

# Výpočet marginálních podmíněných pravděpodobností v bayesovské síti

- Úmluva: Zajímáme se pouze o bayesovské sítě, jejichž graf je spojitý. Jinak uvažujeme každou komponentu zvlášť.
- Notace

## Definition (Pojmy v orientovaném grafu)

- $\text{parents}(\text{'either'}, \text{chestdag})$  - odkud vede hrana do 'either'
- $\text{children}(\text{'tub'}, \text{chestdag})$  - kam vede hrana z 'tub'
- $\text{ancestralSet}(\text{'either'}, \text{chestdag})$  - ancestrální množina, předkové

# Moralizace, sousedé, klika, simplicialní uzly

## Definition (Moralizace)

**Moralizací grafu bayesovské sítě** rozumíme následující dva kroky:

- **Ožnit** (spojit hranou) **rodiče společných dětí**. (Neboli: spojit hranou každé dva vrcholy, které se společně vyskytují v některé tabulce dané bayesovské sítě.)
- **Zapomenout orientaci hran**, tj. vytvořit neorientovaný graf se stejnými hranami.

## Definition (Pojmy v neorientovaném grafu)

Nechť  $X$  je uzel neorientovaného grafu. Pak mají označení následující význam:

|       |  |
|-------|--|
| $N_X$ | sousedé $X$                            |
| $F_X$ | sousedé $X$ včetně $X$ samého (family) |

- Podgraf je **úplný**, pokud jsou každé dva jeho uzly spojeny hranou.
- Podgraf je **klika**, pokud je maximální úplný podgraf.
- Vrchol  $X$  je **simplicialní**, pokud je  $N_X$  úplný podgraf.  
Ekvivalentně: vrchol  $X$  je simplicialní, je-li  $F_X$  klika.

## Definition (Graf domén)

**Graf domén** v konkrétním kroku eliminace proměnných je takový graf, kde

- uzly jsou právě všechny dosud neeliminované proměnné,
  - dva uzly jsou spojeny hranou právě když se odpovídající proměnné vyskytují v aspoň jedné tabulce zároveň.
- 
- **Graf domén** je na počátku moralizovaná bayesovská síť.
  - Po eliminaci každé proměnné odstraníme jí příslušející uzel a spojíme všechny jeho sousedy (nově přidáné hrany se nazývají **doplňené**, fill-in).
  - Naším cílem jsou co nejmenší domény, tj. co nejméně doplněných hran.

# Algoritmus eliminace proměnných

- INIT** Do seznamu  $\Phi_1$  dáme všechny tabulky  $P(A_i, e | pa(A_i))$ , v každé tabulce odstraním "řádky" nekonzistentní s evidencí, tj. s nulovou pravděpodobností. Tj. předem evidencí vynásobíme a marginalizujeme přes proměnné s evidencí.
- ELIM** Postupně budeme eliminovat (následujícím algoritmem) všechny proměnné bez evidence, které nás nezajímají (dostaneme  $P(A, e)$ ).
- NORM** Nakonec eliminujeme i zbývající proměnné bez evidence, čímž spočteme normalizační konstantu  $\alpha = P(e)$ ; touto konstantou vydělíme tabulku z předchozího kroku a dostaneme podmíněnou pravděpodobnost  $P(A|e)$ .

**Eliminace proměnné  $X$  v kroku  $i$  znamená:**

- 1 Vyber z  $\Phi_i$  všechny tabulky, které mají v doméně  $X$ , dej je do  $\Phi_X$ .
- 2 Spočti  $\phi = \sum_X \prod_{T \in \Phi_X} T$
- 3 Nové  $\Phi_{i+1}$  se rovná:  $\Phi_i \setminus \Phi_X \cup \{\phi\}$

- V jakém pořadí eliminovat? Špatné pořadí vede ke zbytečně velkým tabulkám  $\phi$ .
- špatně** Nejdřív nepozorované rodiče (např. v 'minách').
- oprávně** Nejdřív **barren** tj. uzly bez dětí a bez evidence.
- obecně** Nejdřív **simpliciální uzly**.

# Perfektní eliminační posloupnost

## Definition (Perfektní eliminační posloupnost)

**Perfektní eliminační posloupnost** je taková posloupnost eliminace všech proměnných bayesovské sítě, která nevynucuje žádné doplněné hrany v grafu domén.

