

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.



Uncertainty reasoning and machine learning

Some first credal classifiers

Vu-Linh Nguyen

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

AOS4 master courses

Objectives

After this lecture students should be able to

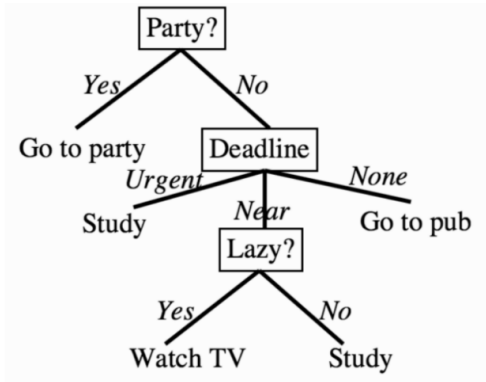
- use IDM and related models in Naïve credal classifier (NCC) [3]
- use IDM and related models in decision trees [8]

Outline

- Graphical Interpretation of Probabilistic Models
- Naïve Bayesian/Credal classifiers
- Decision Trees
- Bayesian Neural Networks

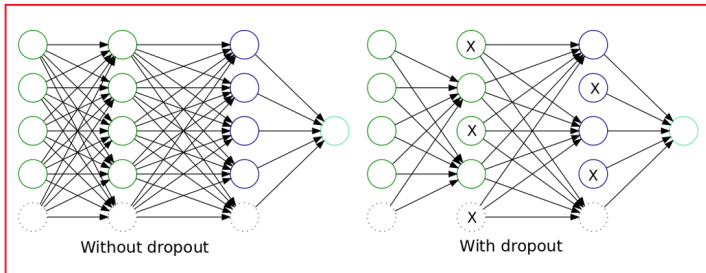
How to interpret a decision tree?

How to interpret a decision tree?



How to interpret a (feedforward) neural network?

How to interpret a (feedforward) neural network?



Probabilistic Models: Graphical Interpretation [5, 9]

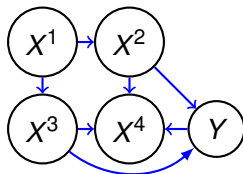
Basic setup

- A set of features $\mathbf{X} = \{X^1, \dots, X^M\}$; $[M] := \{1, \dots, M\}$
- A class variable Y whose outcome $y \in \mathcal{Y}$

Probabilistic Models: Graphical Interpretation [5, 9]

Basic setup

- A set of features $\mathbf{X} = \{X^1, \dots, X^M\}$; $[M] := \{1, \dots, M\}$
- A class variable Y whose outcome $y \in \mathcal{Y}$
- A directed acyclic graph (DAG) connecting Y and X^m



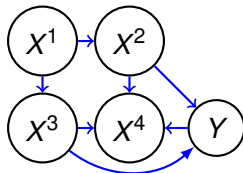
This DAG (model structure) tells us:

- $\text{pa}(Y) = \{X^2, X^3\}$, $\text{pa}(X^1) = \emptyset$
- $\text{pa}(X^2) = \{X^1\}$, $\text{pa}(X^3) = \{X^1\}$
- $\text{pa}(X^4) = \{Y, X^2, X^3\}$

Probabilistic Models: Graphical Interpretation [5, 9]

Basic setup

- A set of features $\mathbf{X} = \{X^1, \dots, X^M\}$; $[M] := \{1, \dots, M\}$
- A class variable Y whose outcome $y \in \mathcal{Y}$
- A directed acyclic graph (DAG) connecting Y and X^m



This DAG (model structure) tells us:

- $\text{pa}(Y) = \{X^2, X^3\}$, $\text{pa}(X^1) = \emptyset$
- $\text{pa}(X^2) = \{X^1\}$, $\text{pa}(X^3) = \{X^1\}$
- $\text{pa}(X^4) = \{Y, X^2, X^3\}$

Probabilistic Models:

- Expressing $P(Y, \mathbf{X})$ using the **chain rule** (probability):

$$P(Y, \mathbf{X}) = P(Y | \text{pa}(Y)) \prod_{m=1}^M P(X^m | \text{pa}(X^m)).$$

Probabilistic Models: Model Families [9]

Probabilistic Models:

- Estimate $P(Y, \mathbf{X})$
- Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y | \text{pa}(Y)) \prod_{m=1}^M P(X^m | \text{pa}(X^m)).$$

Extreme Cases:

- Discriminative models: $Y \notin \text{pa}(X^m)$, $m \in [M]$
- Generative models: $\text{pa}(Y) = \emptyset$ and $Y \in \text{pa}(X^p)$, $m \in [M]$.

Probabilistic Models: Model Families [9]

Probabilistic Models:

- Estimate $P(Y, \mathbf{X})$
- Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y|\text{pa}(Y)) \prod_{m=1}^M P(X^m|\text{pa}(X^m)).$$

Extreme Cases:

- Discriminative models: $Y \notin \text{pa}(X^m)$, $m \in [M]$
- Generative models: $\text{pa}(Y) = \emptyset$ and $Y \in \text{pa}(X^p)$, $m \in [M]$.

Model Families:

- How to encode/parametrize $P(Y|\text{pa}(Y))$ and $P(X^m|\text{pa}(X^m))$.
- How to estimate $P(Y, \mathbf{X})$ from training data.

Credal (Imprecise Probability) Models

Basic setup

- A set of features $\mathbf{X} = \{X^1, \dots, X^M\}$
- A class variable Y whose outcome $y \in \mathcal{Y}$

Credal Models:

- $\mathcal{P} := \{P(Y, \mathbf{X}) \mid P \text{ is compatible with knowledge/data}\}$
- Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y \mid \text{pa}(Y)) \prod_{m=1}^M P(X^m \mid \text{pa}(X^m)).$$

Extreme Cases:

- Discriminative models: $Y \notin \text{pa}(X^m)$, $m \in [M] := \{1, \dots, M\}$
- Generative models: $\text{pa}(Y) = \emptyset$ and $Y \in \text{pa}(X^m)$, $m \in [M]$.

Model Families:

- How to encode/parametrize $P(Y \mid \text{pa}(Y))$ and $P(X^m \mid \text{pa}(X^m))$.
- How to estimate $\mathcal{P}(Y, \mathbf{X})$ from training data.

Assumptions and Questions

Assumption and desirable property:

A1. X^m , $m \in [M] := \{1, \dots, M\}$, are always made available

P1. Best estimates of $P(Y|\text{pa}(Y))$ and $P(X^m|\text{pa}(X^m))$ can be found given (training) data.

Assumptions and Questions

Assumption and desirable property:

A1. X^m , $m \in [M] := \{1, \dots, M\}$, are always made available

P1. Best estimates of $P(Y|pa(Y))$ and $P(X^m|pa(X^m))$ can be found given (training) data.

Questions (Exercise):

- Does the P1 hold for Naïve Bayes Classifier?
- Does the P1 hold for Decision trees?

Assumptions and Questions

Assumption and desirable property:

A1. X^m , $m \in [M] := \{1, \dots, M\}$, are always made available

P1. Best estimates of $P(Y|pa(Y))$ and $P(X^m|pa(X^m))$ can be found given (training) data.

Questions (Exercise):

- Does the P1 hold for Naïve Bayes Classifier?
- Does the P1 hold for Decision trees?

Questions (which will not be discussed in this lecture):

- What may happen if X^m , $m \in [M]$, can be partially given?
- What may happen if best estimates of $P(Y|pa(Y))$ and $P(X^m|pa(X^m))$ may not be found?

The Next Slides

We shall elaborate on how to solve classification task using

- Naïve Bayesian classifier (NBC) (an example of generative model)
- Decision trees (DTs) (examples of discriminative model)

The Next Slides

We shall elaborate on how to solve classification task using

- Naïve Bayesian classifier (NBC) (an example of generative model)
- Decision trees (DTs) (examples of discriminative model)

How IDM (Lecture 3) can be used to generalize NBC and DTs to

- cope with the case of small and partial/missing data
- make set-valued predictions under the presence of uncertainty

The Next Slides

We shall elaborate on how to solve classification task using

- Naïve Bayesian classifier (NBC) (an example of generative model)
- Decision trees (DTs) (examples of discriminative model)

How IDM (Lecture 3) can be used to generalize NBC and DTs to

- cope with the case of small and partial/missing data
- make set-valued predictions under the presence of uncertainty

We would also discuss (if we have time) the cases of

- Ensembles (Trees, Neural Nets, etc.)
- Bayesian Neural Nets

Outline

- Graphical Interpretation of Probabilistic Models
- Naïve Bayesian/Credal classifiers
 - Naïve Bayesian classifier
 - Naïve Credal classifiers
- Decision Trees
- Bayesian Neural Networks

Outline

- Graphical Interpretation of Probabilistic Models
- Naïve Bayesian/Credal classifiers
 - Naïve Bayesian classifier
 - Naïve Credal classifiers
- Decision Trees
- Bayesian Neural Networks

Generative Models

Probabilistic Models:

- Estimate $P(Y, \mathbf{X})$
- Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y|\text{pa}(Y)) \prod_{m=1}^M P(X^m|\text{pa}(X^m)).$$

Extreme Cases:

- Discriminative models: $Y \notin \text{pa}(X^m)$, $m \in [M]$
- Generative models: $\text{pa}(Y) = \emptyset$ and $Y \in \text{pa}(X^m)$, $m \in [M]$.

Model Families:

- How to encode/parametrize $P(Y|\text{pa}(Y))$ and $P(X^m|\text{pa}(X^m))$.
- How to estimate $P(Y, \mathbf{X})$ from training data.

