

Q: HOW MANY PH.D.'S DOES IT TAKE TO GET A POWERPOINT PRESENTATION TO WORK?



ANSWER: $(n+1)$

WHERE n = THE NUMBER OF ACADEMICS IN THE ROOM WHO THINK THEY KNOW HOW TO FIX IT, AND 1 = THE PERSON WHO FINALLY CALLS THE A/V TECHNICIAN.

WWW.PHDCOMICS.COM

Uncertainty reasoning and machine learning

Uncertainty, Decision and Evaluation Revisited

Vu-Linh Nguyen

**Chaire de Professeur Junior, Laboratoire Heudiasyc
Université de technologie de Compiègne**

AOS4 master courses

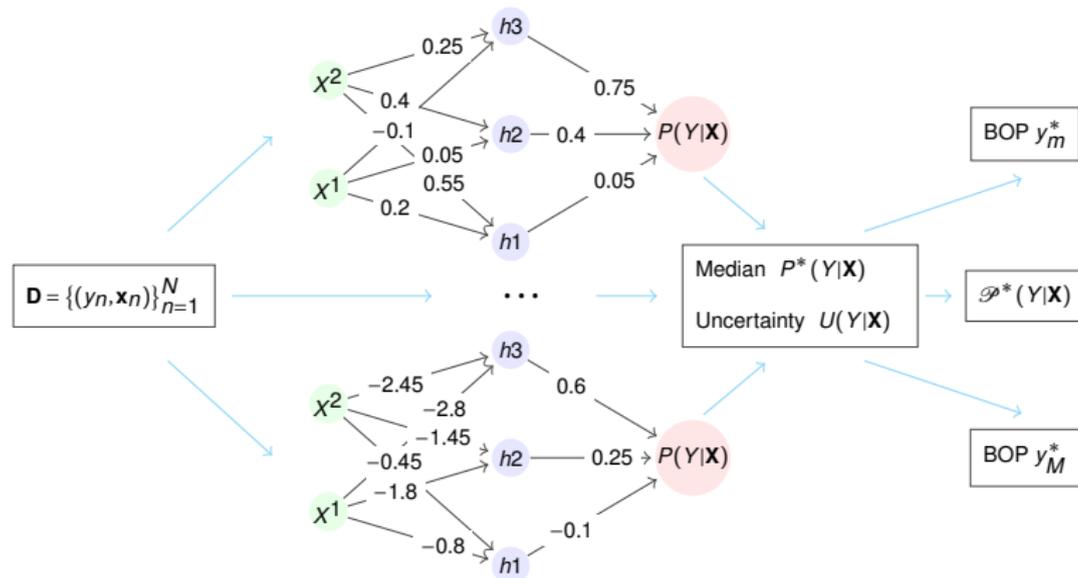
Outline

- Credal ensembling
 - A median classifier: Learning and inference
 - A credal classifier: Learning and inference
- Applications in machine learning
- A few practical aspects
- Exercices on prediction-making

A formal framework [2, 3]

Basic setup:

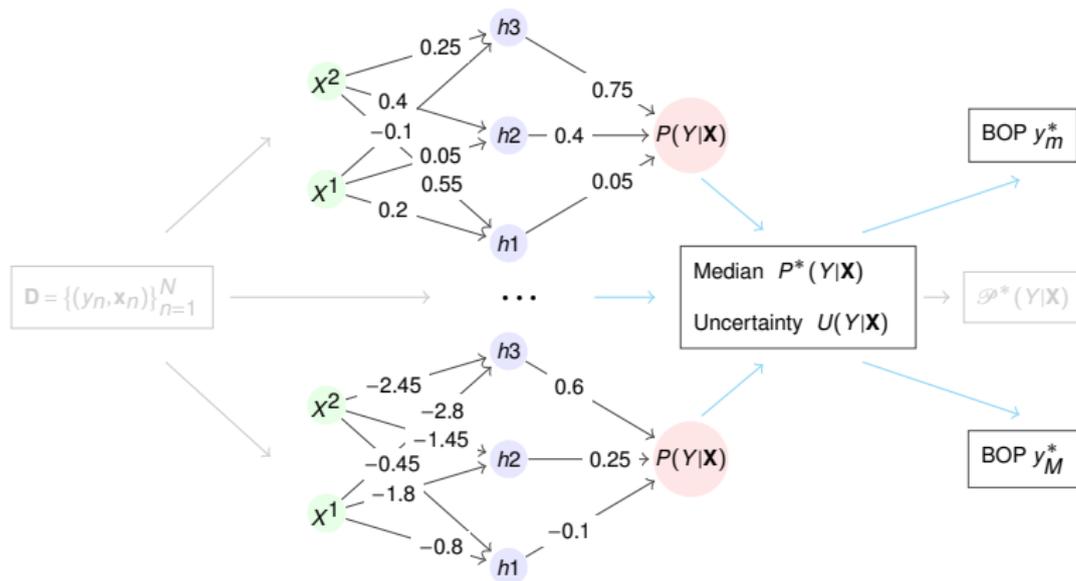
- Features (X^1, \dots, X^P) and a class variables Y
- An finite output space $\mathcal{Y} = \{y^1, \dots, y^C\}$



Outline

- Credal ensembling
 - A median classifier: Learning and inference
 - A credal classifier: Learning and inference
- Applications in machine learning
- A few practical aspects
- Exercices on prediction-making

A median classifier and its predictions [2, 3]



Compute a median classifier

Basic setting:

- An ensemble $\mathbf{H} := \{\mathbf{h}^m | m \in [M] := \{1, \dots, M\}\}$ is made available
- A specified statistical distance d between distributions

A median classifier minimizes the average expected distance:

$$\mathbf{h}_d \in \operatorname{argmin}_{\mathbf{h} \in \mathcal{H}} \mathbf{E} \left[\sum_{m=1}^M d(\mathbf{h}, \mathbf{h}^m) \right] = \operatorname{argmin}_{\mathbf{h} \in \mathcal{H}} \int_{\mathbf{x} \in \mathcal{X}} \left[\sum_{m=1}^M d(\mathbf{h}(\mathbf{x}), \mathbf{h}^m(\mathbf{x})) \right] d\mathbf{x}.$$

Compute a median classifier

Basic setting:

- An ensemble $\mathbf{H} := \{\mathbf{h}^m | m \in [M] := \{1, \dots, M\}\}$ is made available
- A specified statistical distance d between distributions

A median classifier minimizes the average expected distance:

$$\mathbf{h}_d \in \operatorname{argmin}_{\mathbf{h} \in \mathcal{H}} \mathbf{E} \left[\sum_{m=1}^M d(\mathbf{h}, \mathbf{h}^m) \right] = \operatorname{argmin}_{\mathbf{h} \in \mathcal{H}} \int_{\mathbf{x} \in \mathcal{X}} \left[\sum_{m=1}^M d(\mathbf{h}(\mathbf{x}), \mathbf{h}^m(\mathbf{x})) \right] d\mathbf{x}.$$

If no constraint on \mathcal{H} , \mathbf{h}_d can be defined in an instance-wise manner:

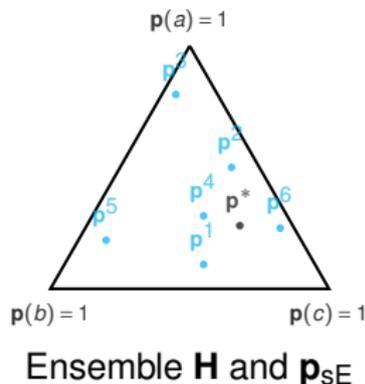
$$\mathbf{h}_d(\mathbf{x}) \in \operatorname{argmin}_{\mathbf{h}(\mathbf{x}) \in \Delta^K} \sum_{m=1}^M d(\mathbf{h}(\mathbf{x}), \mathbf{h}^m(\mathbf{x})). \quad (1)$$

Compute a median classifier (cont.)