## Lemma (6 Simpliciální uzel smí na začátek PEP)

*Nechť je  $X_1, \dots, X_k$  perfektní eliminační posloupnost (PEP),  $X_j$  uzel, jehož každý dva sousedi jsou spojeni hranou. Pak je posloupnost  $X_j, X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_k$  také perfektní.*

Eliminace  $X_j$  nepřidá žádnou hranu, tj. nikomu nepřidá souseda, a při eliminaci se každý stará jen o své sousedy, tj. se na eliminaci ostatních nic nezmění (ledaže by nemuseli přidávat hranu do  $X_j$ , ale i původní posloupnost byla perfektní).

## Definition (Množina maximálních domén)

**Množina maximálních domén** je množina všech domén tabulek, vzniklých během výpočtu, ze které vyřadíme ty domény, které jsou vlastní podmnožinou jiného prvku této množiny.

## Lemma

*Všechny perfektní posloupnosti vytvářejí stejnou množinu maximálních domén, a to **množinu klik** moralizovaného grafu.*

- Kliky tam musí být, neboť tabulka jejich domény vznikne při eliminaci první proměnné z kliky.
- Nic většího tam nemůže být, neboť by to způsobilo doplněné hranu.

# Triangulované grafy

- **POZOR, něco jiného než triangulovaný planární graf!!!**
- Pozn: Pro některé moralizované grafy neexistuje perfektní posloupnost.

## Definition (triangulovaný graf)

Graf je **triangulovaný**, pokud pokud pro něj existuje perfektní eliminační posloupnost.

## Lemma (Alternativní definice triang. grafu)

*Graf je triangulovaný, pokud každý jeho cyklus délky větší než tři má aspoň jednu tětivu.*

- Je 'chestdag' triangulovaný?

## Lemma

Nechť je  $G$  triangulovaný graf a  $X$  simplicialní uzel. Pak graf  $G^1$  získaný eliminací  $X$  z  $G$  je také triangulovaný.

Důsledek lemmatu 6.

## Theorem

- *Triangulovaný graf s aspoň dvěma vrcholy má aspoň dva simplicialní uzly.*
- *Navíc: pokud není úplný, tak má aspoň dva simplicialní uzly, které nejsou spojeny hranou.*

Důkaz indukcí podle počtu vrcholů.

- Pro tři vrcholy platí.
- Pro více: První uzel perfektní posloupnosti je simplicialní, vzniklý graf je triangulovaný.
- Pokud vzniklý graf není úplný,
  - z indukčního předpokladu má aspoň dva nesousední simplicialní uzly,
  - sousedé eliminovaného uzlu byly propojeni (simplicialní uzel)
  - proto aspoň jeden z nových simplicialních uzlů nesousedil s eliminovaným uzlem.





## Lemma

*Pro každý vrchol  $A$  triangulovaného grafu existuje perfektní posloupnost, kde je  $A$  poslední prvek.*

Důkaz: Vždy eliminuj simplicialní uzel různý od  $A$ .

## Theorem

*Neorientovaný graf je triangulovaný právě když můžeme eliminovat všechny uzly tak, že vždy eliminujeme simplicialní uzel.*

- ⇒ Je-li graf triangulovaný, eliminací simplicialního uzlu vznikne opět triangulovaný graf a můžeme pokračovat v eliminaci simplicialních uzlů.
- ⇐ Eliminací simplicialních uzlů tvoříme perfektní posloupnost, tj. graf je triangulovaný.

# Stromy spojení (Join trees)

## Definition (Strom spojení)

Mějme množinu klik neorientovaného grafu  $G$ , kliky jsou organizovány do stromu  $T$ .  $T$  je **strom spojení**, pokud pro každé dva vrcholy  $V, W \in T$  všechny uzly na cestě z  $V$  do  $W$  obsahují průnik  $W \cap V$ . Průnik dvou sousedních uzlů nazveme **separátor** těchto uzlů, separátorem  $V$  a  $W$  je  $S_{V,W} = V \cap W$ .

## Theorem

Pokud kliky grafu  $G$  lze organizovat do stromu spojení, pak je  $G$  triangulovaný.

- Vezměme list stromu spojení  $V$ , který sousedí jen s  $W$ .
- $V$  průnik libovolný uzel je částí  $W$ , proto  $V$  obsahuje uzel, který není v žádné jiné klice. Ten eliminujeme.
- Pokud byl poslední z  $W \setminus V$ , odstraníme uzel  $V$  a dostaneme zase strom spojení, který má list, pokračujeme prvním bodem.

## Theorem

Pokud je  $G$  triangulovaný, pak kliky grafu  $G$  lze organizovat do stromu spojení.