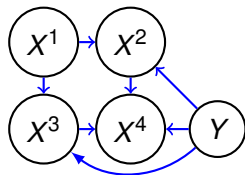
Generative Models: Structure (Exercise 1)

Let's start with an example where one wishes to model

$$P(Y, \mathbf{X}) = P(Y, X^1, X^2, X^3, X^4).$$

Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y|\text{pa}(Y)) \prod_{m=1}^M P(X^m|\text{pa}(X^m)).$$



- $\text{pa}(Y) = \emptyset$, $\text{pa}(X^1) = \emptyset$
- $\text{pa}(X^2) = \{X^1\}$
- $\text{pa}(X^3) = \{X^1\}$
- $\text{pa}(X^4) = \{X^2, X^3, Y\}$

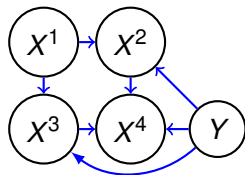
Generative Models: Structure (Exercise 1)

Let's start with an example where one wishes to model

$$P(Y, \mathbf{X}) = P(Y, X^1, X^2, X^3, X^4).$$

Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y | \text{pa}(Y)) \prod_{m=1}^M P(X^m | \text{pa}(X^m)).$$



- $\text{pa}(Y) = \emptyset$, $\text{pa}(X^1) = \emptyset$
- $\text{pa}(X^2) = ???$
- $\text{pa}(X^3) = ???$
- $\text{pa}(X^4) = ???$

Chain rule gives us

$$P(Y, \mathbf{X}) = ???.$$

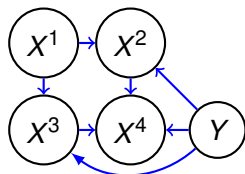
Generative Models: Structure (Solution 1)

Let's start with an example where one wishes to model

$$P(Y, \mathbf{X}) = P(Y, X^1, X^2, X^3, X^4).$$

Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y|\text{pa}(Y)) \prod_{m=1}^M P(X^m|\text{pa}(X^m)).$$



- $\text{pa}(Y) = \emptyset$, $\text{pa}(X^1) = \emptyset$
- $\text{pa}(X^2) = \{Y, X^1\}$
- $\text{pa}(X^3) = \{Y, X^1\}$
- $\text{pa}(X^4) = \{Y, X^2, X^3\}$

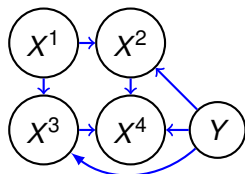
Generative Models: Structure (Solution 1)

Let's start with an example where one wishes to model

$$P(Y, \mathbf{X}) = P(Y, X^1, X^2, X^3, X^4).$$

Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y|\text{pa}(Y)) \prod_{m=1}^M P(X^m|\text{pa}(X^m)).$$



- $\text{pa}(Y) = \emptyset$, $\text{pa}(X^1) = \emptyset$
- $\text{pa}(X^2) = \{Y, X^1\}$
- $\text{pa}(X^3) = \{Y, X^1\}$
- $\text{pa}(X^4) = \{Y, X^2, X^3\}$

Chain rule gives us

$$P(Y, \mathbf{X}) = P(Y)P(X^1)P(X^2|Y, X^1)P(X^3|Y, X^1)P(X^4|Y, X^2, X^3).$$

Naïve Bayesian classifier (NBC)

Comments:

- NBC is a generative model with no arc $X' \rightarrow X$
- Chain rule gives us

$$P(Y, \mathbf{X}) = P(Y) \prod_{m=1}^M P(X^m | Y).$$

Naïve Bayesian classifier (NBC)

Comments:

- NBC is a generative model with no arc $X' \rightarrow X$
- Chain rule gives us

$$P(Y, \mathbf{X}) = P(Y) \prod_{m=1}^M P(X^m | Y).$$

To solve the **classification task**,

- joint probability distribution $P(Y, \mathbf{X})$ is learn from training data **D**
- conditional distribution $P(Y | \mathbf{X})$ is extracted using **Bayes' theorem**

$$P(y | \mathbf{x}) = \frac{P(y, \mathbf{x})}{\sum_{y' \in \mathcal{Y}} P(y', \mathbf{x})} = \frac{P(y) \prod_{m=1}^M P(x^m | y)}{\sum_{y' \in \mathcal{Y}} P(y') \prod_{m=1}^M P(x^m | y')}. \quad (1)$$

Estimate Parameters of NBC

Basic setup:

- A class variable Y with K possible values: $\mathcal{Y} = \{y^1, \dots, y^K\}$
- M discrete features: $\mathbf{X} = (X^1, \dots, X^M)$
- Feature X^m has Q_m possible values: $\mathcal{X}^m = \{x^{m,1}, \dots, x^{m,Q_m}\}$

Estimate Parameters of NBC

Basic setup:

- A class variable Y with K possible values: $\mathcal{Y} = \{y^1, \dots, y^K\}$
- M discrete features: $\mathbf{X} = (X^1, \dots, X^M)$
- Feature X^m has Q_m possible values: $\mathcal{X}^m = \{x^{m,1}, \dots, x^{m,Q_m}\}$

Task: Finding the best estimate of

- $\theta_k := P(y^k), k \in [K]$
- $\theta_k^{m,q_m} := P(x^{q_m,m} | y^k), q_m \in [Q_m], k \in [K], m \in [M]$

Estimate Parameters of NBC

Basic setup:

- A class variable Y with K possible values: $\mathcal{Y} = \{y^1, \dots, y^K\}$
- M discrete features: $\mathbf{X} = (X^1, \dots, X^M)$
- Feature X^m has Q_m possible values: $\mathcal{X}^m = \{x^{m,1}, \dots, x^{m,Q_m}\}$

Task: Finding the best estimate of

- $\theta_k := P(y^k), k \in [K]$
- $\theta_k^{m,q_m} := P(x^{q_m,m} | y^k), q_m \in [Q_m], k \in [K], m \in [M]$

Probability axioms:

- $\sum_{k=1}^K \theta_k = 1$
- $\sum_{q_m=1}^{Q_m} \theta_k^{m,q_m} = 1$ when fixing k and m

Maximum Likelihood Estimate

Basic setup: Given training data $\mathbf{D} = \{(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)\}$, count

- n_k : Number of training instances with label y^k
- n_k^{m, q_m} : Number of training instances with label y^k and feature X^m takes value x^{m, q_m}

Maximum Likelihood Estimate

Basic setup: Given training data $\mathbf{D} = \{(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)\}$, count

- n_k : Number of training instances with label y^k
- n_k^{m,q_m} : Number of training instances with label y^k and feature X^m takes value x^{m,q_m}

MLE gives us the best estimates

$$\theta_k := n_k / N \quad (2)$$

$$\theta_k^{m,q_m} := n_k^{m,q_m} / n_k \quad (3)$$

MLE (Exercise 2)

- $\mathcal{Y} = \{A, B, C\}$
- $\mathcal{X}^1 = \{d, e\}$
- $\mathcal{X}^2 = \{f, g, h\}$

n	Y	X^1	X^2
1	A	d	f
2	A	d	g
3	A	e	g
4	B	d	f
5	B	e	g
6	C	d	f
7	C	e	f
8	C	e	g

MLE (Exercise 2)

- $\mathcal{Y} = \{A, B, C\}$
- $\mathcal{X}^1 = \{d, e\}$
- $\mathcal{X}^2 = \{f, g, h\}$

$$n_A = 3 \quad n_B = 2 \quad n_C = 3$$
$$\theta_A = 3/8 \quad \theta_B = 1/4 \quad \theta_C = 3/8$$

n	Y	X^1	X^2
1	A	d	f
2	A	d	g
3	A	e	g
4	B	d	f
5	B	e	g
6	C	d	f
7	C	e	f
8	C	e	g

MLE (Exercise 2)

- $\mathcal{Y} = \{A, B, C\}$
- $\mathcal{X}^1 = \{d, e\}$
- $\mathcal{X}^2 = \{f, g, h\}$

$$n_A = 3 \quad n_B = 2 \quad n_C = 3$$

$$\theta_A = 3/8 \quad \theta_B = 1/4 \quad \theta_C = 3/8$$

n	Y	X^1	X^2
1	A	d	f
2	A	d	g
3	A	e	g
4	B	d	f
5	B	e	g
6	C	d	f
7	C	e	f
8	C	e	g

$n_A^{1,d} = 2$	$n_A^{1,e} = 1$	$\theta_A^{1,d} = 2/3$	$\theta_A^{1,e} = 1/3$
$n_B^{1,d} = 1$	$n_B^{1,e} = 1$	$\theta_B^{1,d} = 1/2$	$\theta_B^{1,e} = 1/2$
$n_C^{1,d} = 1$	$n_C^{1,e} = 2$	$\theta_C^{1,d} = 1/3$	$\theta_C^{1,e} = 2/3$
$n_A^{2,f} = ???$	$n_A^{2,g} = ???$	$\theta_A^{2,f} = ???$	$\theta_A^{2,g} = ???$
$n_B^{2,f} = ???$	$n_B^{2,g} = ???$	$\theta_B^{2,f} = ???$	$\theta_B^{2,g} = ???$
$n_C^{2,f} = ???$	$n_C^{2,g} = ???$	$\theta_C^{2,f} = ???$	$\theta_C^{2,g} = ???$

MLE (Exercise 2)

- $\mathcal{Y} = \{A, B, C\}$
- $\mathcal{X}^1 = \{d, e\}$
- $\mathcal{X}^2 = \{f, g, h\}$

$$n_A = 3 \quad n_B = 2 \quad n_C = 3$$

$$\theta_A = 3/8 \quad \theta_B = 1/4 \quad \theta_C = 3/8$$

n	Y	X^1	X^2
1	A	d	f
2	A	d	g
3	A	e	g
4	B	d	f
5	B	e	g
6	C	d	f
7	C	e	f
8	C	e	g

$n_A^{1,d} = 2$	$n_A^{1,e} = 1$	$\theta_A^{1,d} = 2/3$	$\theta_A^{1,e} = 1/3$
$n_B^{1,d} = 1$	$n_B^{1,e} = 1$	$\theta_B^{1,d} = 1/2$	$\theta_B^{1,e} = 1/2$
$n_C^{1,d} = 1$	$n_C^{1,e} = 2$	$\theta_C^{1,d} = 1/3$	$\theta_C^{1,e} = 2/3$
$n_A^{2,f} = ???$	$n_A^{2,g} = ???$	$\theta_A^{2,f} = ???$	$\theta_A^{2,g} = ???$
$n_B^{2,f} = ???$	$n_B^{2,g} = ???$	$\theta_B^{2,f} = ???$	$\theta_B^{2,g} = ???$
$n_C^{2,f} = ???$	$n_C^{2,g} = ???$	$\theta_C^{2,f} = ???$	$\theta_C^{2,g} = ???$

$$n_A^{2,h} = ??? \quad n_B^{2,h} = ??? \quad n_C^{2,h} = ???$$

$$\theta_A^{2,h} = ??? \quad \theta_B^{2,h} = ??? \quad \theta_C^{2,h} = ???$$

MLE (Solution 2)