For each \mathbf{x} , dropping \mathbf{x} and denoting $\mathbf{p} = \mathbf{h}$ give

$$\mathbf{p}_d \in \operatorname{argmin}_{\mathbf{p} \in \Delta^K} \sum_{m=1}^M d(\mathbf{p}, \mathbf{p}^m). \quad (2)$$

Examples of d are squared Euclidean distance (sE), L_1 distance, and KL divergence.

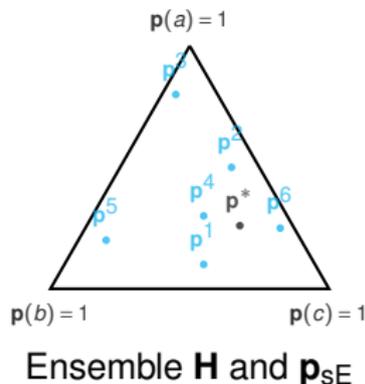


Compute a median classifier (cont.)

For each \mathbf{x} , dropping \mathbf{x} and denoting $\mathbf{p} = \mathbf{h}$ give

$$\mathbf{p}_d \in \operatorname{argmin}_{\mathbf{p} \in \Delta^K} \sum_{m=1}^M d(\mathbf{p}, \mathbf{p}^m). \quad (2)$$

Examples of d are squared Euclidean distance (sE), L_1 distance, and KL divergence.



For any convex distance d :

- The convex optimization problem (2) can be solved using any solver.
- Close-form solution $\mathbf{p}_{sE} =$ averaging the distributions class-wise.

Bayesian-optimal predictions

Basic set (instance-wise manner):

- The median distribution \mathbf{p}_d is given.
- A the higher the better utility $u : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}_+$

A **Bayesian-optimal prediction (BOP)** of u is

$$y_d^u \in \operatorname{argmax}_{y' \in \mathcal{Y}} \mathbf{E}[u(y', y)] = \operatorname{argmax}_{y' \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} u(y', y) \mathbf{p}_d(y). \quad (3)$$

Bayesian-optimal predictions

Basic set (instance-wise manner):

- The median distribution \mathbf{p}_d is given.
- A the higher the better utility $u : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}_+$

A **Bayesian-optimal prediction (BOP)** of u is

$$y_d^u \in \operatorname{argmax}_{y' \in \mathcal{Y}} \mathbf{E}[u(y', y)] = \operatorname{argmax}_{y' \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} u(y', y) \mathbf{p}_d(y). \quad (3)$$

Commonly used utilities, such as 0/1 and cost-sensitive accuracies:

- Find a BOP (10) takes from $O(K)$ to $O(K^2)$
- A BOP $y_d^{0/1}$ (10) of 0/1 accuracy = a most probable class

Bayesian-optimal set-valued predictions

Basic set (instance-wise manner):

- The median distribution \mathbf{p}_d is given.
- A the higher the better utility $U: \mathcal{Y} \times 2^{\mathcal{Y}} \mapsto \mathbb{R}_+$

A **Bayesian-optimal prediction (BOP)** of U is

$$Y_d^U \in \operatorname{argmax}_{Y' \subset \mathcal{Y}} \mathbf{E}[U(Y', y)] = \operatorname{argmax}_{Y' \subset \mathcal{Y}} \sum_{y \in \mathcal{Y}} U(Y', y) \mathbf{p}_d(y). \quad (4)$$

Bayesian-optimal set-valued predictions

Basic set (instance-wise manner):

- The median distribution \mathbf{p}_d is given.
- A the higher the better utility $U: \mathcal{Y} \times 2^{\mathcal{Y}} \mapsto \mathbb{R}_+$

A **Bayesian-optimal prediction (BOP)** of U is

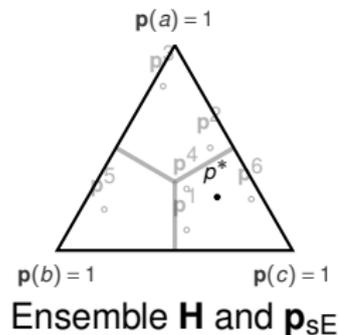
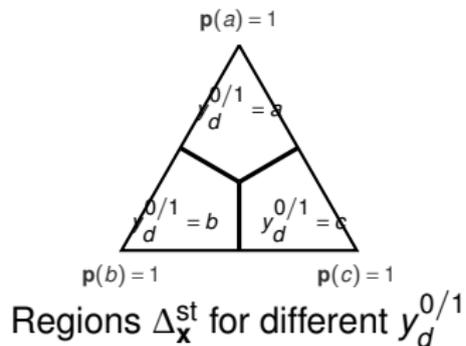
$$Y_d^U \in \operatorname{argmax}_{Y' \subset \mathcal{Y}} \mathbf{E}[U(Y', y)] = \operatorname{argmax}_{Y' \subset \mathcal{Y}} \sum_{y \in \mathcal{Y}} U(Y', y) \mathbf{p}_d(y). \quad (4)$$

Commonly used utilities, such as utility-discounted accuracies:

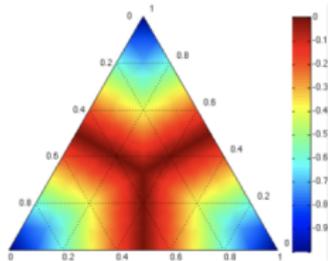
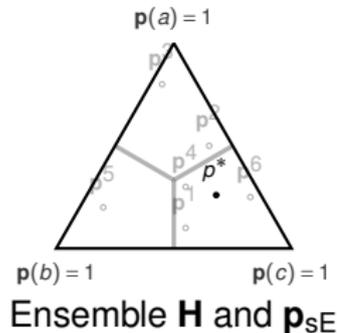
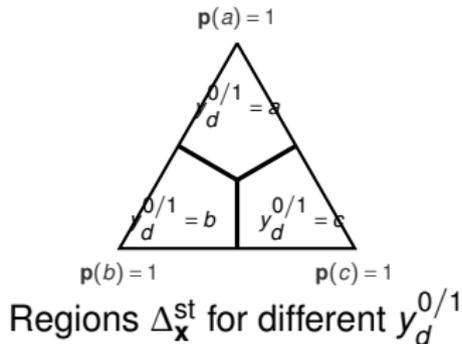
$$U(Y', y) = \frac{1}{g(|Y'|)} \llbracket y \in Y' \rrbracket, \quad (5)$$

- Find a BOP Y_d^U (11) takes $O(K \log(K))$.
- A BOP Y_d^U (11) consists of the most probable classes on \mathbf{p}_d .

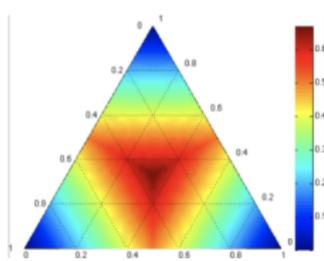
Probabilistic uncertainty scores



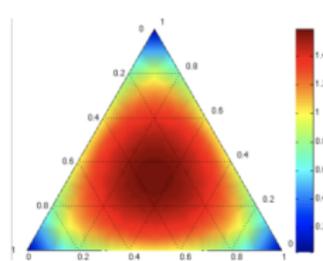
Probabilistic uncertainty scores



(a) Smallest margin (↑)



(b) Least confidence (↓)



(c) Entropy (↓)

Heatmaps illustrating the **behavior of probabilistic uncertainty scores**

Probabilistic uncertainty scores (Cont.)

Smallest margin (\uparrow) is defined as

$$S_{SM}(\mathbf{p}_d) = \mathbf{p}_d(y^{st}) - \mathbf{p}_d(y^{nd}). \quad (6)$$

Example: A classification problem with $\mathcal{Y} = \{a, b, c\}$:

	\mathbf{x}_1	\mathbf{x}_2
	50 \rightarrow (0.6, 0.4, 0.0)	100 \rightarrow (0.3, 0.4, 0.3)
	50 \rightarrow (0.0, 0.4, 0.6)	
\mathbf{p}_{sE}	(0.3, 0.4, 0.3)	
$S_{SM}(\uparrow)$	0.1	
	Should we consider \mathbf{x}_1 and \mathbf{x}_2 the same?	