Důkaz indukcí, pro grafy s jedním vrcholem platí.

- Eliminuji simplicialní uzel  $X$ , jeho rodina  $F_X$  je klika (označíme jí  $C$ ).
- Pro vzniklý graf  $G' = G \setminus \{X\}$  najdu strom spojení  $T'$  dle indukčního předpokladu.
- Pokud je  $C \setminus \{X\}$  klika  $G'$ , k uzlu odpovídajícímu této klice v  $T'$  přidám 'popisku'  $X$  a mám strom spojení grafu  $G$ .
- Pokud  $C \setminus \{X\}$  není klika  $G'$ :
  - musí být částí kliky  $C_?$  grafu  $G'$ ,
  - Ke stromu  $T'$  přidáme uzel  $C$  a připojíme separátorem  $C \setminus \{X\}$  k uzlu kliky  $C_?$ . Vzniklý strom je strom spojení pro  $G$ .



# Alternativní konstrukce

- Najdi kliky.
- Vytvoř graf, uzly=kliky, hrany váhy počtu veličin v průniku.
- Najdi kostru (spanning tree) nejvyšší váhy (Prim's or Kruskal's algorithm).

Tato kostra je strom spojení, protože

- Je-li proměnná  $X$  v  $j$  klikách, může být maximálně v  $j - 1$  separátorech stromu spojení.
- Číslo  $j - 1$  dosáhneme jen v případě, že všechny kliky obsahující  $X$  budou spojeny separátorem obsahujícím  $X$ .
- Proto má strom spojení nejvyšší možný součet velikostí separátorů přes všechny kostry grafu.

## Strom spojení

Používám termín **strom spojení** ve třech významech:

- viz definice výše, strom klik splňující vlastnost průniků
- strom dle definice výše, kde jsou navíc hrany označeny separátory
- strom dle definice výše, kde je navíc v každé klice "schránka" na seznam pravděpodobnostních tabulek a v každém separátoru jsou dvě schránky na zprávy – tabulky – jdoucí jednotlivými směry. Tomuto se říká **junction tree**.

# Strom spojení reprezentující bayesovskou síť

Mějme bayesovskou síť s množinou pravděpodobnostních tabulek  $\Phi$  a evidenci  $e$ . Necht množina tabulek  $\Phi_e$  vznikne z  $\Phi$  vložením evidence  $e$  do příslušných tabulek, tj. "vyříznutím" konkrétních "řádků" v pravděpodobnostních tabulkách.

**Strom spojení reprezentuje bayesovskou síť s evidencí  $e$ ,** pokud každou tabulku  $\phi \in \Phi_e$  přiřadíme do schránky některé z klik  $C_i$  takových, že  $\text{dom}(\phi) \subseteq C_i$ .

- Pozn: pokud strom spojení vznikl z moralizovaného a triangularizovaného grafu bayesovské sítě, tak takové klika vždy existuje.
- Pokud moralizovaný graf není triangulovaný, doplníme ho hranami na triangulovaný a z něj vytvoříme strom spojení.
- Pokud některý uzel stromu spojení nemá žádnou tabulku, přiřadíme tabulku dávající identicky 1 na doméně dané kliky.

# Propagace ve stromu spojení

- Propagace (výpočet) ve stromu spojení spočívá v posílání zpráv, kterými se postupně plní schránky separátorů.
- Každý uzel (klika) posílá v každém směru právě jednu zprávu.
- Uzel (klika) může poslat zprávu v daném směru, pokud už ze všech ostatních směrů zprávy dostala.
- Protože se jedná o strom, vždycky někdo může poslat zprávu, nebo jsou již všechny schránky plné.

## Poslání zprávy

Uvažujme kliku  $C$  se sousedními separátory  $S_1, \dots, S_k$ , směr separátoru  $S_1$  (bez újmy na obecnosti). **Poslat zprávu** z  $C$  do  $S_1$  znamená zapsat do odchozí schránky  $S_1$  tabulku, která vznikne součinem příchozích zpráv v separátorech  $S_2, \dots, S_k$  a tabulek obsažených v  $C$ . Tento součin marginalizujeme přes všechny veličiny  $C \setminus S_1$  a výsledek zapíšeme do  $S_1$ .