- $\mathcal{Y} = \{A, B, C\}$
- $\mathcal{X}^1 = \{d, e\}$
- $\mathcal{X}^2 = \{f, g, h\}$

$$n_A = 3 \quad n_B = 2 \quad n_C = 3$$

$$\theta_A = 3/8 \quad \theta_B = 1/4 \quad \theta_C = 3/8$$

n	Y	X^1	X^2
1	A	d	f
2	A	d	g
3	A	e	g
4	B	d	f
5	B	e	g
6	C	d	f
7	C	e	f
8	C	e	g

$n_A^{1,d} = 2$	$n_A^{1,e} = 1$	$\theta_A^{1,d} = 2/3$	$\theta_A^{1,e} = 1/3$
$n_B^{1,d} = 1$	$n_B^{1,e} = 1$	$\theta_B^{1,d} = 1/2$	$\theta_B^{1,e} = 1/2$
$n_C^{1,d} = 1$	$n_C^{1,e} = 2$	$\theta_C^{1,d} = 1/3$	$\theta_C^{1,e} = 2/3$
$n_A^{2,f} = 1$	$n_A^{2,g} = 2$	$\theta_A^{2,f} = 1/3$	$\theta_A^{2,g} = 2/3$
$n_B^{2,f} = 1$	$n_B^{2,g} = 1$	$\theta_B^{2,f} = 1/2$	$\theta_B^{2,g} = 1/2$
$n_C^{2,f} = 2$	$n_C^{2,g} = 1$	$\theta_C^{2,f} = 2/3$	$\theta_C^{2,g} = 1/3$

$$n_A^{2,h} = 0 \quad n_B^{2,h} = 0 \quad n_C^{2,h} = 0$$

$$\theta_A^{2,h} = 0 \quad \theta_B^{2,h} = 0 \quad \theta_C^{2,h} = 0$$

Conditional Probabilities (Exercise 3)

Given $\mathbf{x} = (x^{1,q_1}, \dots, x^{M,q_M})$, for any $y^k \in \mathcal{Y}$:

$$P(y^k | \mathbf{x}) = \frac{\theta_k \prod_{m=1}^M \theta_k^{m,q_m}}{\sum_{y^{k'} \in \mathcal{Y}} \theta_{k'} \prod_{m=1}^M \theta_{k'}^{m,q_m}} \propto P'(y^k | \mathbf{x}) = \theta_k \prod_{m=1}^M \theta_k^{m,q_m}. \quad (4)$$

Conditional Probabilities (Exercise 3)

Given $\mathbf{x} = (x^{1,q_1}, \dots, x^{M,q_M})$, for any $y^k \in \mathcal{Y}$:

$$P(y^k | \mathbf{x}) = \frac{\theta_k \prod_{m=1}^M \theta_k^{m,q_m}}{\sum_{y^{k'} \in \mathcal{Y}} \theta_{k'} \prod_{m=1}^M \theta_{k'}^{m,q_m}} \propto P'(y^k | \mathbf{x}) = \theta_k \prod_{m=1}^M \theta_k^{m,q_m}. \quad (4)$$

$$\theta_A = 3/8 \quad \theta_B = 1/4 \quad \theta_C = 3/8$$

$$\theta_A^{1,d} = 2/3 \quad \theta_A^{1,e} = 1/3 \quad \theta_A^{2,f} = 1/3 \quad \theta_A^{2,g} = 2/3$$

$$\theta_B^{1,d} = 1/2 \quad \theta_B^{1,e} = 1/2 \quad \theta_B^{2,f} = 1/2 \quad \theta_B^{2,g} = 1/2$$

$$\theta_C^{1,d} = 1/3 \quad \theta_C^{1,e} = 2/3 \quad \theta_C^{2,f} = 2/3 \quad \theta_C^{2,g} = 1/3$$

$$\theta_A^{2,h} = 0 \quad \theta_B^{2,h} = 0 \quad \theta_C^{2,h} = 0$$

\mathbf{x}	$P'(A \mathbf{x})$	$P'(B \mathbf{x})$	$P'(C \mathbf{x})$
(d, f)	???	???	???
(e, h)	???	???	???

Conditional Probabilities (Solution 3)

Given $\mathbf{x} = (x^{1,q_1}, \dots, x^{M,q_M})$, for any $y^k \in \mathcal{Y}$:

$$P(y^k | \mathbf{x}) = \frac{\theta_k \prod_{m=1}^M \theta_k^{m,q_m}}{\sum_{y^{k'} \in \mathcal{Y}} \theta_{k'} \prod_{m=1}^M \theta_{k'}^{m,q_m}} \propto P'(y^k | \mathbf{x}) = \theta_k \prod_{m=1}^M \theta_k^{m,q_m}. \quad (5)$$

$$\theta_A = 3/8 \quad \theta_B = 1/4 \quad \theta_C = 3/8$$

$$\theta_A^{1,d} = 2/3 \quad \theta_A^{1,e} = 1/3 \quad \theta_A^{2,f} = 1/3 \quad \theta_A^{2,g} = 2/3$$

$$\theta_B^{1,d} = 1/2 \quad \theta_B^{1,e} = 1/2 \quad \theta_B^{2,f} = 1/2 \quad \theta_B^{2,g} = 1/2$$

$$\theta_C^{1,d} = 1/3 \quad \theta_C^{1,e} = 2/3 \quad \theta_C^{2,f} = 2/3 \quad \theta_C^{2,g} = 1/3$$

$$\theta_A^{2,h} = 0 \quad \theta_B^{2,h} = 0 \quad \theta_C^{2,h} = 0$$

\mathbf{x}	$P'(A \mathbf{x})$	$P'(B \mathbf{x})$	$P'(C \mathbf{x})$
(d, f)	1/12	1/16	1/12
(e, h)	0	0	0

Optimal Decision Rules (Exercise 4)

If $\ell(y^{k'}, y^k) = \mathbb{1}(y^{k'} \neq y^k)$, then (See Lecture 3+Check!)

$$y_{\ell}^{\theta}(\mathbf{x}) = \operatorname{argmax}_{y^k \in \mathcal{Y}} P'(y^k | \mathbf{x})$$

Optimal Decision Rules (Exercise 4)

If $\ell(y^{k'}, y^k) = \mathbb{1}(y^{k'} \neq y^k)$, then (See Lecture 3+Check!)

$$y_\ell^\theta(\mathbf{x}) = \operatorname{argmax}_{y^k \in \mathcal{Y}} P'(y^k | \mathbf{x})$$

\mathbf{x}	$P'(A \mathbf{x})$	$P'(B \mathbf{x})$	$P'(C \mathbf{x})$
(d, f)	$1/12$	$1/16$	$1/12$
(e, h)	0	0	0

Optimal Decision Rules (Exercise 4)

If $\ell(y^{k'}, y^k) = \mathbb{1}(y^{k'} \neq y^k)$, then (See Lecture 3+Check!)

$$y_\ell^\theta(\mathbf{x}) = \operatorname{argmax}_{y^k \in \mathcal{Y}} P'(y^k | \mathbf{x})$$

\mathbf{x}	$P'(A \mathbf{x})$	$P'(B \mathbf{x})$	$P'(C \mathbf{x})$
(d, f)	$1/12$	$1/16$	$1/12$
(e, h)	0	0	0

- If $\mathbf{x} = (d, f)$, then

$$y_\ell^\theta(\mathbf{x}) = ???, \quad (6)$$

- If $\mathbf{x} = (e, h)$, then

$$y_\ell^\theta(\mathbf{x}) = ???, \quad (7)$$

Optimal Decision Rules (Solution 4)

If $\ell(y^{k'}, y^k) = \mathbb{1}(y^{k'} \neq y^k)$, then (See Lecture 3+Check!)

$$y_\ell^\theta(\mathbf{x}) = \operatorname{argmax}_{y^k \in \mathcal{Y}} P'(y^k | \mathbf{x})$$

\mathbf{x}	$P'(A \mathbf{x})$	$P'(B \mathbf{x})$	$P'(C \mathbf{x})$
(d, f)	$1/12$	$1/16$	$1/12$
(e, h)	0	0	0

- If $\mathbf{x} = (d, f)$, then

$$y_\ell^\theta(\mathbf{x}) = \text{either } A \text{ or } C, \quad (8)$$

- If $\mathbf{x} = (e, h)$, then

$$y_\ell^\theta(\mathbf{x}) = \text{not well-defined}, \quad (9)$$

NBC + MLE: Comments

- May lead to **indecision** and **not well-defined** $P(y|x)$

NBC + MLE: Comments

- May lead to **indecision** and **not well-defined** $P(y|\mathbf{x})$

\mathbf{x}	$P'(A \mathbf{x})$	$P'(B \mathbf{x})$	$P'(C \mathbf{x})$
(d, f)	$1/12$	$1/16$	$1/12$
(e, h)	0	0	0

- May suffer from small numbers of observations
 - n_k : Number of training instances with label y^k
 - n_k^{m,q_m} : Number of training instances with label y^k and feature X^m takes value x^{m,q_m}

NBC + MLE: Comments

- May lead to **indecision** and **not well-defined** $P(y|\mathbf{x})$

\mathbf{x}	$P'(A \mathbf{x})$	$P'(B \mathbf{x})$	$P'(C \mathbf{x})$
(d, f)	$1/12$	$1/16$	$1/12$
(e, h)	0	0	0

- May suffer from small numbers of observations
 - n_k : Number of training instances with label y^k
 - n_k^{m,q_m} : Number of training instances with label y^k and feature X^m takes value x^{m,q_m}
- Does not (naturally) take into account missing/partial data

NBC + Dirichlet Model (DM)

Basic setup: Given training data $\mathbf{D} = \{(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)\}$, count

- n_k : Number of training instances with label y^k
- n_k^{m, q_m} : Number of instances with $Y = y^k$ and feature $X^m = x^{m, q_m}$

NBC + Dirichlet Model (DM)

Basic setup: Given training data $\mathbf{D} = \{(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)\}$, count

- n_k : Number of training instances with label y^k
- $n_k^{m,qm}$: Number of instances with $Y = y^k$ and feature $X^m = x^{m,qm}$

DM gives Bayesian estimates

$$\theta_k := (n_k + \alpha_k) / (N + s) = (n_k + sf_k) / (N + s) \quad (10)$$

$$\theta_k^{m,qm} := (n_k^{m,qm} + \alpha_k^{m,qm}) / (n_k + s) = (n_k^{m,qm} + sf_k^{m,qm}) / (n_k + s) \quad (11)$$