Probabilistic uncertainty scores (Cont.)

Smallest margin (\uparrow) is defined as

$$S_{SM}(\mathbf{p}_d) = \mathbf{p}_d(y^{st}) - \mathbf{p}_d(y^{nd}). \quad (6)$$

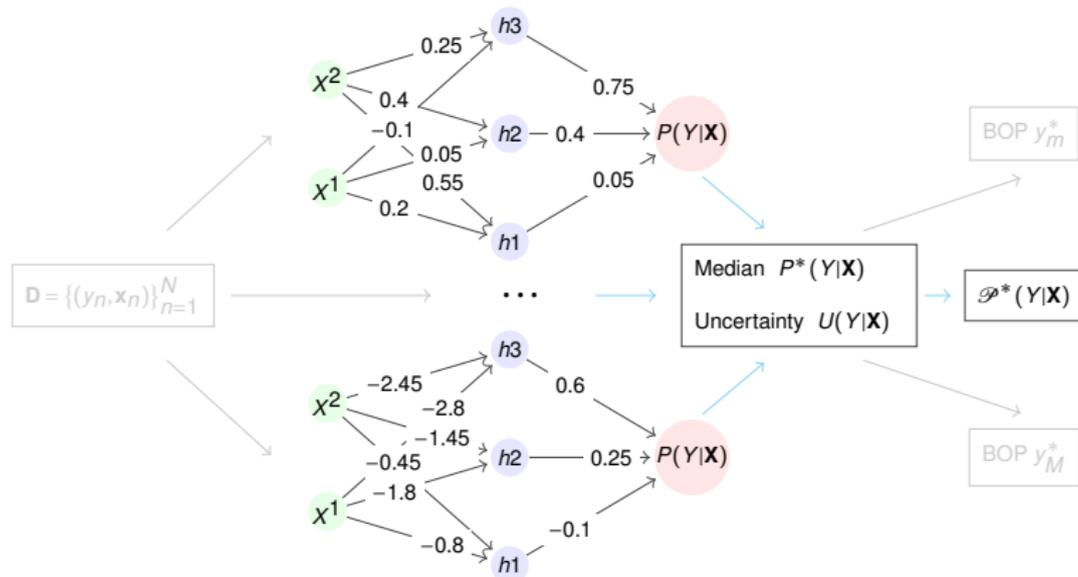
Example: A classification problem with $\mathcal{Y} = \{a, b, c\}$:

	\mathbf{x}_3	\mathbf{x}_4
	80 \rightarrow (1.0, 0.0, 0.0)	100 \rightarrow (0.8, 0.2, 0.0)
	20 \rightarrow (0.0, 1.0, 0.0)	
\mathbf{p}_{sE}	(0.8, 0.2, 0.0)	
$S_{SM}(\uparrow)$	0.6	
	Should we consider \mathbf{x}_3 and \mathbf{x}_4 the same?	

Outline

- Credal ensembling
 - A median classifier: Learning and inference
 - A credal classifier: Learning and inference
- Applications in machine learning
- A few practical aspects
- Exercices on prediction-making

A credal classifier and its predictions [2]



For any query instance, once $\mathcal{P}^*(\mathcal{Y}|\mathbf{x})$ is estimated:

- IP decision rules can be called to make set-valued predictions
- uncertainty scores defined for credal sets can be computed.

Estimate a credal classifier

Each credal classifier \mathbf{CH}_α^d is defined in a point-wise manner:

$$\mathbf{CH}_\alpha^d := \left\{ \mathbf{p} := \sum_{m=1}^{M_\alpha} \gamma_m \mathbf{p}^{(m)} \mid \gamma_m \geq 0, m \in [M_\alpha], \sum_{m=1}^{M_\alpha} \gamma_m = 1 \right\}, \quad (7)$$

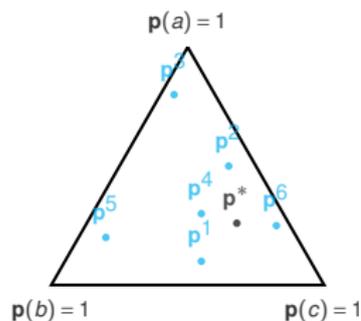
where $\mathbf{p}^{(m)}$ is the m -th closet point to \mathbf{p}_d according to the distance d .

Estimate a credal classifier

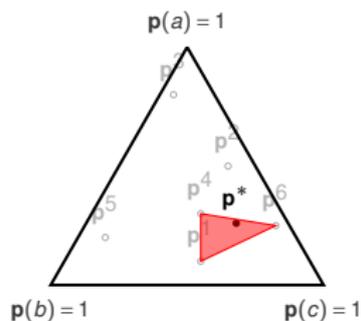
Each credal classifier \mathbf{CH}_α^d is defined in a point-wise manner:

$$\mathbf{CH}_\alpha^d := \left\{ \mathbf{p} := \sum_{m=1}^{M_\alpha} \gamma_m \mathbf{p}^{(m)} \mid \gamma_m \geq 0, m \in [M_\alpha], \sum_{m=1}^{M_\alpha} \gamma_m = 1 \right\}, \quad (7)$$

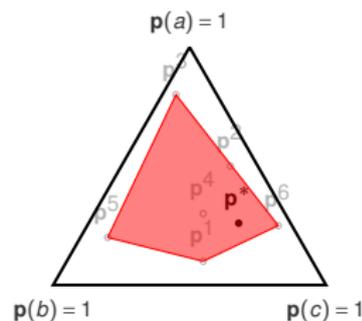
where $\mathbf{p}^{(m)}$ is the m -th closet point to \mathbf{p}_d according to the distance d .



Ensemble \mathbf{H} and \mathbf{p}_{sE}^*



Credal set $\mathbf{CH}_{0.5}^{sE}$



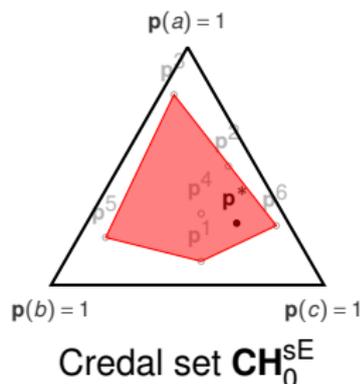
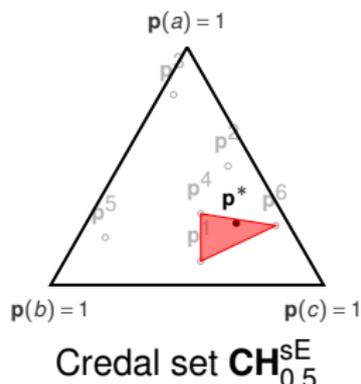
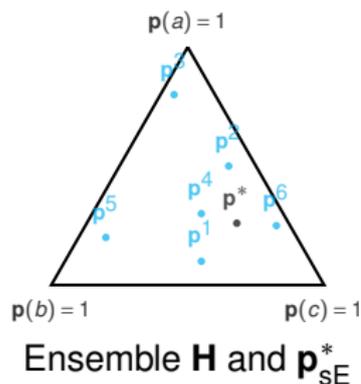
Credal set \mathbf{CH}_0^{sE}

Estimate a credal classifier

Each credal classifier \mathbf{CH}_α^d is defined in a point-wise manner:

$$\mathbf{CH}_\alpha^d := \left\{ \mathbf{p} := \sum_{m=1}^{M_\alpha} \gamma_m \mathbf{p}^{(m)} \mid \gamma_m \geq 0, m \in [M_\alpha], \sum_{m=1}^{M_\alpha} \gamma_m = 1 \right\}, \quad (7)$$

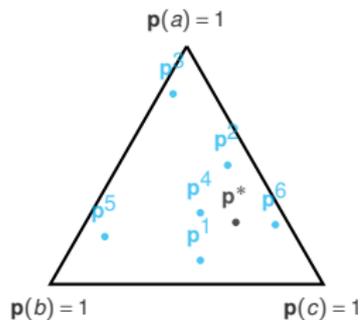
where $\mathbf{p}^{(m)}$ is the m -th closet point to \mathbf{p}_d according to the distance d .



