

Exercises

Exercise 1 – IDM

- Consider the following data set (assume Boolean variables)
- | # | C | A |
|---|---|---|
| 1 | 0 | 1 |
| 2 | 0 | 1 |
| 3 | 0 | 1 |
- Use the IDM($s=2$) sequentially to
 - compute the lower and upper probabilities of ($C=1|A=1$) for the observations 0, 1, 2, 3
 - classify the same observations
 - Repeat the above using the precise Dirichlet model with Bayes-Laplace's uniform prior, and compare the results

Exercise 2 – NCC, complete data

- Consider the weather data set

Day	Outlook	Temperature	Humidity	Wind	Play/Tennis
d1	sunny	hot	high	strong	no
d2	sunny	hot	high	strong	no
d3	overcast	hot	high	weak	yes
d4	rainy	mild	high	weak	yes
d5	rainy	cool	normal	weak	yes
d6	rainy	cool	normal	strong	no
d7	overcast	cool	normal	strong	yes
d8	overcast	cool	high	strong	yes
d9	sunny	cool	normal	weak	yes
d10	rainy	mild	normal	weak	yes
d11	rainy	mild	normal	strong	yes
d12	overcast	mild	high	strong	yes
d13	overcast	hot	normal	weak	yes
d14	rainy	mild	high	strong	no

- Use the NCC($s=1$) to classify the instance (overcast, cool, high, strong)
- What if Outlook = sunny?

Exercise 3 – NCC, incomplete data

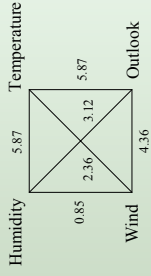
- Consider the incomplete weather data set

Day	Outlook	Temperature	Humidity	Wind	Play/Tennis
d1	sunny	hot	high	weak	no
d2	overcast	hot	high	strong	yes
d3	overcast	hot	high	*	yes
d4	rainy	mild	high	weak	yes
d5	rainy	cool	normal	weak	yes
d6	rainy	cool	normal	strong	no
d7	overcast	cool	normal	strong	yes
d8	sunny	cool	high	weak	no
d9	overcast	cool	normal	strong	yes
d10	rainy	mild	normal	weak	yes
d11	*	mild	normal	strong	yes
d12	overcast	mild	high	strong	yes
d13	overcast	hot	normal	weak	yes
d14	rainy	mild	high	strong	no

- Use the NCC($s=1$) to test whether or not 'no' dominates 'yes' with the instance (sunny, *, high, *)
 - Assume that 'Temperature' is missing at random, while 'Outlook' and 'Wind' are missing in non-ignorable way
 - Tip: minimum of function to optimize in $t_{yes} = 1$

Exercise 4 – TANC

- Consider the (complete) weather data set
- Use the TANC($s=1$) to classify the instance (sunny, mild, high, strong)
 - Assuming the following values of conditional mutual information



- and taking Humidity as root of the tree

5

Marco Zaffalon: Knowledge discovery from data sets under testable assumptions. School on Imprecise Probabilities, Lugano (CH), 27-31 July 2004.

Exercise 5 – mixed rules

- Sketch the theoretical derivation of the mixed rule in the learning setting
 - Stating new assumptions if needed

6

The multinomial case

- Traditional choice for $p(\theta)$ is Dirichlet(s, \mathbf{t}): $p(\theta) \propto \prod_{d \in \mathcal{D}} \theta_d^{s t_d - 1}$
 - s and $\mathbf{t} = (t_d)_{d \in \mathcal{D}}$ are hyperparameters
 - t_d = expectation of θ_d w.r.t. the prior
 - Informal view (prior is conjugate)
 - $s > 0$ as the number of additional (fractional) *virtual* observations
 - t_d as the proportion of observations equal to d in the virtual sample
 - $0 < t_d < 1, \sum_{d \in \mathcal{D}} t_d = 1$
 - e.g., uniform prior: $s = |\mathcal{D}|, t_d = 1/|\mathcal{D}|$; Perks prior: $s = 1, t_d = 1/|\mathcal{D}|$
- Posterior is Dirichlet($N+s, \mathbf{t}^*$): $p(\theta|\mathbf{d}) \propto \prod_{d \in \mathcal{D}} \theta_d^{n_d + s t_d - 1}$
 - $\mathbf{t}^* = (t_d^*)_{d \in \mathcal{D}}, t_d^* = (n_d + s t_d)/(N + s)$
 - t_d^* = expectation of θ_d w.r.t. the posterior
- The needed probability is then

$$p(c, x|\mathbf{d}) = E[\theta_{c,x}|\mathbf{d}] = \frac{n_{c,x} + s t_{c,x}}{N + s}$$