NBC + Dirichlet Model (DM)

Basic setup: Given training data $\mathbf{D} = \{(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)\}$, count

- n_k : Number of training instances with label y^k
- n_k^{m,q_m} : Number of instances with $Y = y^k$ and feature $X^m = x^{m,q_m}$

DM gives Bayesian estimates

$$\theta_k := (n_k + \alpha_k) / (N + s) = (n_k + sf_k) / (N + s) \quad (10)$$

$$\theta_k^{m,q_m} := (n_k^{m,q_m} + \alpha_k^{m,q_m}) / (n_k + s) = (n_k^{m,q_m} + sf_k^{m,q_m}) / (n_k + s) \quad (11)$$

Advocators	$\alpha_v (= y^k \text{ or } x^{m,q_m})$	s
Haldane (1948)	0	0
Perks (1947)	$1/ \mathcal{V} $	1
Jeffreys (1946, 1961)	1/2	$ \mathcal{V} /2$
Bayes-Laplace	1	$ \mathcal{V} $

NBC + DM (Exercise 5)

$$\theta_k := (n_k + 1/3)/(N+1),$$

$$\theta_k^{m, q_m} := (n_k^{m, q_m} + 1/|X^m|)/(n_k + 1).$$

NBC + DM (Exercise 5)

$$\theta_k := (n_k + 1/3)/(N+1),$$

$$\theta_k^{m,qm} := (n_k^{m,qm} + 1/|X^m|)/(n_k + 1).$$

- $\mathcal{Y} = \{A, B, C\}$
- $\mathcal{X}^1 = \{d, e\}$
- $\mathcal{X}^2 = \{f, g, h\}$

$$n_A = 3 \quad n_B = 2 \quad n_C = 3$$

$$\theta_A = 10/27 \quad \theta_B = 7/27 \quad \theta_C = 10/27$$

n	Y	X^1	X^2
1	A	d	f
2	A	d	g
3	A	e	g
4	B	d	f
5	B	e	g
6	C	d	f
7	C	e	f
8	C	e	g

$n_A^{1,d} = 2$	$n_A^{1,e} = 1$	$\theta_A^{1,d} = 5/8$	$\theta_A^{1,e} = 3/8$
$n_B^{1,d} = 1$	$n_B^{1,e} = 1$	$\theta_B^{1,d} = 3/6$	$\theta_B^{1,e} = 3/6$
$n_C^{1,d} = 1$	$n_C^{1,e} = 2$	$\theta_C^{1,d} = 3/8$	$\theta_C^{1,e} = 5/8$
$n_A^{2,f} = 1$	$n_A^{2,g} = 2$	$\theta_A^{2,f} = ???$	$\theta_A^{2,g} = ???$
$n_B^{2,f} = 1$	$n_B^{2,g} = 1$	$\theta_B^{2,f} = ???$	$\theta_B^{2,g} = ???$
$n_C^{2,f} = 2$	$n_C^{2,g} = 1$	$\theta_C^{2,f} = ???$	$\theta_C^{2,g} = ???$

$$n_A^{2,h} = 0 \quad n_B^{2,h} = 0 \quad n_C^{2,h} = 0$$

$$\theta_A^{2,h} = ??? \quad \theta_B^{2,h} = ??? \quad \theta_C^{2,h} = ???$$

NBC + DM (Solution 5)

$$\theta_k := (n_k + 1/3)/(N + 1),$$

$$\theta_k^{m,qm} := (n_k^{m,qm} + 1/|X^m|)/(n_k + 1).$$

- $\mathcal{Y} = \{A, B, C\}$

- $\mathcal{X}^1 = \{d, e\}$

- $\mathcal{X}^2 = \{f, g, h\}$

$$n_A = 3 \quad n_B = 2 \quad n_C = 3$$

$$\theta_A = 10/27 \quad \theta_B = 7/27 \quad \theta_C = 10/27$$

n	Y	X^1	X^2
1	A	d	f
2	A	d	g
3	A	e	g
4	B	d	f
5	B	e	g
6	C	d	f
7	C	e	f
8	C	e	g

$n_A^{1,d} = 2$	$n_A^{1,e} = 1$	$\theta_A^{1,d} = 5/8$	$\theta_A^{1,e} = 3/8$
$n_B^{1,d} = 1$	$n_B^{1,e} = 1$	$\theta_B^{1,d} = 3/6$	$\theta_B^{1,e} = 3/6$
$n_C^{1,d} = 1$	$n_C^{1,e} = 2$	$\theta_C^{1,d} = 3/8$	$\theta_C^{1,e} = 5/8$
$n_A^{2,f} = 1$	$n_A^{2,g} = 2$	$\theta_A^{2,f} = 4/12$	$\theta_A^{2,g} = 7/12$
$n_B^{2,f} = 1$	$n_B^{2,g} = 1$	$\theta_B^{2,f} = 4/9$	$\theta_B^{2,g} = 4/9$
$n_C^{2,f} = 2$	$n_C^{2,g} = 1$	$\theta_C^{2,f} = 7/12$	$\theta_C^{2,g} = 4/12$

$$n_A^{2,h} = 0 \quad n_B^{2,h} = 0 \quad n_C^{2,h} = 0$$

$$\theta_A^{2,h} = 1/12 \quad \theta_B^{2,h} = 1/9 \quad \theta_C^{2,h} = 1/12$$

Conditional Probabilities (Exercise 6)

Given $\mathbf{x} = (x^{1,q_1}, \dots, x^{M,q_M})$, for any $y^k \in \mathcal{Y}$:

$$P(y^k | \mathbf{x}) = \frac{\theta_k \prod_{m=1}^M \theta_k^{m,q_m}}{\sum_{y^{k'} \in \mathcal{Y}} \theta_{k'} \prod_{m=1}^M \theta_{k'}^{m,q_m}} \propto P'(y^k | \mathbf{x}) = \theta_k \prod_{m=1}^M \theta_k^{m,q_m}. \quad (12)$$

$$\theta_A = 10/27 \quad \theta_B = 7/27 \quad \theta_C = 10/27$$

$$\theta_A^{1,d} = 5/8 \quad \theta_A^{1,e} = 3/8 \quad \theta_A^{2,f} = 4/12 \quad \theta_A^{2,g} = 7/12$$

$$\theta_B^{1,d} = 3/6 \quad \theta_B^{1,e} = 3/6 \quad \theta_B^{2,f} = 4/9 \quad \theta_B^{2,g} = 4/9$$

$$\theta_C^{1,d} = 3/8 \quad \theta_C^{1,e} = 5/8 \quad \theta_C^{2,f} = 7/12 \quad \theta_C^{2,g} = 4/12$$

$$\theta_A^{2,h} = 1/12 \quad \theta_B^{2,h} = 1/9 \quad \theta_C^{2,h} = 1/12$$

\mathbf{x}	$P'(A \mathbf{x})$	$P'(B \mathbf{x})$	$P'(C \mathbf{x})$	$y_\ell^\theta(\mathbf{x})$
(d, f)	???	???	???	???
(e, h)	???	???	???	???

Conditional Probabilities (Solution 6)

Given $\mathbf{x} = (x^{1,q_1}, \dots, x^{M,q_M})$, for any $y^k \in \mathcal{Y}$:

$$P(y^k|\mathbf{x}) = \frac{\theta_k \prod_{m=1}^M \theta_k^{m,q_m}}{\sum_{y^{k'} \in \mathcal{Y}} \theta_{k'} \prod_{m=1}^M \theta_{k'}^{m,q_m}} \propto P'(y^k|\mathbf{x}) = \theta_k \prod_{m=1}^M \theta_k^{m,q_m}. \quad (13)$$

$$\theta_A = 10/27 \quad \theta_B = 7/27 \quad \theta_C = 10/27$$

$$\theta_A^{1,d} = 5/8 \quad \theta_A^{1,e} = 3/8 \quad \theta_A^{2,f} = 4/12 \quad \theta_A^{2,g} = 7/12$$

$$\theta_B^{1,d} = 3/6 \quad \theta_B^{1,e} = 3/6 \quad \theta_B^{2,f} = 4/9 \quad \theta_B^{2,g} = 4/9$$

$$\theta_C^{1,d} = 3/8 \quad \theta_C^{1,e} = 5/8 \quad \theta_C^{2,f} = 7/12 \quad \theta_C^{2,g} = 4/12$$

$$\theta_A^{2,h} = 1/12 \quad \theta_B^{2,h} = 1/9 \quad \theta_C^{2,h} = 1/12$$

\mathbf{x}	$P'(A \mathbf{x})$	$P'(B \mathbf{x})$	$P'(C \mathbf{x})$	$y_\ell^\theta(\mathbf{x})$
(d, f)	$\frac{10}{27} \frac{5}{8} \frac{4}{12}$	$\frac{7}{27} \frac{3}{6} \frac{4}{9}$	$\frac{10}{27} \frac{3}{8} \frac{7}{12}$	C
(e, h)	$\frac{10}{27} \frac{3}{8} \frac{1}{12}$	$\frac{7}{27} \frac{3}{6} \frac{1}{9}$	$\frac{10}{27} \frac{5}{8} \frac{1}{12}$	C

NBC + DM: Comments

- May lead to **indecision**, but can avoid **not well-defined** $P(y|\mathbf{x})$

NBC + DM: Comments

- May lead to **indecision**, but can avoid **not well-defined** $P(y|\mathbf{x})$

\mathbf{x}	$P'(A \mathbf{x})$	$P'(B \mathbf{x})$	$P'(C \mathbf{x})$	$y_\ell^\theta(\mathbf{x})$
(d, f)	$\frac{10}{27} \frac{5}{8} \frac{4}{12}$	$\frac{7}{27} \frac{3}{6} \frac{4}{9}$	$\frac{10}{27} \frac{3}{8} \frac{7}{12}$	C
(e, h)	$\frac{10}{27} \frac{3}{8} \frac{1}{12}$	$\frac{7}{27} \frac{3}{6} \frac{1}{9}$	$\frac{10}{27} \frac{5}{8} \frac{1}{12}$	C

- May suffer from small numbers of observations
 - n_k : Number of training instances with label y^k
 - n_k^{m, q_m} : Number of training instances with label y^k and feature X^m takes value x^{m, q_m}