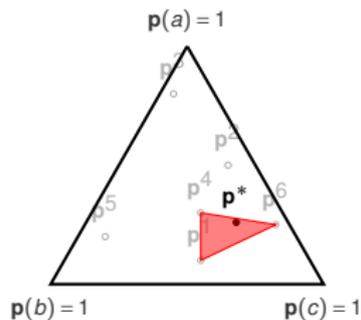
The hyperparameter $\alpha^* \leftarrow$ nested cross validation or a validation set.

Optimal set-valued predictions under IP decision rules

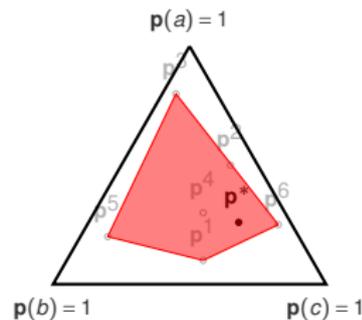
Basic set (instance-wise manner): The credal set $\mathbf{CH}_{\alpha^*}^d$ is given.



Ensemble \mathbf{H} and \mathbf{p}_{sE}



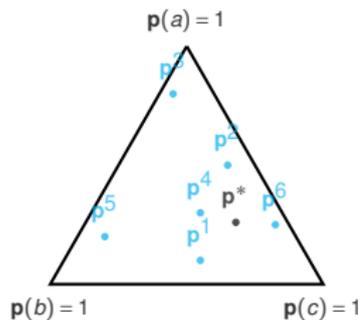
Credal set $\mathbf{CH}_{0.5}^{sE}$



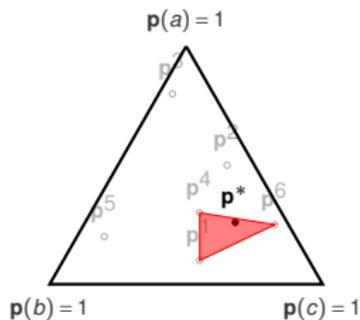
Credal set \mathbf{CH}_0^{sE}

Optimal set-valued predictions under IP decision rules

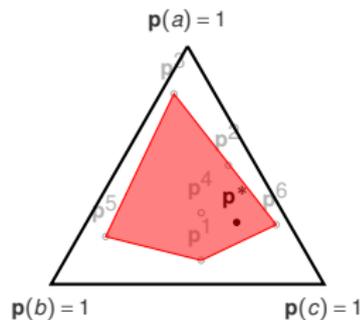
Basic set (instance-wise manner): The credal set $\mathbf{CH}_{\alpha^*}^d$ is given.



Ensemble \mathbf{H} and \mathbf{p}_{sE}



Credal set $\mathbf{CH}_{0.5}^{SE}$



Credal set \mathbf{CH}_0^{SE}

- Any IP decision rule $R_{IP} : 2^{\Delta^K} \mapsto 2^{\mathcal{Y}}$ can be applied.
- Any related algorithmic solutions can be leveraged.

Optimal set-valued predictions under IP decision rules (Cont.)

Basic set (instance-wise manner):

- The credal set $\mathbf{CH}_{\alpha^*}^d$ is given.
- A the higher the better utility $u: \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}_+$

E-admissibility under u :

- A class y is E-admissible if there exist $\mathbf{p} \in \mathbf{CH}_{\alpha^*}^d$ so that $y = y^u$.
- This can be checked by solving a linear program.

Optimal set-valued predictions under IP decision rules (Cont.)

Basic set (instance-wise manner):

- The credal set $\mathbf{CH}_{\alpha^*}^d$ is given.
- A the higher the better utility $u: \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}_+$

E-admissibility under u :

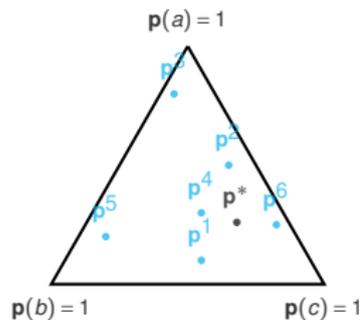
- A class y is E-admissible if there exist $\mathbf{p} \in \mathbf{CH}_{\alpha^*}^d$ so that $y = y^u$.
- This can be checked by solving a linear program.

Maximality under u :

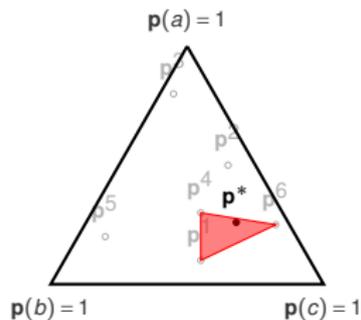
- A class y is maximal if there doesn't exist $y' \neq y$ such that y' dominates y on all $\mathbf{p} \in \mathbf{CH}_{\alpha^*}^d$ (w.r.t. u).
- This can be checked by solving $K - 1$ linear programs.
- We can also enumerate all the distributions \mathbf{p}^m , $m \in [M]$.

Credal set-based uncertainty scores

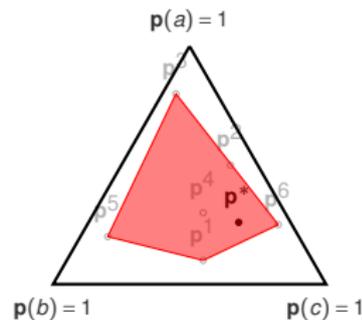
Basic set (instance-wise manner): The credal set $\mathbf{CH}_{\alpha^*}^d$ is given.



Ensemble \mathbf{H} and \mathbf{p}_{sE}



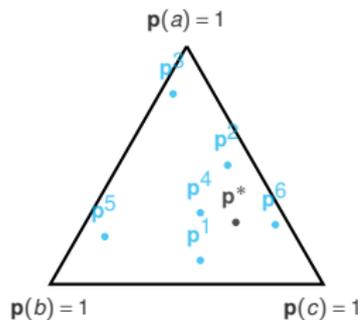
Credal set $\mathbf{CH}_{0.5}^{sE}$



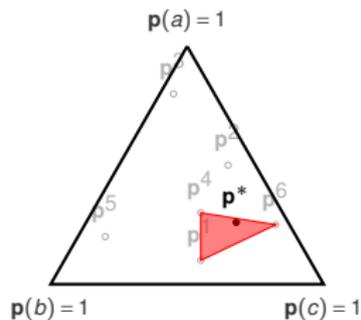
Credal set \mathbf{CH}_0^{sE}

Credal set-based uncertainty scores

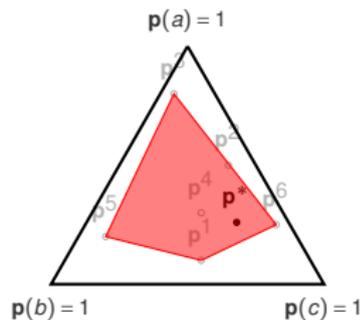
Basic set (instance-wise manner): The credal set $\mathbf{CH}_{\alpha^*}^d$ is given.



Ensemble \mathbf{H} and \mathbf{p}_{sE}



Credal set $\mathbf{CH}_{0.5}^{sE}$



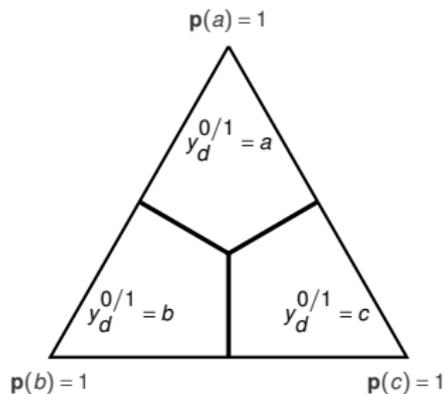
Credal set \mathbf{CH}_0^{sE}

- Any credal set-based uncertainty score can be used.
- Any related algorithmic solutions can be leveraged.