7

Appendixes

8

Marco Zaffalon: Knowledge discovery from data sets under testable assumptions. School on Imprecise Probabilities, Lugano (CH), 27-31 July 2004.

The imprecise Dirichlet model (IDM)

- An imprecise probability method for multinomial sampling

- Generalization of Bayesian approach

- Model prior ignorance by the set of all the Dirichlet(s, t) priors with s fixed
- s is a *degree of caution*

- From this, a set of posterior densities

- Lower and upper probabilities

$$p(c, x | d, s, t) = \frac{n_{c,x} + st_{c,x}}{N + s}$$

$$p(c, x | d, s) \in \left[\frac{n_{c,x}}{N + s}, \frac{n_{c,x} + s}{N + s} \right]$$

- Imprecision (upper – lower) is $\frac{s}{N+s}$

$$p(d) \propto \prod_{d \in \mathcal{D}} \theta_d^{t_d-1} \\ \sum_{d \in \mathcal{D}} t_d = 1 \\ 0 < t_d < 1 \quad \forall d \in \mathcal{D}$$

Credal classification

- There is generally no single *optimal* class with imprecise models
 - e.g., $p(c' | x) \in [0.5, 0.8]$, $p(c'' | x) \in [0.3, 0.6]$
 - No class *dominates* the other (they are *incomparable*); both are plausible options
 - Opposite approach: rather than looking for the optimal class, the focus is on **discarding the dominated classes**

- Credal dominance (strict preference)**

- $c' > c''$ iff $p(c' | d, x, s, t) > p(c'' | d, x, s, t)$ for all t in the IDM iff $\inf_t p(c' | d, s, t) / p(c'' | d, s, t) > 1$

- Classification**

- Compare all the classes and output the undominated ones
 - The output is a set of classes! The set shrinks with sample size; **reliability**

- Credal classification**

A credal classifier is a function $f: \mathcal{A}_1 \times \dots \times \mathcal{A}_m \rightarrow \wp(\mathcal{C})$

10

The naive credal classifier (NCC)

- Recall that $p(c, x | d, s, t) = \frac{n_{c,x} + st_{c,x}}{N + s} \prod_{j=1}^m \frac{n_{c,a_j} + st_{c,a_j}}{n_{c,x} + st_{c,x}}$
- The optimization problem for testing credal dominance is

$$\inf_t \frac{p(c', x | d, s, t)}{p(c'', x | d, s, t)} = \inf \left[\frac{(n_{c'} + st_{c'})}{(n_{c''} + st_{c''})} \prod_{j=1}^{m-1} \frac{n_{c', a_j} + st_{c', a_j}}{n_{c'', a_j} + st_{c'', a_j}} \right]$$

$$\text{s.t. } \sum_{c \in \mathcal{C}} t_c = 1$$

$$0 < t_c < 1 \quad \forall c \in \mathcal{C}$$

$$0 < t_{c,a_j} < t_c \quad \forall c \in \mathcal{C}, j \in \{1, \dots, m\}$$

$$= \inf \left[\frac{(n_{c'} + st_{c'})}{(n_{c''} + st_{c''})} \prod_{j=1}^{m-1} \frac{n_{c', a_j}}{n_{c'', a_j} + st_{c'', a_j}} \right]$$

$$\text{s.t. } t_{c'} + t_{c''} = 1$$

$$t_{c'}, t_{c''} > 0$$

- The latter is a convex optimization problem in a single variable!