NBC + DM: Comments

- May lead to **indecision**, but can avoid **not well-defined** $P(y|\mathbf{x})$

\mathbf{x}	$P'(A \mathbf{x})$	$P'(B \mathbf{x})$	$P'(C \mathbf{x})$	$y_\ell^\theta(\mathbf{x})$
(d, f)	$\frac{10}{27} \frac{5}{8} \frac{4}{12}$	$\frac{7}{27} \frac{3}{6} \frac{4}{9}$	$\frac{10}{27} \frac{3}{8} \frac{7}{12}$	C
(e, h)	$\frac{10}{27} \frac{3}{8} \frac{1}{12}$	$\frac{7}{27} \frac{3}{6} \frac{1}{9}$	$\frac{10}{27} \frac{5}{8} \frac{1}{12}$	C

- May suffer from small numbers of observations
 - n_k : Number of training instances with label y^k
 - n_k^{m,q_m} : Number of training instances with label y^k and feature X^m takes value x^{m,q_m}
- Does not (naturally) take into account missing/partial data

Outline

- Graphical Interpretation of Probabilistic Models
- Naïve Bayesian/Credal classifiers
 - Naïve Bayesian classifier
 - Naïve Credal classifiers
- Decision Trees
- Bayesian Neural Networks

Naïve Credal classifiers (NCC)

Basic setup: Given training data $\mathbf{D} = \{(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)\}$, count

- n_k : Number of training instances with label y^k
- n_k^{m, q_m} : Number of instances with $Y = y^k$ and feature $X^m = x^{m, q_m}$

Naïve Credal classifiers (NCC)

Basic setup: Given training data $\mathbf{D} = \{(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)\}$, count

- n_k : Number of training instances with label y^k
- n_k^{m, q_m} : Number of instances with $Y = y^k$ and feature $X^m = x^{m, q_m}$

Imprecise Dirichlet model (IDM) gives

$$\underline{\theta}_k := n_k / (N + s) \quad (14)$$

$$\underline{\theta}_k^{m, q_m} := n_k^{m, q_m} / (n_k + s) \quad (15)$$

$$\bar{\theta}_k := (n_k + s) / (N + s) \quad (16)$$

$$\bar{\theta}_k^{m, q_m} := (n_k^{m, q_m} + s) / (n_k + s) \quad (17)$$

Naïve Credal classifiers (NCC)

Basic setup: Given training data $\mathbf{D} = \{(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)\}$, count

- n_k : Number of training instances with label y^k
- n_k^{m, q_m} : Number of instances with $Y = y^k$ and feature $X^m = x^{m, q_m}$

Imprecise Dirichlet model (IDM) gives

$$\underline{\theta}_k := n_k / (N + s) \quad (14) \quad \bar{\theta}_k := (n_k + s) / (N + s) \quad (16)$$

$$\underline{\theta}_k^{m, q_m} := n_k^{m, q_m} / (n_k + s) \quad (15) \quad \bar{\theta}_k^{m, q_m} := (n_k^{m, q_m} + s) / (n_k + s) \quad (17)$$

IDM + ϵ regularization [2]

$$\underline{\theta}_k := (n_k + s\underline{\epsilon}_k) / (N + s) \quad (18) \quad \bar{\theta}_k := (n_k + s\bar{\epsilon}_k) / (N + s) \quad (20)$$

$$\underline{\theta}_k^{m, q_m} := (n_k^{m, q_m} + s\underline{\epsilon}_k^{m, q_m}) / (n_k + s) \quad (19) \quad \bar{\theta}_k^{m, q_m} := (n_k^{m, q_m} + s\bar{\epsilon}_k^{m, q_m}) / (n_k + s) \quad (21)$$

Interval Conditional Probabilities

Given a query instance $\mathbf{x} = (x^{q_1,1}, x^{q_2,2}, \dots, x^{q_M,M})$, we have

$$\frac{1}{\overline{P}(y^k|\mathbf{x})} - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s \underline{\epsilon}_k}{n_k + s \overline{\epsilon}_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_m,m} + s \underline{\epsilon}_k^{m,q_m}}{n_k^{q_m,m} + s \overline{\epsilon}_k^{m,q_m}} \right),$$

$$\frac{1}{\underline{P}(y^k|\mathbf{x})} - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s \overline{\epsilon}_k}{n_k + s \underline{\epsilon}_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_m,m} + s \overline{\epsilon}_k^{m,q_m}}{n_k^{q_m,m} + s \underline{\epsilon}_k^{m,q_m}} \right).$$

Interval Conditional Probabilities

Given a query instance $\mathbf{x} = (x^{q_1,1}, x^{q_2,2}, \dots, x^{q_M,M})$, we have

$$\frac{1}{\bar{P}(y^k|\mathbf{x})} - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s \underline{\epsilon}_k}{n_k + s \bar{\epsilon}_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_m,m} + s \underline{\epsilon}_k^{m,q_m}}{n_k^{q_m,m} + s \bar{\epsilon}_k^{m,q_m}} \right),$$

$$\frac{1}{\underline{P}(y^k|\mathbf{x})} - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s \bar{\epsilon}_k}{n_k + s \underline{\epsilon}_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_m,m} + s \bar{\epsilon}_k^{m,q_m}}{n_k^{q_m,m} + s \underline{\epsilon}_k^{m,q_m}} \right).$$

$$\mathcal{P}(\mathcal{Y}|\mathbf{x}) := \left\{ P(\mathcal{Y}|\mathbf{x}) \mid P(y^k|\mathbf{x}) \in [\underline{P}(y^k|\mathbf{x}), \bar{P}(y^k|\mathbf{x})], \sum_{k=1}^K P(y^k|\mathbf{x}) = 1 \right\}.$$

Interval Conditional Probabilities (Exercise 7)

$$1/\overline{P}(y^k|\mathbf{x}) - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s \underline{\epsilon}_k}{n_k + s \overline{\epsilon}_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_{m,m}} + s \underline{\epsilon}_k^{m, q_m}}{n_k^{q_{m,m}} + s \overline{\epsilon}_k^{m, q_m}} \right),$$

$$1/\underline{P}(y^k|\mathbf{x}) - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s \overline{\epsilon}_k}{n_k + s \underline{\epsilon}_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_{m,m}} + s \overline{\epsilon}_k^{m, q_m}}{n_k^{q_{m,m}} + s \underline{\epsilon}_k^{m, q_m}} \right).$$

Interval Conditional Probabilities (Exercise 7)

$$1/\overline{P}(y^k|\mathbf{x}) - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s \underline{\epsilon}_k}{n_k + s \overline{\epsilon}_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_{m,m}} + s \underline{\epsilon}_k^{m, q_m}}{n_k^{q_{m,m}} + s \overline{\epsilon}_k^{m, q_m}} \right),$$

$$1/\underline{P}(y^k|\mathbf{x}) - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s \overline{\epsilon}_k}{n_k + s \underline{\epsilon}_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_{m,m}} + s \overline{\epsilon}_k^{m, q_m}}{n_k^{q_{m,m}} + s \underline{\epsilon}_k^{m, q_m}} \right).$$

$$s = 1 \quad \underline{\epsilon}_k = 0.01 \quad \overline{\epsilon}_k = 0.99 \quad n_A = 3 \quad n_B = 2 \quad n_C = 3$$

$$n_A^{1,d} = 2 \quad n_A^{1,e} = 1 \quad n_A^{2,f} = 1 \quad n_A^{2,g} = 2$$

$$n_B^{1,d} = 1 \quad n_B^{1,e} = 1 \quad n_B^{2,f} = 1 \quad n_B^{2,g} = 1$$

$$n_C^{1,d} = 1 \quad n_C^{1,e} = 2 \quad n_C^{2,f} = 2 \quad n_C^{2,g} = 1$$

$$n_A^{2,h} = 0 \quad n_B^{2,h} = 0 \quad n_C^{2,h} = 0$$

\mathbf{x}	$\underline{P}(A \mathbf{x})$	$\underline{P}(B \mathbf{x})$	$\underline{P}(C \mathbf{x})$	$\overline{P}(A \mathbf{x})$	$\overline{P}(B \mathbf{x})$	$\overline{P}(C \mathbf{x})$
(d, f)	???	???	???	???	???	???
(e, h)	???	???	???	???	???	???

Interval Conditional Probabilities (Solution 7)

$$1/\overline{P}(y^k|\mathbf{x}) - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s\epsilon_k}{n_k + s\bar{\epsilon}_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_{m,m}} + s\epsilon_k^{m,q_m}}{n_k^{q_{m,m}} + s\bar{\epsilon}_k^{m,q_m}} \right),$$

$$1/\underline{P}(y^k|\mathbf{x}) - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s\bar{\epsilon}_k}{n_k + s\epsilon_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_{m,m}} + s\bar{\epsilon}_k^{m,q_m}}{n_k^{q_{m,m}} + s\epsilon_k^{m,q_m}} \right).$$

Interval Conditional Probabilities (Solution 7)

$$1/\overline{P}(y^k|\mathbf{x}) - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s \underline{\epsilon}_k}{n_k + s \overline{\epsilon}_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_{m,m}} + s \underline{\epsilon}_k^{m, q_m}}{n_k^{q_{m,m}} + s \overline{\epsilon}_k^{m, q_m}} \right),$$

$$1/\underline{P}(y^k|\mathbf{x}) - 1 = \sum_{k' \neq k} \left(\frac{n_{k'} + s \overline{\epsilon}_k}{n_k + s \underline{\epsilon}_k} \left(\frac{n_k + s}{n_{k'} + s} \right)^M \prod_{m=1}^M \frac{n_{k'}^{q_{m,m}} + s \overline{\epsilon}_k^{m, q_m}}{n_k^{q_{m,m}} + s \underline{\epsilon}_k^{m, q_m}} \right).$$

$$s = 1 \quad \underline{\epsilon}_k = 0.01 \quad \overline{\epsilon}_k = 0.99 \quad n_A = 3 \quad n_B = 2 \quad n_C = 3$$

$$n_A^{1,d} = 2 \quad n_A^{1,e} = 1 \quad n_A^{2,f} = 1 \quad n_A^{2,g} = 2$$

$$n_B^{1,d} = 1 \quad n_B^{1,e} = 1 \quad n_B^{2,f} = 1 \quad n_B^{2,g} = 1$$

$$n_C^{1,d} = 1 \quad n_C^{1,e} = 2 \quad n_C^{2,f} = 2 \quad n_C^{2,g} = 1$$

$$n_A^{2,h} = 0 \quad n_B^{2,h} = 0 \quad n_C^{2,h} = 0$$

\mathbf{x}	$\underline{P}(A \mathbf{x})$	$\underline{P}(B \mathbf{x})$	$\underline{P}(C \mathbf{x})$	$\overline{P}(A \mathbf{x})$	$\overline{P}(B \mathbf{x})$	$\overline{P}(C \mathbf{x})$
(d, f)	???	???	???	???	???	???
(e, h)	???	???	???	???	???	???