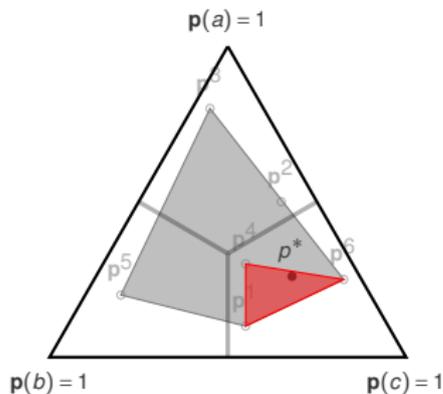
Credal set-based uncertainty scores (Cont.)

Decision-related uncertainty scores:

- How certain the ensemble \mathbf{H} is about y_d^u ?
- How consensus of the ensemble members is about y_d^u ?



Regions $\Delta_{\mathbf{x}}^{\text{st}}$ for different $y_d^{0/1}$



Credal set $\mathbf{CH}_0^{\text{SE}}(\mathbf{x})$ with p_{SE}

Decision-related uncertainty scores

Basic setting (instance-wise manner):

- The median distribution \mathbf{p}_d is given.
- The predictions $\{\mathbf{p}^m | m \in [M]\}$ are given.
- A probabilistic uncertainty score $S: \Delta^K \mapsto \mathbb{R}$ is given.

Decision-related uncertainty scores

Basic setting (instance-wise manner):

- The median distribution \mathbf{p}_d is given.
- The predictions $\{\mathbf{p}^m | m \in [M]\}$ are given.
- A probabilistic uncertainty score $S: \Delta^K \mapsto \mathbb{R}$ is given.

A **decision-related uncertainty** version of S is (defined as its empirical expectation)

$$RS(\mathbf{p}_d^u) := \frac{1}{M+1} \left(\sum_{m=1}^M \mathbb{I}[\mathbf{p}^m \in \mathbf{CH}_{y_d^u}^d] S(\mathbf{p}^m) + S(\mathbf{p}_d) \right), \quad (8)$$

where $\mathbb{I}[\mathbf{p}^m \in \mathbf{CH}_{y_d^u}^d] = 1$ implies y_d^u is a best solution on \mathbf{p}^m under u .

Decision-related uncertainty scores (Cont.)

Smallest margin (\uparrow) is defined as

$$S_{SM}(\mathbf{p}_d) = \mathbf{p}_d(y^{st}) - \mathbf{p}_d(y^{nd}). \quad (9)$$

Example: A classification problem with $\mathcal{Y} = \{a, b, c\}$:

	\mathbf{x}_1	\mathbf{x}_2
	50 \rightarrow (0.6, 0.4, 0.0)	100 \rightarrow (0.3, 0.4, 0.3)
	50 \rightarrow (0.0, 0.4, 0.6)	
\mathbf{p}_{sE}	(0.3, 0.4, 0.3)	
$S_{SM}(\uparrow)$	0.1	
$RS_{SM}(\uparrow)$	0.0	0.1

Decision-related uncertainty scores (Cont.)

Smallest margin (\uparrow) is defined as

$$S_{SM}(\mathbf{p}_d) = \mathbf{p}_d(y^{st}) - \mathbf{p}_d(y^{nd}). \quad (9)$$

Example: A classification problem with $\mathcal{Y} = \{a, b, c\}$:

	\mathbf{x}_3	\mathbf{x}_4
	$80 \rightarrow (1.0, 0.0, 0.0)$	$100 \rightarrow (0.8, 0.2, 0.0)$
	$20 \rightarrow (0.0, 1.0, 0.0)$	
\mathbf{p}_{sE}	$(0.8, 0.2, 0.0)$	
$S_{SM}(\uparrow)$	0.6	
$RS_{SM}(\uparrow)$	0.798	0.6

Should we put weights on the impact of ensemble members?

Outline

- Credal ensembling
- Applications in machine learning
 - Prediction making
 - Classification with rejection
 - Uncertainty sampling
- A few practical aspects
- Exercises on prediction-making

Outline

- Credal ensembling
- Applications in machine learning
 - Prediction making
 - Classification with rejection
 - Uncertainty sampling
- A few practical aspects
- Exercices on prediction-making

Experimental setting

Basic setting:

- Use random forests of cardinality 100 as the ensembles
- Follow a 10-cross validation protocol.
- Use hyperparameter α^* ← nested 10 fold cross validation

Assess the impact of \mathbf{p}_{SE} , \mathbf{p}_{L1} and \mathbf{p}_{KL} on

- the clean version of the data sets
- noisy version (randomly flip the class of 25% of training instances)

Experimental setting

Basic setting:

- Use random forests of cardinality 100 as the ensembles
- Follow a 10-cross validation protocol.
- Use hyperparameter α^* ← nested 10 fold cross validation

Assess the impact of \mathbf{p}_{SE} , \mathbf{p}_{L1} and \mathbf{p}_{KL} on

- the clean version of the data sets
- noisy version (randomly flip the class of 25% of training instances)

Once credal set $\mathbf{CH}_{\alpha^*}^d$ is computed, it is used to

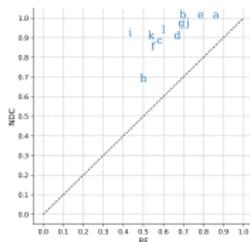
- find the set-valued prediction under the E-admissibility rule.

Results on clean data sets: U_{65} scores (in %) [3]

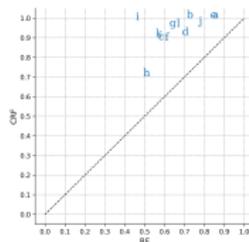
Data set: (N,P,K)	NDC	SQE-E	L1-E	KL-E	CRF	CH_0
eco.: (336,7,8)	85.51	86.07	85.81	87.07	84.46	43.60
der.: (358,34,6)	97.18	97.05	97.22	98.59	96.19	51.74
lib.: (360, 90, 15)	76.58	73.35	75.24	79.41	73.45	14.60
vow.: (990, 10, 11)	86.63	86.35	87.65	92.35	82.68	17.75
win.: (1599, 11, 6)	68.66	68.32	68.39	68.63	67.35	36.53
seg.: (2300, 19, 7)	97.17	97.12	96.99	97.64	96.73	71.00

Results on clean data sets: U_{65} scores (in %) [3]

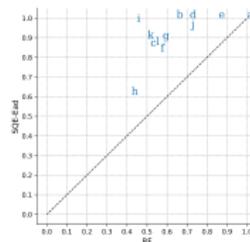
Data set: (N,P,K)	NDC	SQE-E	L1-E	KL-E	CRF	CH_0
eco.: (336,7,8)	85.51	86.07	85.81	87.07	84.46	43.60
der.: (358,34,6)	97.18	97.05	97.22	98.59	96.19	51.74
lib.: (360, 90, 15)	76.58	73.35	75.24	79.41	73.45	14.60
vow.: (990, 10, 11)	86.63	86.35	87.65	92.35	82.68	17.75
win.: (1599, 11, 6)	68.66	68.32	68.39	68.63	67.35	36.53
seg.: (2300, 19, 7)	97.17	97.12	96.99	97.64	96.73	71.00



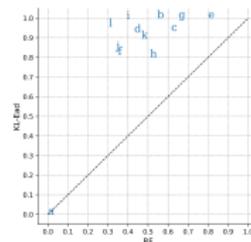
(a) NDC vs RF



(b) CRF vs RF



(e) SQE-Ead vs RF



(f) KL-Ead vs RF

Correctness of cautious predictors (vertical) vs accuracy of RF (horizontal)

Outline

- Credal ensembling
- Applications in machine learning
 - Prediction making
 - **Classification with rejection**
 - Uncertainty sampling
- A few practical aspects
- Exercises on prediction-making

Experimental setting [3]

Basic setting:

- Randomly flip the class of 25% of training instances
- Using random forests of cardinality 100 as the ensembles
- The median \mathbf{p}_{SE} is employed
- Assess smallest margin S_{SM} (\uparrow) and \mathbf{RS}_{SM} (\uparrow)

Experimental setting [3]

Basic setting:

- Randomly flip the class of 25% of training instances
- Using random forests of cardinality 100 as the ensembles
- The median p_{SE} is employed
- Assess smallest margin S_{SM} (\uparrow) and RS_{SM} (\uparrow)

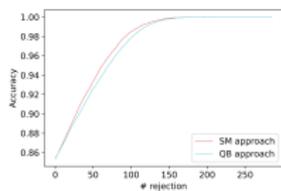
Budget based rejection protocol requires

- a sufficiently large number of test instances,
- a predefined number (or proportion) of rejections.