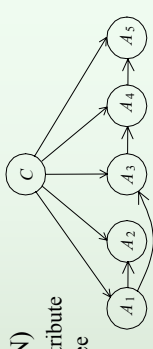
- Global* optimum found rapidly via numerical method, e.g., Newton-Raphson

11

Relaxing the independence assumption: the TAN model

- Tree augmented naive Bayes (TAN)**

- The dependence structure between attribute variables, given class variable, is a tree



- Learning problem**

- Learning the tree structure + probabilities
- Mix of *maximum likelihood* and Bayesian techniques

- Classification**

$$p(c, x | d) = p(c | d) \prod_{j=1}^m p(a_j | \pi_{A_j}, d)$$

12

Learning TAN structure

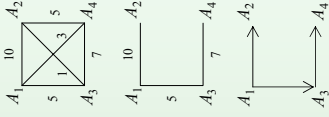
- Consider the *conditional mutual information* (MI)

$$MI(A_j, A_k | C) = \sum_{c \in \mathcal{C}} \theta_c \sum_{(a_j, a_k) \in A_j \times A_k} \theta_{a_j, a_k | c} \log \frac{\theta_{a_j, a_k | c}}{\theta_{a_j | c} \theta_{a_k | c}}$$

- Procedure

- Set up a fully connected undirected graph
 - Node \Leftrightarrow attribute variable
 - Weight each edge (A_j, A_k) by the *empirical* $MI(A_j, A_k | C)$
 - Empirical MI = replace chances with relative frequencies
- Compute the *maximum weight spanning tree*
- Arbitrarily choose root node

- The TAN T constructed this way maximizes the likelihood $p(\mathcal{d} | T)$



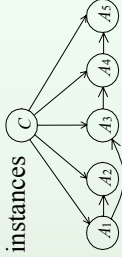
13

TANC classification

- $\underline{p}(c, a_1, \dots, a_m) = \underline{p}(c) \prod_{j=1}^m \underline{p}(a_j | \pi_{A_j})$
- Credal dominance test becomes (complete instance)

$$\frac{\underline{p}(c', a_1, \dots, a_m)}{\overline{p}(c'', a_1, \dots, a_m)} = \frac{\underline{p}(c')}{\overline{p}(c'')} \prod_{j=1}^m \frac{\underline{p}(a_j | \pi'_{A_j})}{\overline{p}(a_j | \pi''_{A_j})} > 1$$

- Linear time
- More involved with incomplete (ignorable) instances
 - $\underline{p}(c, a_5)$?
 - A_1, \dots, A_4 must be “marginalized out”
 - The expression to minimize involves a sum
- Exact minimization by propagating intervals over the tree
 - Linear time again



15

Tree-augmented naive credal classifier (TANC)

- Straightforward extension of TAN to credal classification
- Structure learning as in the precise case
- IDM-based learning of probabilities

- Replace $p(a_j | c, a_k, \mathbf{d}, s, \mathbf{t}) = \frac{n_{a_j, c, a_k} + s t_{a_j, c, a_k}}{n_{c, a_k} + s}$ with

$$p(a_j | c, a_k, \mathbf{d}, s) \in \left[\frac{n_{a_j, c, a_k}}{n_{c, a_k} + s}, \frac{n_{a_j, c, a_k} + s}{n_{c, a_k} + s} \right]$$

that is, $p(a_j | \pi_{A_j}) \in [\underline{p}(a_j | \pi_{A_j}), \overline{p}(a_j | \pi_{A_j})]$

- This way makes different credal sets be *separately specified* (Conditional) Credal sets in the same node and in different nodes

14

Learning – ideal variables

- Ideal learning data are regarded as an instance of the random matrix

$$\begin{bmatrix} C_1 & A_{11} & \dots & A_{1j} & \dots & A_{1k} \\ \vdots & \vdots & & \vdots & & \vdots \\ C_i & A_{i1} & \dots & A_{ij} & \dots & A_{ik} \\ \vdots & \vdots & & \vdots & & \vdots \\ C_N & A_{N1} & \dots & A_{Nj} & \dots & A_{Nk} \end{bmatrix} = \begin{bmatrix} C_1 & X_1 \\ \vdots & \vdots \\ C_i & X_i \\ \vdots & \vdots \\ C_N & X_N \end{bmatrix} = \begin{bmatrix} D_1 \\ \vdots \\ D_i \\ \vdots \\ D_N \end{bmatrix} = \mathcal{D}$$