Making Set-Valued Predictions (Recap)

For each instance \mathbf{x} , let

- $\theta \leftarrow P(\mathcal{Y}|\mathbf{x})$ and $\Theta \leftarrow \mathcal{P}(\mathcal{Y}|\mathbf{x})$

E-admissibility Rule:

- An optimal prediction is

$$\mathbf{Y}_{\ell, \Theta}^E = \{y \in \mathcal{Y} \mid \exists \theta \in \Theta \text{ s.t. } y = y_{\ell}^{\theta}\}.$$

- Computation: Solving linear programs [10], etc.

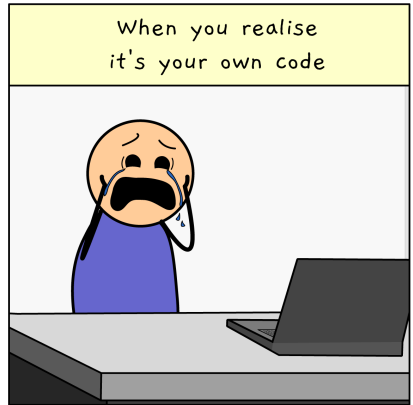
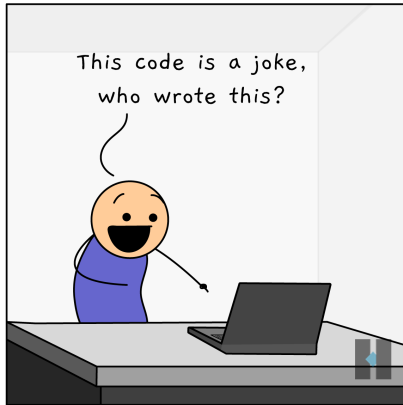
Maximality Rule:

- An optimal prediction is

$$\mathbf{Y}_{\ell, \Theta}^M = \{y \in \mathcal{Y} \mid \nexists y' \text{ s.t. } y' >_{\ell, \Theta} y\}.$$

- Computation: Solving linear programs [10], Iterating over the extreme points of Θ [10], **exploiting the properties of NCC [3]**.

Illustrative Examples → Lines of Code Would be Useful!



f /techindustan

🐦 /techindustan

📷 /techindustan

NCC: Comments

NCC inherits properties of IDM [3]:

- May lead to **set-valued predictions**
- ϵ -regularization can avoid **not well-defined** $P(y|\mathbf{x})$
- May provide reliable interval probabilities when seeing small numbers of observations
 - n_k : Number of training instances with label y^k
 - n_k^{m,q_m} : Number of training instances with label y^k and feature X^m takes value x^{m,q_m}

NCC: Comments

NCC inherits properties of IDM [3]:

- May lead to **set-valued predictions**
- ϵ -regularization can avoid **not well-defined** $P(y|\mathbf{x})$
- May provide reliable interval probabilities when seeing small numbers of observations
 - n_k : Number of training instances with label y^k
 - n_k^{m,q_m} : Number of training instances with label y^k and feature X^m takes value x^{m,q_m}
- Provide tools to (naturally) take into account missing/partial data
 - Naive solutions are computationally expensive (in exponential time)
 - **More efficient (polynomial-time) procedure exists**

NCC: Technical Details + Performance

Journal of Machine Learning Research 9 (2008) 581-621

Submitted 1/07; Revised 2/08; Published 4/08

Learning Reliable Classifiers From Small or Incomplete Data Sets: The Naïve Credal Classifier 2

Giorgio Corani
Marco Zaffalon

IDSIA

*Istituto Dalle Molle di Studi sull'Intelligenza Artificiale
CH-6928 Manno (Lugano), Switzerland*

GIORGIO@IDSIA.CH
ZAFFALON@IDSIA.CH

Editor: Charles Elkan

Outline

- Graphical Interpretation of Probabilistic Models
- Naïve Bayesian/Credal classifiers
- **Decision Trees**
 - Decision Trees
 - Credal Decision Trees
- Bayesian Neural Networks

Outline

- Graphical Interpretation of Probabilistic Models
- Naïve Bayesian/Credal classifiers
- **Decision Trees**
 - Decision Trees
 - Credal Decision Trees
- Bayesian Neural Networks

Discriminative Models

Probabilistic Models:

- Estimate $P(Y, \mathbf{X})$
- Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y|\text{pa}(Y)) \prod_{m=1}^M P(X^m|\text{pa}(X^m)).$$

Extreme Cases:

- Discriminative models: $Y \notin \text{pa}(X^m)$, $m \in [M]$
- Generative models: $\text{pa}(Y) = \emptyset$ and $Y \in \text{pa}(X^m)$, $m \in [M]$.

Model Families:

- How to encode/parametrize $P(Y|\text{pa}(Y))$ and $P(X^m|\text{pa}(X^m))$.
- How to estimate $P(Y, \mathbf{X}) = P(Y|\mathbf{X})P(\mathbf{X})$ from training data.

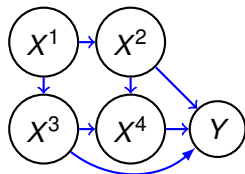
Discriminative Models: Structure (Exercise 8)

Let's start with an example where one wishes to model

$$P(Y, \mathbf{X}) = P(Y, X^1, X^2, X^3, X^4).$$

Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y | \text{pa}(Y)) \prod_{m=1}^M P(X^m | \text{pa}(X^m)).$$



- $\text{pa}(Y) = ???$,
- $\text{pa}(X^1) = \emptyset$, $\text{pa}(X^2) = ???$
- $\text{pa}(X^3) = ???$
- $\text{pa}(X^4) = ???$

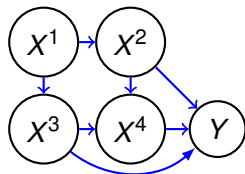
Discriminative Models: Structure (Exercise 8)

Let's start with an example where one wishes to model

$$P(Y, \mathbf{X}) = P(Y, X^1, X^2, X^3, X^4).$$

Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y | \text{pa}(Y)) \prod_{m=1}^M P(X^m | \text{pa}(X^m)).$$



- $\text{pa}(Y) = ???$,
- $\text{pa}(X^1) = \emptyset$, $\text{pa}(X^2) = ???$
- $\text{pa}(X^3) = ???$
- $\text{pa}(X^4) = ???$

Chain rule gives us

$$P(Y, \mathbf{X}) = ???.$$

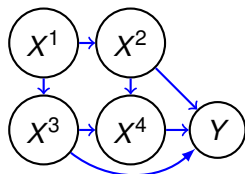
Discriminative Models: Structure (Solution 8)

Let's start with an example where one wishes to model

$$P(Y, \mathbf{X}) = P(Y, X^1, X^2, X^3, X^4).$$

Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y|\text{pa}(Y)) \prod_{m=1}^M P(X^m|\text{pa}(X^m)).$$



- $\text{pa}(Y) = \{X^2, X^3, X^4\}$
- $\text{pa}(X^1) = \emptyset$, $\text{pa}(X^2) = \{X^1\}$
- $\text{pa}(X^3) = \{X^1\}$
- $\text{pa}(X^4) = \{X^2, X^3\}$

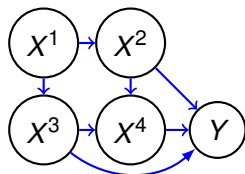
Discriminative Models: Structure (Solution 8)

Let's start with an example where one wishes to model

$$P(Y, \mathbf{X}) = P(Y, X^1, X^2, X^3, X^4).$$

Chain rule (probability):

$$P(Y, \mathbf{X}) = P(Y|\text{pa}(Y)) \prod_{m=1}^M P(X^m|\text{pa}(X^m)).$$



- $\text{pa}(Y) = \{X^2, X^3, X^4\}$
- $\text{pa}(X^1) = \emptyset$, $\text{pa}(X^2) = \{X^1\}$
- $\text{pa}(X^3) = \{X^1\}$
- $\text{pa}(X^4) = \{X^2, X^3\}$

Chain rule gives us

$$P(Y, \mathbf{X}) = P(Y|X^2, X^3, X^4) P(X^1) P(X^2|X^1) P(X^3|X^1) P(X^4|X^2, X^3).$$

Classification Task (Exercise + Solution 9)

Prove that

$$P(y|\mathbf{x}) = P(y|pa(y)). \quad (22)$$

Classification Task (Exercise + Solution 9)

Prove that

$$P(y|\mathbf{x}) = P(y|pa(y)). \quad (22)$$

We have

$$P(y|\mathbf{x}) = \frac{P(y, \mathbf{x})}{\sum_{y' \in \mathcal{Y}} P(y', \mathbf{x})} \quad (23)$$

$$= \frac{P(y|pa(y)) \prod_{m=1}^M P(x^m|pa(x^m))}{\sum_{y' \in \mathcal{Y}} P(y'|pa(y')) \prod_{m=1}^M P(x^m|pa(x^m))} \quad (24)$$

$$= \frac{\prod_{m=1}^M P(x^m|pa(x^m)) P(y|pa(y))}{\prod_{m=1}^M P(x^m|pa(x^m)) \sum_{y' \in \mathcal{Y}} P(y'|pa(y'))} \quad (25)$$

$$= \frac{P(y|pa(y))}{\sum_{y' \in \mathcal{Y}} P(y'|pa(y'))} \quad (26)$$

$$= P(y|pa(y)). \quad (27)$$

Classification Task: Comments

$P(Y|\mathbf{X})$ is extracted using **Bayes' theorem**

$$P(y|\mathbf{x}) = P(y|pa(y)). \quad (28)$$

Classification Task: Comments

$P(Y|\mathbf{X})$ is extracted using **Bayes' theorem**

$$P(y|\mathbf{x}) = P(y|pa(y)). \quad (28)$$

- Features outside $pa(Y)$ are redundant
- To solve the classification task, we only need $P(Y|pa(Y))$
- Commonly used assumption: $pa(Y) = (X^1, \dots, X^M)$