## Theorem

Nechť strom spojení reprezentuje bayesovskou síť a evidenci  $e$ , všechny schánky byly naplněny. Potom:

- Nechť  $V$  je klika obsahující tabulky  $\Phi_V$  a  $k$  ní směřující separátory  $S_1, \dots, S_k$  obsahují zprávy  $\phi_1, \dots, \phi_k$ .

$$P(V, e) = \prod_{\phi \in \Phi_V} \phi \cdot \prod_{i=1}^k \phi_i$$

- Nechť  $S$  je separátor se zprávami  $\phi_1, \phi_2$ .

$$P(S, e) = \phi_1 \cdot \phi_2$$

Zprávy směřující do  $V$  odpovídají perfektní elim. posl., která má  $V$  na svém konci. Pro separátor, odchozí zpráva vznikla marginalizací z  $V$ , jen tam nebyla započtena zpráva přicházející z tohoto směru.

$$\begin{aligned} P(S_1, e) &= \sum_{V \setminus S_1} P(V, e) = \sum_{V \setminus S_1} \left( \prod_{\phi \in \Phi_V} \phi \cdot \prod_{i=1}^k \phi_i \right) \\ &= \sum_{V \setminus S_1} \left( \prod_{\phi \in \Phi_V} \phi \cdot \prod_{i=2}^k \phi_i \cdot \phi_1 \right) = \left( \sum_{V \setminus S_1} \prod_{\phi \in \Phi_V} \phi \cdot \prod_{i=2}^k \phi_i \right) \cdot \phi_1 \end{aligned}$$

což je odchozí krát přichodí zpráva. Poslední řádek plyne z toho, že  $\text{dom}(\phi_1) = S_1$ .

# Výpočet pomocí stromu spojení (shrnutí)

- BN moralizujeme
- doplníme hrany na triangulovaný graf
- vytvoříme strom spojení
- naplníme tabulkami
- vypočteme posíláním zpráv
- pravděpodobnost na veličině  $A$  zjistíme tak, že najdeme libovolnou kliku  $C$  obsahující  $A$  a marginalizujeme, tj.  $P(A, e) = \sum_{C \setminus \{A\}} P(C, e)$
- pokud nás zajímá sdružená distribuce na množině, která není částí žádné kliky, musíme použít Eliminaci proměnných.  
Nebo předem zajistit výskyt v jedné klice:  
`m3=compile(grain(plist),root=c('lung','bronc','tub'), propagate=TRUE).`

# Separace, rozložitelnost

## Definition (Separace)

Buď  $G$  neorientovaný graf nad  $V$ . Řekneme  $C$  separuje  $A$  a  $B$ , psáno  $A \perp\!\!\!\perp_G B | C$  v  $G$ , pokud každá cesta z  $A$  do  $B$  v  $G$  vede přes  $C$ .

## Definition (Rozklad)

Buď  $G$  neorientovaný graf nad  $V$ . Jestliže  $S, T \subseteq V$  jsou takové, že

- 1  $S \cup T = V$
- 2  $S \cap T$  je úplná množina v  $G$  a
- 3  $S \setminus T \perp\!\!\!\perp_G T \setminus S | S \cap T$ ,

Pak dvojici indukovaných podgrafů  $G_S, G_T$  nazveme **rozkladem**  $G$ . Tento rozklad nazveme netriviální, jestli že  $S \setminus T \neq \emptyset \neq T \setminus S$ .

Neorientovaný graf  $G$  nazveme **rozložitelný**

- buď  $G$  je úplný
- nebo existuje netriviální rozklad  $G_S, G_T$  grafu  $G$  takový, že  $G_S$  a  $G_T$  jsou rozložitelné.

# Markovské vlastnosti

## Definition

Markovská vlastnost, Globální, lokální, párová Buď  $G$  neorientovaný graf nad  $V$ . Pravděpodobnostní míra  $P$  nad  $V$  je **(globálně) markovská** vzhledem k  $G$ , jestliže:

$$\forall (A, B \in V, C \subseteq V) \quad A \perp\!\!\!\perp_G B | C \Rightarrow A \perp\!\!\!\perp B | C \text{ v } P.$$

- Míra je lokálně markovská, pokud  $\forall A \in V \quad A \perp\!\!\!\perp V \setminus Fa_A | N_A [P]$
- Míra je párově markovská, pokud  $\forall A, B \in V, A \neq B$ , nespojené hranou,  $A \perp\!\!\!\perp B | V \setminus \{A, B\} [P]$

Pro striktně pozitivní míry (bez nul) jsou všechny tři vlastnosti ekvivalentní. Jinak protipříklady s mírou v bodech  $(0, 0, 0)$  a  $(1, 1, 1)$  viz Studený.