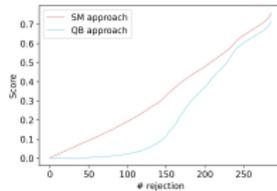
Threshold-based rejection protocol

- requires a predefined threshold on uncertainty score (\uparrow),
- rejects instances whose scores are lower than the threshold.

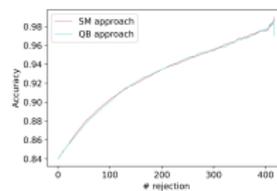
Results on noisy data sets [3]



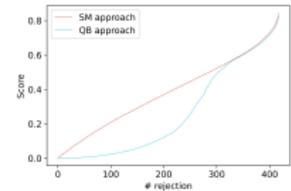
(a) derma. + SM



(b) derma. + SM

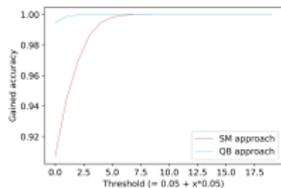


(c) forest + SM

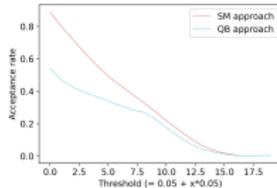


(d) forest + SM

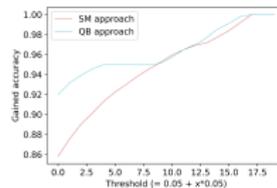
Test accuracy and chosen score as the functions of the number of rejections
 20×5 cross-validation with (train, test) = (20%, 80%)



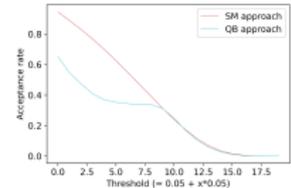
(a) derma. + SM



(b) derma. + SM



(c) forest + SM



(d) forest + SM

Test accuracy and acceptance rate as the functions of the threshold
 20×5 cross-validation with (train, test) = (20%, 80%)

Outline

- Credal ensembling
- Applications in machine learning
 - Prediction making
 - Classification with rejection
 - **Uncertainty sampling**
- A few practical aspects
- Exercices on prediction-making

Experimental setting

Basic setting:

- Randomly flip the class of 25% of training + pool instances
- Using random forests of cardinality 100 as the ensembles
- The median \mathbf{p}_{SE} is employed
- Assess smallest margin S_{SM} (\uparrow) and \mathbf{RS}_{SM} (\uparrow)

Experimental setting

Basic setting:

- Randomly flip the class of 25% of training + pool instances
- Using random forests of cardinality 100 as the ensembles
- The median p_{SE} is employed
- Assess smallest margin S_{SM} (\uparrow) and RS_{SM} (\uparrow)

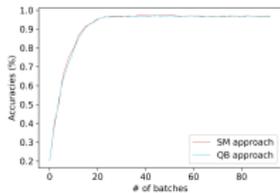
Budget based sampling protocol

- requires a predefined number (or proportion) of queries,
- stops when the predefined number is reached.

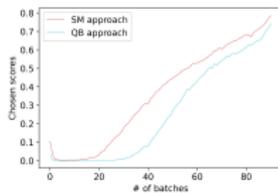
Threshold-based sampling protocol

- requires a predefined threshold on uncertainty score (\uparrow),
- stops when the predefined threshold is reached.

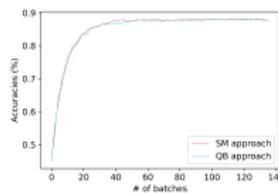
Results on noisy data sets [3]



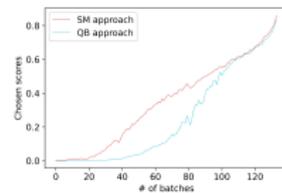
(a) derma. + SM



(b) derma. + SM

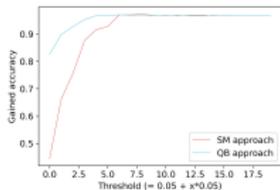


(c) forest + SM

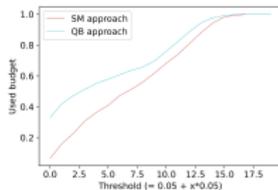


(d) forest + SM

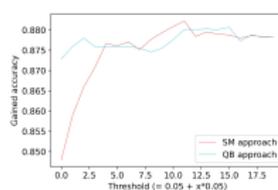
Test accuracy and chosen score as the functions of the number of queries
 10×5 cross-validation with (train, pool, test) = (3%, 77%, 20%)



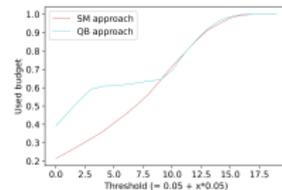
(a) derma. + SM



(b) derma. + SM



(c) forest + SM



(d) forest + SM

Test accuracy and used budget as the functions of the threshold
 10×5 cross-validation with (train, pool, test) = (3%, 77%, 20%)

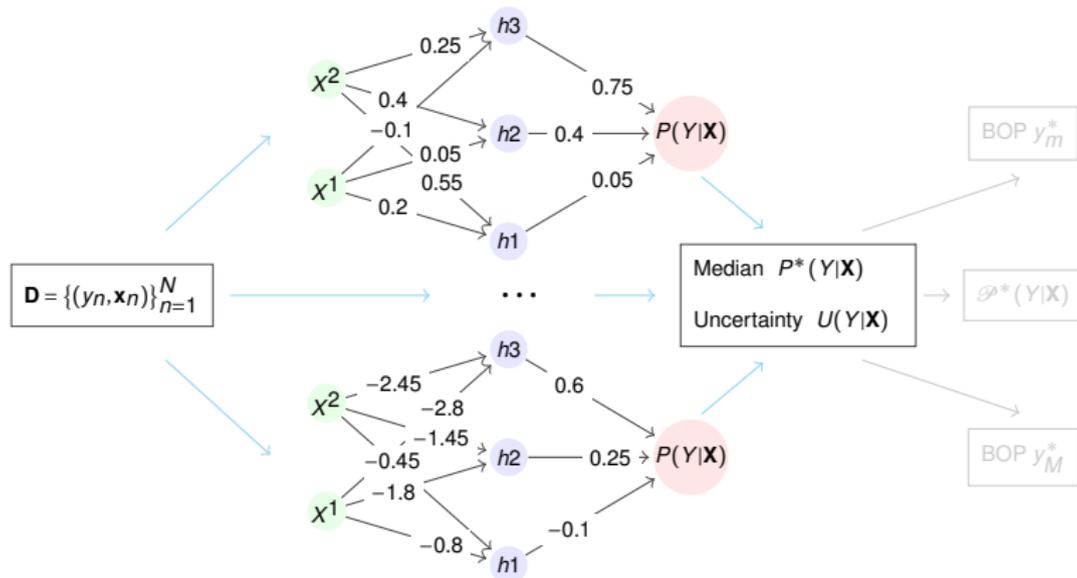
Outline

- Credal ensembling
- Applications in machine learning
- **A few practical aspects**
 - Compact deep ensembles
- Exercises on prediction-making