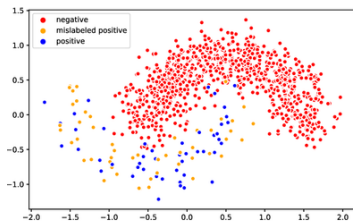
Classification Task: Comments

$P(Y|\mathbf{X})$ is extracted using **Bayes' theorem**

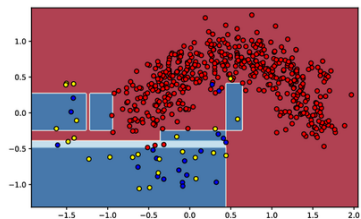
$$P(y|\mathbf{x}) = P(y|pa(y)). \quad (28)$$

- Features outside $pa(Y)$ are redundant
- To solve the classification task, we only need $P(Y|pa(Y))$
- Commonly used assumption: $pa(Y) = (X^1, \dots, X^M)$
- $P(Y|pa(Y))$ can be defined either **globally** or **locally**:
 - Logistic regression, neural nets, etc., define $P(Y|pa(Y))$ **globally**
 - Decision tree, model trees, etc., define $P(Y|pa(Y))$ **locally**
 - Decision tree **does not** require $pa(Y) = (X^1, \dots, X^M)$

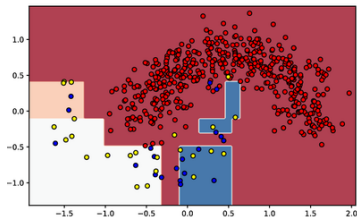
Decision Trees: Example [11]



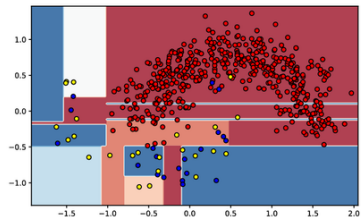
(a) Half-moons data set (ground truth)



(b) PU-Hellinger Decision Tree on the test set



(c) Hellinger Decision Tree on the test set



(d) CART on the test set

Decision Trees: (Informal+Probabilistic) Definition

A decision trees is

- a collection of non-overlapping leaves L_1, \dots, L_H
- where $L_1 \cap \dots \cap L_H = \mathcal{X}$
- and each leaf L_h has its own $P_h(Y|pa(Y))$

Decision Trees: (Informal+Probabilistic) Definition

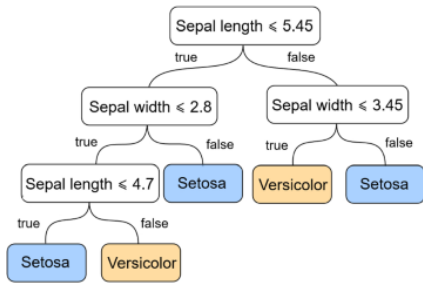
A decision trees is

- a collection of non-overlapping leaves L_1, \dots, L_H
- where $L_1 \cap \dots \cap L_H = \mathcal{X}$
- and each leaf L_h has its own $P_h(Y|pa(Y))$

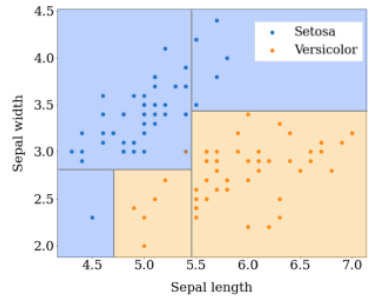
Learning an optimal decision tree from training data

- can be extremely hard (due to huge numbers of possible trees)
- and is often done approximately (top-down induction, bottom-up induction, etc.)

Top-Down Induction (Example) [4]



(a) Tree visualization



(b) Partitioning visualization

Top-down induction: Steps

Basic Setup:

- Training data $\mathbf{D} = \{(y^n, \mathbf{x}^n) | n \in [N]\}$
- Local hypothesis space $P(Y|pa(Y)) \in \mathcal{P}(Y|pa(Y))$
- An uncertainty measure U or a loss function ℓ ← assess how good/bad each **local classifier** is

Top-down induction: Steps

Basic Setup:

- Training data $\mathbf{D} = \{(y^n, \mathbf{x}^n) | n \in [N]\}$
- Local hypothesis space $P(Y | \text{pa}(Y)) \in \mathcal{P}(Y | \text{pa}(Y))$
- An uncertainty measure U or a loss function ℓ \leftarrow assess how good/bad each **local classifier** is

Induction protocol:

- Recursively partition the feature space \mathcal{X}
- From the current node, choose the best split which improves the evaluation criterion
- Evaluation criteria: Information gain, entropy, Gini score, etc.,
- Stopping criteria: No more gain on evaluation criterion U or ℓ

Splitting Criteria: Entropy (Frequentist)

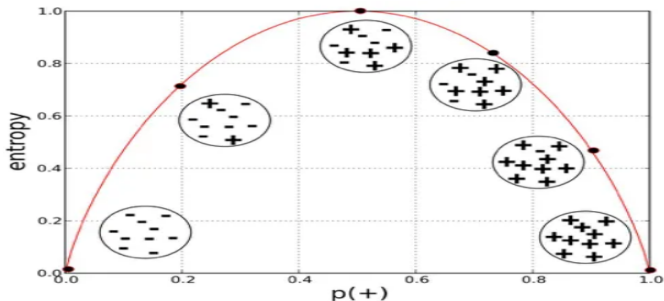
- Entropy of a node $\mathbf{D}_h \subset \mathbf{D}$ with $P(\mathcal{Y}|\mathbf{D}_h)$

$$U_E(P(\mathcal{Y}|\mathbf{D}_h)) = - \sum_{y \in \mathcal{Y}} P(y|\mathbf{D}_h) \log_2(P(y|\mathbf{D}_h)).$$

Splitting Criteria: Entropy (Frequentist)

- Entropy of a node $\mathbf{D}_h \subset \mathbf{D}$ with $P(\mathcal{Y}|\mathbf{D}_h)$

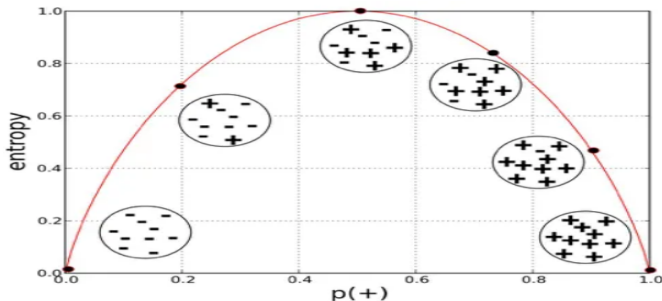
$$U_E(P(\mathcal{Y}|\mathbf{D}_h)) = - \sum_{y \in \mathcal{Y}} P(y|\mathbf{D}_h) \log_2(P(y|\mathbf{D}_h)).$$



Splitting Criteria: Entropy (Frequentist)

- Entropy of a node $\mathbf{D}_h \subset \mathbf{D}$ with $P(\mathcal{Y}|\mathbf{D}_h)$

$$U_E(P(\mathcal{Y}|\mathbf{D}_h)) = - \sum_{y \in \mathcal{Y}} P(y|\mathbf{D}_h) \log_2(P(y|\mathbf{D}_h)).$$



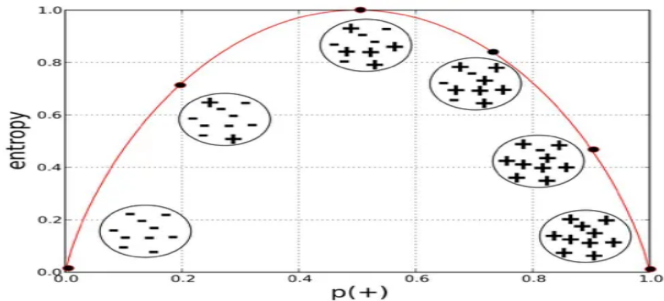
- For each possible split $\mathbf{D}_h = \mathbf{D}_h^1 \cup \mathbf{D}_h^2$, its entropy is

$$U_E(\mathbf{D}_h^1 \cup \mathbf{D}_h^2) = U_E(P(\mathcal{Y}|\mathbf{D}_h^1)) + U_E(P(\mathcal{Y}|\mathbf{D}_h^2)).$$

Compute Entropy (Exercise 10)

- Entropy of a node $\mathbf{D}_h \subset \mathbf{D}$ with $P(\mathcal{Y}|\mathbf{D}_h)$

$$U_E(P(y|\mathbf{D}_h)) = - \sum_{y \in \mathcal{Y}} P(y|\mathbf{D}_h) \log_2(P(y|\mathbf{D}_h)).$$

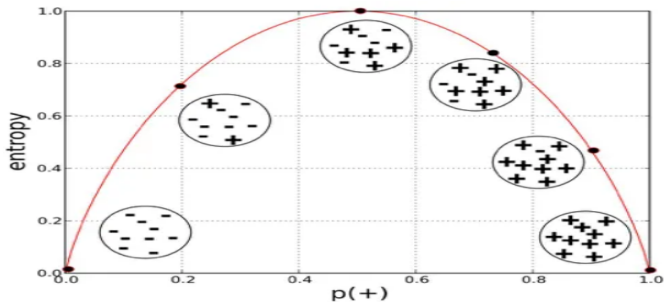


- Entropy of the bottom left node is ???
- Entropy of the top node is ???

Compute Entropy (Solution 10)

- Entropy of a node $\mathbf{D}_h \subset \mathbf{D}$ with $P(\mathcal{Y}|\mathbf{D}_h)$

$$U_E(P(y|\mathbf{D}_h)) = - \sum_{y \in \mathcal{Y}} P(y|\mathbf{D}_h) \log_2(P(y|\mathbf{D}_h)).$$



- Entropy of the bottom left node is 0
- Entropy of the top node is $-0.5 \log_2(0.5) + 0.5 \log_2(0.5) = 1$

Splitting Criteria: Bayesian

In principle, we can employ Dirichlet models (DM) to

- derive Bayesian estimates of $P(\mathcal{Y}|\mathbf{D}_h)$ and/or $U(P(y|\mathbf{D}_h))$
- and modify the top-down induction steps.

