the entropy of a probability distribution and the uncertainty around the parameters
 - Can be used for upper estimate densities without assuming a particular shape
- Application to decision and regression trees :
 - Provides well balanced and well sized trees
 - Avoids over-fitting
 - Simple and efficient online algorithm
- Easy extension to :
 - Bayesian networks
 - Density estimation
 - Bandwidth selection in knn