Outline

- Credal ensembling
- Applications in machine learning
- A few practical aspects
 - Compact deep ensembles
- Exercices on prediction-making

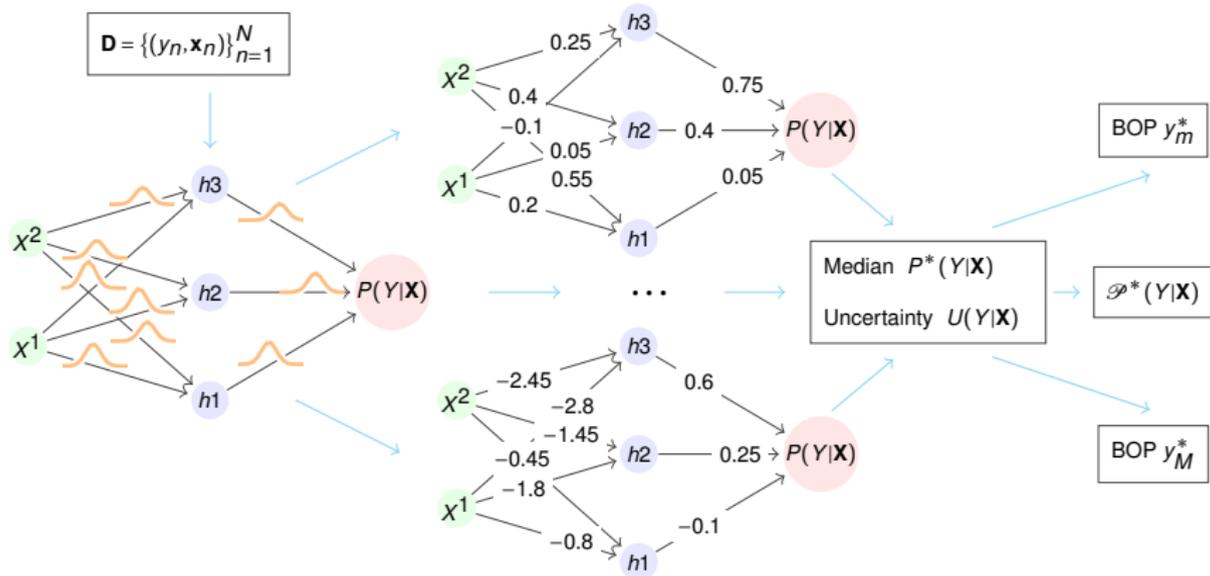
Conventional deep ensembles



Compared to the use of a single network:

- Much longer training time + Much larger storage memory
- Longer inference time

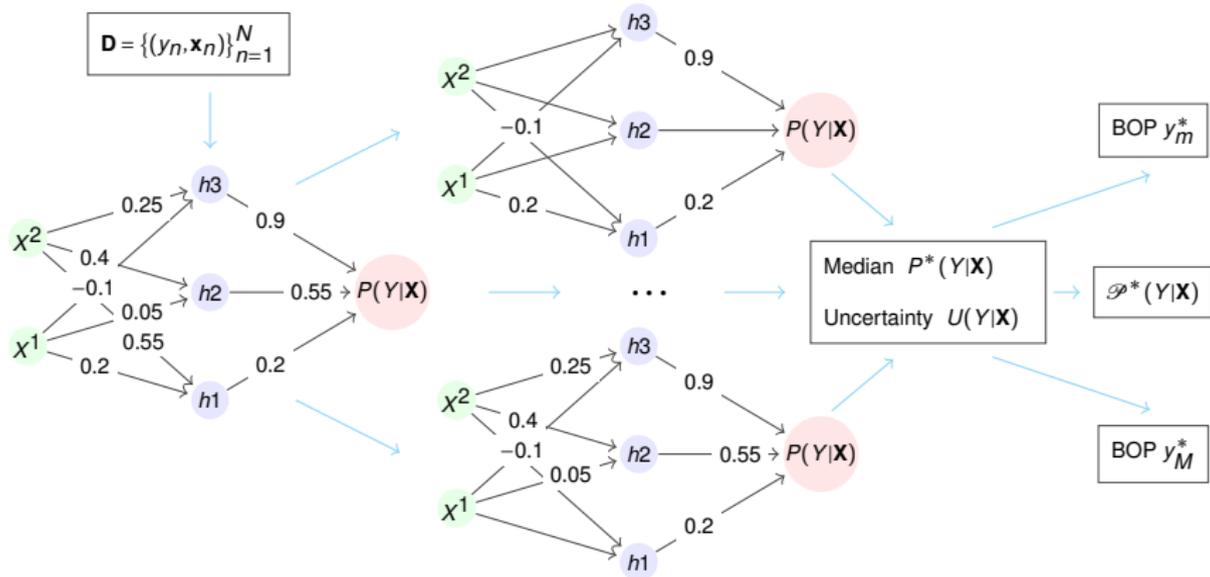
A BNN as an ensemble [5]



Compared to the use of a single network:

- A bit longer training time + A bit larger storage memory
- Longer inference time

A CNN with dropout predictions as an ensemble [4]



Compared to the use of a single network:

- Similar training time + Similar storage memory
- Longer inference time

Experimental setting

Basic setting:

- Use BNNs with 100 Monte Carlo runs as the ensembles
- Use the clean version of the data sets

Assess the impact of p_{SE} , p_{L1} and p_{KL} on

	Image	train/test	# classes
CIFAR-10	32x32 color	50,000/10,000	10
Fashion-MNIST	grayscale	60,000/10,000	10

Experimental setting

Basic setting:

- Use BNNs with 100 Monte Carlo runs as the ensembles
- Use the clean version of the data sets

Assess the impact of \mathbf{p}_{SE} , \mathbf{p}_{L1} and \mathbf{p}_{KL} on

	Image	train/test	# classes
CIFAR-10	32x32 color	50,000/10,000	10
Fashion-MNIST	grayscale	60,000/10,000	10

Once \mathbf{p}_d is computed, it is used to

- find precise prediction optimizing the $u_{0,1}$,
- find set-valued predictions optimizing the u_{65} and u_{80} [1].

Average $u_{0,1}$, u_{65} and u_{80} on the test set

Results [5]	CIFAR-10			Fashion MNIST		
	sE	L1	KL	sE	L1	KL
$u_{0/1}$ (\uparrow)	90.04	90.10	90.14	93.07	93.11	93.08
opt_u65_eva_u65 (\uparrow)	90.47	90.51	90.46	93.38	93.31	93.26
opt_u80_eva_u80 (\uparrow)	91.77	91.76	91.76	94.41	94.39	94.27
u65_set_size (\downarrow)	2.03	2.02	2.03	2.02	2.02	2.02
u80_set_size (\downarrow)	2.04	2.02	2.03	2.02	2.02	2.02

A closer look at $u_{0,1}$ and u_{65}

Results [5]	CIFAR-10			Fashion MNIST		
	sE	L1	KL	sE	L1	KL
c_pr_u65_c_si (\uparrow)	94.91	95.91	97.53	97.53	97.19	98.43
c_pr_u65_c_se (\downarrow)	5.08	4.08	2.46	2.46	2.80	1.56
w_pr_u65_c_se (\uparrow)	32.12	26.86	17.64	24.96	25.39	15.75
w_pr_u65_w_se (\downarrow)	15.26	11.81	7.50	5.05	5.95	4.62
w_pr_u65_w_si (\downarrow)	52.61	61.31	74.84	69.98	68.65	79.62

A closer look at $u_{0,1}$ and u_{80}

Results [5]	CIFAR-10			Fashion MNIST		
	sE	L1	KL	sE	L1	KL
c_pr_u80_c_si (\uparrow)	86.89	93.22	94.28	94.34	94.03	95.47
c_pr_u80_c_se (\downarrow)	13.10	6.77	5.71	5.65	5.96	4.52
w_pr_u80_c_se (\uparrow)	53.21	37.07	34.38	43.86	44.26	37.28
w_pr_u80_w_se (\downarrow)	23.89	19.19	16.32	10.82	10.44	8.67
w_pr_u80_w_si (\downarrow)	22.89	43.73	49.29	45.31	45.28	54.04