Splitting Criteria: Bayesian

In principle, we can employ Dirichlet models (DM) to

- derive Bayesian estimates of $P(\mathcal{Y}|\mathbf{D}_h)$ and/or $U(P(y|\mathbf{D}_h))$
- and modify the top-down induction steps.

- I haven't seen such decision decision trees
- I hope I can find some reference soon ...

Outline

- Graphical Interpretation of Probabilistic Models
- Naïve Bayesian/Credal classifiers
- **Decision Trees**
 - Decision Trees
 - **Credal Decision Trees**
- Bayesian Neural Networks

Where and How to be Imprecise?

In principle, we can employ Imprecise Dirichlet models (IDM) to

- derive interval estimates of $P(\mathcal{Y}|\mathbf{D}_h)$ and/or $U(P(y|\mathbf{D}_h))$
- and modify the top-down induction steps.

Where and How to be Imprecise?

In principle, we can employ Imprecise Dirichlet models (IDM) to

- derive interval estimates of $P(\mathcal{Y}|\mathbf{D}_h)$ and/or $U(P(y|\mathbf{D}_h))$
- and modify the top-down induction steps.

Credal Decision Trees [1, 7]

- Use IDM to derive interval estimates $\mathcal{P}(\mathcal{Y}|\mathbf{D}_h)$ of $P(\mathcal{Y}|\mathbf{D}_h)$
- Seek the highest entropy

$$U(\mathcal{P}(\mathcal{Y}|\mathbf{D}_h)) = \max_{P \in \mathcal{P}} U(P(\mathcal{Y}|\mathbf{D}_h)). \quad (29)$$

- Each leaf is equipped a $\mathcal{P}(\mathcal{Y}|\mathbf{L}_h) \rightarrow$ Precise predictions.

Credal Decision Trees: Performance [7]

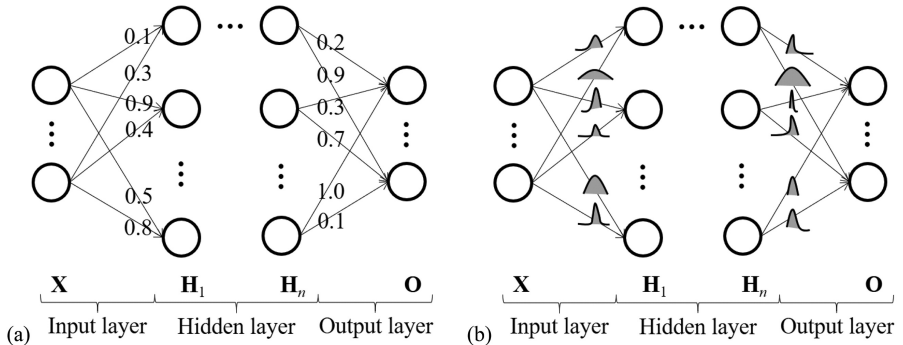
Splitting criterion	0% noise	10% noise	20% noise
Info-Gain (IG)	78.96	77.49	74.76
Info-Gain Ratio (IGR)	78.97	77.66	75.14
Imprecise Info-Gain (IIG)	79.56	78.65	76.72
Complete IIG (CIIG)	79.63	78.66	76.74

Table: 10×10 -fold cross-validation procedure: Average accuracies (on 60 data sets) with random noise to the features and the class variable

Outline

- Graphical Interpretation of Probabilistic Models
- Naïve Bayesian/Credal classifiers
- Decision Trees
- Bayesian Neural Networks

Artificial neural networks vs Bayesian Neural Networks



Graphical interpretation of (a) ANN and (b) BNN

Inference Problems [6]

Algorithm 1 Inference procedure for a BNN.

$$\text{Define } p(\theta|D) = \frac{p(D_Y|D_X, \theta) p(\theta)}{\int_{\theta} p(D_Y|D_X, \theta') p(\theta') d\theta'};$$

for $i = 0$ **to** N **do**

 Draw $\theta_i \sim p(\theta|D)$;

$\mathbf{y}_i = \Phi_{\theta_i}(\mathbf{x})$;

end for

return $Y = \{\mathbf{y}_i | i \in [0, N)\}$, $\Theta = \{\theta_i | i \in [0, N)\}$;

Inference Problems [6]

Algorithm 1 Inference procedure for a BNN.

$$\text{Define } p(\theta|D) = \frac{p(D_Y|D_X, \theta) p(\theta)}{\int_{\theta} p(D_Y|D_X, \theta') p(\theta') d\theta'};$$

for $i = 0$ **to** N **do**

 Draw $\theta_i \sim p(\theta|D)$;

$\mathbf{y}_i = \Phi_{\theta_i}(\mathbf{x})$;

end for

return $Y = \{\mathbf{y}_i | i \in [0, N)\}$, $\Theta = \{\theta_i | i \in [0, N)\}$;

- We need some way to aggregate the set of outputs

Aggregation procedures

Predict-then-aggregate (You can try it yourself):

- For each Monte Carlo sample, turn $\Theta_{\theta}(\mathbf{x})$ into a hard prediction y .
- Aggregate the set of hard predictions into the final hard prediction.
- You might want to try with MLE (Frequentist), DM (Bayesian), IDM (IP), etc.

Aggregation procedures

Predict-then-aggregate (You can try it yourself):

- For each Monte Carlo sample, turn $\Theta_{\theta}(\mathbf{x})$ into a hard prediction y .
- Aggregate the set of hard predictions into the final hard prediction.
- You might want to try with MLE (Frequentist), DM (Bayesian), IDM (IP), etc.

Aggregation-then-predict (See Lecture 6):

- For each Monte Carlo sample, compute a soft prediction $\Theta_{\theta}(\mathbf{x})$.
- Aggregate the set of soft predictions into either
 - the final hard prediction
 - or a credal set, from which IP decision rules can be applied.



ELSEVIER

Contents lists available at ScienceDirect

Software Impacts

journal homepage: www.journals.elsevier.com/software-impacts

Original software publication

BNNpriors: A library for Bayesian neural network inference with different prior distributions

Vincent Fortuin^{a,1,*}, Adrià Garriga-Alonso^{b,1}, Mark van der Wilk^{c,2}, Laurence Aitchison^{d,2}^a ETH Zürich, Zürich, Switzerland^b University of Cambridge, Cambridge, UK^c Imperial College London, London, UK^d University of Bristol, Bristol, UK

ARTICLE INFO

Keywords:

Machine learning
Bayesian neural networks
Prior distributions

ABSTRACT

Bayesian neural networks have shown great promise in many applications where calibrated uncertainty estimates are crucial and can often also lead to a higher predictive performance. However, it remains challenging to choose a good prior distribution over their weights. While isotropic Gaussian priors are often chosen in practice due to their simplicity, they do not reflect our true prior beliefs well and can lead to suboptimal performance. Our new library, *BNNpriors*, enables state-of-the-art Markov Chain Monte Carlo inference on Bayesian neural networks with a wide range of predefined priors, including heavy-tailed ones, hierarchical ones, and mixture priors. Moreover, it follows a modular approach that eases the design and implementation of new custom priors. It has facilitated foundational discoveries on the nature of the cold posterior effect in Bayesian neural networks and will hopefully catalyze future research as well as practical applications in this area.

Improve Trustworthiness in Deep Learning Models with Bayesian-Torch



What is Bayesian Deep Learning?

- › Uncertainty Estimation in Deep Learning

Creating the Foundation for Robust, Trustworthy AI

A Framework for Seamless Bayesian Model Development

- › How to Use Bayesian-Torch
- › Model Inferencing and Uncertainty Estimation

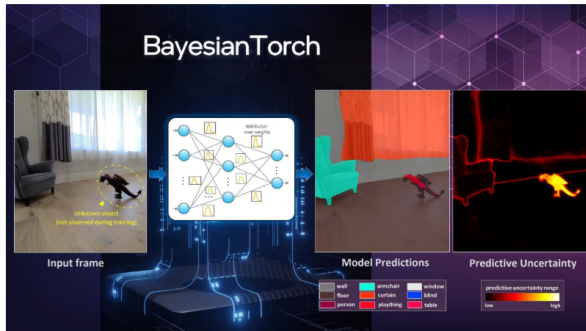
Use Case: Medical Application (Colorectal Histology Diagnosis)

Accounting for Distributional Shifts
Advancing Real-World Benchmarks

- › Developing Efficient Computing Systems for BDL Models

Get Involved

About the Author



References I

- [1] J. Abellán and S. Moral.
Building classification trees using the total uncertainty criterion.
International Journal of Intelligent Systems, 18(12):1215–1225, 2003.
- [2] A. Antonucci, G. Corani, and S. Bernaschina.
Active learning by the naive credal classifier.
In *Proceedings of the Sixth European Workshop on Probabilistic Graphical Models (PGM)*, pages 3–10, 2012.
- [3] G. Corani and M. Zaffalon.
Learning reliable classifiers from small or incomplete data sets: The naive credal classifier 2.
Journal of Machine Learning Research, 9(4), 2008.
- [4] V. G. Costa and C. E. Pedreira.
Recent advances in decision trees: An updated survey.
Artificial Intelligence Review, 56(5):4765–4800, 2023.
- [5] N. Friedman, D. Geiger, and M. Goldszmidt.
Bayesian network classifiers.
Machine learning, 29:131–163, 1997.
- [6] L. V. Jospin, H. Laga, F. Boussaid, W. Buntine, and M. Bennamoun.
Hands-on bayesian neural networksa tutorial for deep learning users.
IEEE Computational Intelligence Magazine, 17(2):29–48, 2022.
- [7] C. J. Mantas and J. Abellán.
Analysis and extension of decision trees based on imprecise probabilities: Application on noisy data.
Expert Systems with Applications, 41(5):2514–2525, 2014.

References II

- [8] C. J. Mantas and J. Abellan.
Credal-c4. 5: Decision tree based on imprecise probabilities to classify noisy data.
Expert Systems with Applications, 41(10):4625–4637, 2014.
- [9] V.-L. Nguyen, Y. Yang, and C. P. de Campos.
Probabilistic multi-dimensional classification.
In *Proceedings of the 39th Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 1522–1533, 2023.
- [10] 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)*, 2023.
- [11] C. Ortega Vázquez, S. vanden Broucke, and J. De Weerd.
Hellinger distance decision trees for pu learning in imbalanced data sets.
Machine Learning, pages 1–32, 2023.