- C_i 's are *class variables*, with values in \mathcal{C}
- A_{ij} 's are *attribute variables*, with values in \mathcal{A}_j for each (ij)
- Each row D_i in the matrix is called *unit*, with values in $\mathcal{D} = \mathcal{C} \times \mathcal{X} = \mathcal{C} \times \mathcal{A}_1 \times \dots \times \mathcal{A}_m$

16

Learning – actual variables

- Consider one more random variable
- \mathbf{O} is the *actual observation* of \mathbf{D}
 - \mathbf{D} takes values \mathbf{d} from \mathcal{D}^N
 - \mathbf{O} takes values \mathbf{o} from $\varphi(\mathcal{D}^N)$
 - Rather than observing \mathbf{d} , you observe a set \mathbf{o} that contains it
 - This is called *coarsening*, i.e., looking at \mathbf{d} with different levels of detail
 - Note that \mathbf{o} is simply a symbol, not a set, when regarded as the value of \mathbf{O}

17

Learning – problem

- Formulation of the learning problem
 - Using observed data \mathbf{o} to update beliefs about a function $f: \Theta \rightarrow \mathbb{R}$
 - e.g., θ real number, $f(\theta) = \theta$
 - In the precise framework, one would compute $E(f|\mathbf{o})$

$$\begin{aligned}
 E(f|\mathbf{o}) &= \frac{\int f(\theta) p(\theta, \mathbf{o}) d\theta}{p(\mathbf{o})} \\
 &= \frac{\int f(\theta) \sum_{\mathbf{d} \in \mathbf{o}} p(\theta) p(\mathbf{d}|\theta) p(\mathbf{o}|\mathbf{d}) d\theta}{\sum_{\mathbf{d} \in \mathbf{o}} p(\mathbf{o}|\mathbf{d}) p(\mathbf{d})}
 \end{aligned}$$

Using Factorization, Accuracy and Positivity

- What about $p(\mathbf{o}|\mathbf{d})$ and $p(\theta)$?
 - Imprecise knowledge

19

Learning – assumptions

- Overall model = joint density $p(T, \mathbf{D}, \mathbf{O})$
 - T is a random parameter with values $\theta \in \Theta$
- Factorization
 - $p(T, \mathbf{D}, \mathbf{O}) = p(T) p(\mathbf{D}|T) p(\mathbf{O}|\mathbf{D})$
 - $p(T)$ is an imprecise prior for T , i.e., it belongs to a certain non-empty set
 - The IM depends only on \mathbf{D}
- Accuracy (of mechanism)
 - $p(\mathbf{o}|\mathbf{d}) = 0$ if $\mathbf{d} \notin \mathbf{o}$
- Connection between complete and incomplete observations
- Positivity
 - Ideal data: $p(\mathbf{D}) > 0$
 - Actual observation: $p(\mathbf{o}) > 0$
 - There exists \mathbf{d} s.t. $p(\mathbf{o}|\mathbf{d}) > 0$

18

Learning – express ignorance about the IM

- Focus on $p(\mathbf{o}|\mathbf{d})$, $\mathbf{d} \in \mathbf{o}$
 - Call $p(\mathbf{o}|\mathbf{D})$ the vector with elements $p(\mathbf{o}|\mathbf{d})$, $\mathbf{d} \in \mathbf{o}$
- We can only constrain the set of possible vectors $p(\mathbf{o}|\mathbf{D})$
 - $\sum_{\mathbf{d} \in \mathbf{o}} p(\mathbf{o}|\mathbf{d}) > 0$
 - $0 \leq p(\mathbf{o}|\mathbf{d}) \leq 1$, $\mathbf{d} \in \mathbf{o}$
- The restriction on the sum is due to Accuracy
- The inequality following the sum is due to Positivity
- The inequalities define an open linear set called $\mathcal{P}(\mathbf{o}|\mathbf{D})$
- Call $\mathcal{P}_\varepsilon(\mathbf{o}|\mathbf{D})$ the approximating closed set

$$\begin{aligned}
 \sum_{\mathbf{d} \in \mathbf{o}} p(\mathbf{o}|\mathbf{d}) &\geq \varepsilon \\
 0 &\leq p(\mathbf{o}|\mathbf{d}) \leq 1, \mathbf{d} \in \mathbf{o}
 \end{aligned}$$