Results [5]	CIFAR-10			Fashion MNIST		
	sE	L1	KL	sE	L1	KL
$u_{0/1}$ (\uparrow)	90.04	90.10	90.14	93.07	93.11	93.08
u65_set_size (\downarrow)	2.03	2.02	2.03	2.02	2.02	2.02
u80_set_size (\downarrow)	2.04	2.02	2.03	2.02	2.02	2.02
c_pr_u65_c_si (\uparrow)	94.91	95.91	97.53	97.53	97.19	98.43
c_pr_u65_c_se (\downarrow)	5.08	4.08	2.46	2.46	2.80	1.56
w_pr_u65_c_se (\uparrow)	32.12	26.86	17.64	24.96	25.39	15.75
w_pr_u65_w_se (\downarrow)	15.26	11.81	7.50	5.05	5.95	4.62
w_pr_u65_w_si (\downarrow)	52.61	61.31	74.84	69.98	68.65	79.62
c_pr_u80_c_si (\uparrow)	86.89	93.22	94.28	94.34	94.03	95.47
c_pr_u80_c_se (\downarrow)	13.10	6.77	5.71	5.65	5.96	4.52
w_pr_u80_c_se (\uparrow)	53.21	37.07	34.38	43.86	44.26	37.28
w_pr_u80_w_se (\downarrow)	23.89	19.19	16.32	10.82	10.44	8.67
w_pr_u80_w_si (\downarrow)	22.89	43.73	49.29	45.31	45.28	54.04

Outline

- Credal ensembling
- Applications in machine learning
- A few practical aspects
- **Exercices on prediction-making**
 - Probabilistic classifiers
 - Credal classifiers

Outline

- Credal ensembling
- Applications in machine learning
- A few practical aspects
- Exercices on prediction-making
 - Probabilistic classifiers
 - Credal classifiers

Bayesian-optimal predictions

A **Bayesian-optimal prediction (BOP)** of u is

$$y_d^u \in \operatorname{argmax}_{y' \in \mathcal{Y}} \mathbf{E}[u(y', y)] = \operatorname{argmax}_{y' \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} u(y', y) \mathbf{p}_d(y). \quad (10)$$

Bayesian-optimal predictions

A **Bayesian-optimal prediction (BOP)** of u is

$$y_d^u \in \operatorname{argmax}_{y' \in \mathcal{Y}} \mathbf{E}[u(y', y)] = \operatorname{argmax}_{y' \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} u(y', y) \mathbf{p}_d(y). \quad (10)$$

Question: Prove that a BOP (10) takes from $O(K)$ to $O(K^2)$.

Bayesian-optimal set-valued predictions

A **Bayesian-optimal prediction (BOP)** of U is

$$Y_d^U \in \operatorname{argmax}_{Y' \subset \mathcal{Y}} \mathbf{E}[U(Y', y)] = \operatorname{argmax}_{Y' \subset \mathcal{Y}} \sum_{y \in \mathcal{Y}} U(Y', y) \mathbf{p}_d(y). \quad (11)$$

Bayesian-optimal set-valued predictions

A **Bayesian-optimal prediction (BOP)** of U is

$$Y_d^U \in \operatorname{argmax}_{Y' \subset \mathcal{Y}} \mathbf{E}[U(Y', y)] = \operatorname{argmax}_{Y' \subset \mathcal{Y}} \sum_{y \in \mathcal{Y}} U(Y', y) \mathbf{p}_d(y). \quad (11)$$

Question: Prove that for any utility-discounted accuracy

$$U(Y', y) = \frac{1}{g(|Y'|)} \mathbb{I}[y \in Y'], \quad (12)$$

finding a BOP Y_d^U (11) takes $O(K \log(K))$.

Bayesian-optimal set-valued predictions

A **Bayesian-optimal prediction (BOP)** of U is

$$Y_d^U \in \operatorname{argmax}_{Y' \subset \mathcal{Y}} \mathbf{E}[U(Y', y)] = \operatorname{argmax}_{Y' \subset \mathcal{Y}} \sum_{y \in \mathcal{Y}} U(Y', y) \mathbf{p}_d(y). \quad (11)$$

Question: Prove that for any utility-discounted accuracy

$$U(Y', y) = \frac{1}{g(|Y'|)} \mathbb{I}[y \in Y'], \quad (12)$$

finding a BOP Y_d^U (11) takes $O(K \log(K))$.

Hint: First show that a BOP Y_d^U (11) consists of the most probable classes on \mathbf{p}_d .

Outline

- Credal ensembling
- Applications in machine learning
- A few practical aspects
- Exercices on prediction-making
 - Probabilistic classifiers
 - Credal classifiers

E-admissible and maximal sets (Recap)

- **E-admissibility** under u : A class y is E-admissible if there exist $\mathbf{p} \in \mathbf{CH}_{\alpha^*}^d$ so that $y = y^u$.
- **Maximality** under u : A class y is maximal if there doesn't exist $y' \neq y$ such that y' dominates y on all $\mathbf{p} \in \mathbf{CH}_{\alpha^*}^d$ (w.r.t. u).

Check if a class is E -admissible

Question: Prove that checking whether a given class y is E -admissible can be done by solving a linear program.

Check if a class is maximal

Question: Prove that checking whether a given class y is maximal can be done by solving $K - 1$ linear program.

Properties of E-admissible and maximal sets

- **Question 1:** Prove that the E-admissible set is a subset of the maximal set.
- **Question 2:** Show that the E-admissible set can be a strict subset of the maximal set.
- **Question 3:** Show that the two sets can be identical.
- **Question 4:** Show that the cardinality of the E-admissible set can be larger than the number of extreme points on the credal set.

Properties of E-admissible and maximal sets: Hints

- **Question 1:** Prove that the E-admissible set is a subset of the maximal set. → We did it during the last lecture.
- **Question 2:** Show that the E-admissible set can be a strict subset of the maximal set. → Consider the credal set with two extreme points $\{(0.35, 0.4, 0.25), (0.3, 0.2, 0.5)\}$.
- **Question 3:** Show that the two sets can be identical. → Consider the credal set with two extreme points $\{(0.3, 0.5, 0.2), (0.2, 0.7, 0.1)\}$
- **Question 4:** Show that the cardinality of the E-admissible set can be larger than the number of extreme points on the credal set. → Consider the credal set with two extreme points $\{(0.6, 0.4, 0.0), (0.0, 0.4, 0.6)\}$

References I

- [1] T. Mortier, M. Wydmuch, K. Dembczyński, E. Hüllermeier, and W. Waegeman.
Efficient set-valued prediction in multi-class classification.
Data Mining and Knowledge Discovery, 35(4):1435–1469, 2021.
- [2] V.-L. Nguyen, H. Zhang, and S. Destercke.
Learning sets of probabilities through ensemble methods.
In *Proceedings of the 17th European Conference on Symbolic and Quantitative Approaches to Reasoning with Uncertainty (ECSQARU)*, pages 270–283, 2023.
- [3] V.-L. Nguyen, H. Zhang, and S. Destercke.
Credal ensembling in multi-class classification.
Machine Learning, pages 1–64, 2024.
- [4] K.-D. Tran, X.-T. Hoang, D.-M. Nguyen, V.-L. Nguyen, S. Destercke, and V.-N. Huynh.
Compact probabilistic ensemble learning in multi-class classification.
Under preparation, pages 1–20, 2024.
- [5] K.-D. Tran, X.-T. Hoang, V.-L. Nguyen, S. Destercke, and V.-N. Huynh.
Robust classification in bayesian neural networks.
In *Submitted to the Eleventh International Symposium on Integrated Uncertainty in Knowledge Modelling and Decision Making (IUKM)*, pages 1–12, 2024.