20

Learning – goal

- Goal: $\underline{E}(f|\mathbf{o}) = \inf_{p(T) \in \mathcal{P}(T)} \inf_{p(\mathbf{o}|D) \in \mathcal{P}(\mathbf{o}|D)} E(f|\mathbf{o})$
- Focus on $\min_{p(\mathbf{o}|D) \in \mathcal{P}_e(\mathbf{o}|D)} E(f|\mathbf{o})$
 - i.e., $\min_{p(\mathbf{o}|D) \in \mathcal{P}_e(\mathbf{o}|D)} \frac{\sum_{d \in \mathbf{o}} p(\mathbf{o}|d) \int f(\theta) p(\theta) p(d|\theta) d\theta}{\sum_{d \in \mathbf{o}} p(\mathbf{o}|d) p(d)}$
 - Objective function is ratio of linear functions
- Fractional programming theorem
 - Consider $\min_{x \in S} \frac{q(x)}{r(x)}$ where S = compact subset of \mathbb{R}^v
 - q, r continuous and r positive on S
 - Define $h : \mathbb{R} \rightarrow \mathbb{R}$, $h(\mu) = \min_{x \in S} [q(x) - \mu r(x)]$
 - Then $\mu^* = \arg\min_{x \in S} \frac{q(x)}{r(x)} \iff h(\mu^*) = 0$
 - h has a single zero

21

Obtaining stronger conclusions: the mixed rules

- CAR might hold in some cases
 - This should not be neglected in order to strengthen conclusions
- Some attribute variables are CAR, the others non-ignorable
- Mixed rules
 - Learning: $\underline{E}(f|\mathbf{o}, \hat{\mathbf{o}}) = \inf_{p(T) \in \mathcal{P}(T)} \min_{d \in \hat{\mathbf{o}}} E(f|d, \hat{d} \in \hat{\mathbf{o}})$
 - Classification: $\underline{E}_C(g|\mathbf{o}^+, \hat{\mathbf{o}}^+) = \inf_{p(T) \in \mathcal{P}(T)} \min_{d^- \in \hat{\mathbf{o}}^-} E_C(g|d^-, \hat{d}^- \in \hat{\mathbf{o}}^-)$

23

Learning – solution: conservative learning rule

- $h(\mu) = \min_{p(\mathbf{o}|D) \in \mathcal{P}_e(\mathbf{o}|D)} \left[\sum_{d \in \mathbf{o}} p(\mathbf{o}|d) \int f(\theta) p(\theta) p(d|\theta) d\theta - \mu \sum_{d \in \mathbf{o}} p(\mathbf{o}|d) p(d) \right]$
 - $= \min_{p(\mathbf{o}|D) \in \mathcal{P}_e(\mathbf{o}|D)} \sum_{d \in \mathbf{o}} p(\mathbf{o}|d) p(d) \left[\int f(\theta) p(\theta) p(d|\theta) d\theta - \mu \right]$
 - $= \min_{p(\mathbf{o}|D) \in \mathcal{P}_e(\mathbf{o}|D)} \sum_{d \in \mathbf{o}} p(\mathbf{o}|d) p(d) [E(f|d) - \mu]$
 - $\mu^* = \min_{d \in \mathbf{o}} E(f|d) \implies h(\mu^*) = 0$
 - $\min_{p(\mathbf{o}|D) \in \mathcal{P}_e(\mathbf{o}|D)} E(f|\mathbf{o}) = \min_{d \in \mathbf{o}} E(f|d)$
- ↓
- $$\underline{E}(f|\mathbf{o}) = \inf_{p(T) \in \mathcal{P}(T)} \min_{d \in \mathbf{o}} E(f|d)$$

22

Inference of the NCC from incomplete data

- Test of credal dominance with an incomplete learning sample

$$\min_{d \in \mathbf{o}} \inf_t \frac{p(c', x|d, s, t)}{p(c'', x|d, s, t)} = \inf \left[\left(\frac{n_{c'} + st_{c'}}{n_{c''} + st_{c''}} \right)^{m-1} \prod_{j=1}^m \frac{n_{c', a_j}}{n_{c'', a_j} + st_{c''}} \right]$$

s.t. $t_{c'} + t_{c''} = 1$
 $t_{c'}, t_{c''} > 0$
- Same complexity as with complete data

24