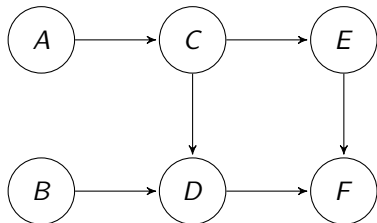
## Definition (d-separace)

Dvě veličiny  $A, B \in V$  bayesovské sítě  $G = (V, E)$  jsou **d-separované**  $A \perp\!\!\!\perp_d B | C$  množinou  $C \subseteq V \setminus \{A, B\}$  právě když pro každou (neorientovanou) cestu z  $A$  do  $B$  platí aspoň jedno z následujících:

- cesta obsahuje uzel  $Blocking \in C$  a hrany se v  $Blocking$  **nesetkávají** 'head-to-head',
- cesta obsahuje uzel  $Blocking$  kde se hrany **setkávají** 'head-to-head' a ani on ani nikdo z jeho následníků není v  $C$ ,  $\{Blocking, succ(Blocking)\} \cap C = \emptyset$ .

## Theorem (d-separace)

*Pokud jsou  $A, B$  d-separované dáno  $C$  ( $A \perp\!\!\!\perp_d B | C$ ) v BN  $B$ , pak jsou i podmíněně nezávislé ( $A \perp\!\!\!\perp B | C$ ).*



Platí?

- $E \perp\!\!\!\perp_d B$  ano
- $E \perp\!\!\!\perp_d D$  ne
- $E \perp\!\!\!\perp_d D | A$  ne
- $E \perp\!\!\!\perp_d D | C$  ano
- $E \perp\!\!\!\perp_d D | \{C, F\}$  ne
- $E \perp\!\!\!\perp_d B | F$  ne

## Podmíněná nezávislost

Tabulka ukazuje zadané hodnoty  $P(A, B, C)$ . Pro která  $x, y, v, w$  platí podmíněná nezávislost  $A \perp\!\!\!\perp B | C$ ?

|    | c1 |     |     | c2  |
|----|----|-----|-----|-----|
|    | b1 | b2  | b1  | b2  |
| a1 | x  | 0,2 | v   | w   |
| a2 | y  | 0,1 | 0,1 | 0,1 |

$$\frac{P(a_1, b_1, c_1)}{P(a_2, b_1, c_1)} = \frac{P(a_1|c_1) \cdot P(b_1|c_1) \cdot P(c_1)}{P(a_2|c_1) \cdot P(b_1|c_1) \cdot P(c_1)} = \frac{P(a_1, b_2, c_1)}{P(a_2, b_2, c_1)}$$

tedy  $x = 2 \cdot y$

obdobně  $v = w$ .

Navíc celkový součet musí být 1, tedy:

$$0,5 + 3y + 2v = 1$$

$$v = 0,25 - 1,5y$$

## Přibližný výpočet bayesovské sítě

- Základní myšlenkou je vygenerovat data dle zadaných podmíněných pravděpodobností a z nich spočítat pravděpodobnosti, které nás zajímají.
- Přesnost výpočtu samozřejmě závisí na počtu vygenerovaných vzorků.
- Metody generující náhodné vzorky se nazývají metody **Monte Carlo**.
- Základem je generátor náhodného výsledku podle zadané pravděpodobnosti, např.  $\langle \frac{1}{4}, \frac{1}{2}, \frac{1}{4} \rangle$ .

## Přímé vzorkování bez evidence

- Uspořádáme vrcholy BN tak, aby každá hrana začínala v uzlu menšího čísla než končí.
- Vytvoříme  $N$  vzorků, každý následovně
  - Pro první uzel  $A_1$  vygenerujeme náhodně výsledek  $a_1$  podle  $P(A_1)$ .
  - Pro druhý uzel  $A_2$  vygenerujeme náhodně výsledek  $a_2$  podle  $P(A_2|A_1 = a_1)$  (je-li hrana, jinak nepodmíněně)
  - Pro  $n$ -tý uzel vygenerujeme výsledek podle  $P(A_n|pa(A_n))$ , na rodičích už známe konkrétní hodnoty.
- Z  $N$  vzorků spočteme pravděpodobnost jevu, který nás zajímá. Pro  $N$  jdoucích k nekonečnu podíl výskytu jevu konverguje k správné pravděpodobnosti.



## Přímé vzorkování s evidencí $e$ (rejection sampling)

- $N(e)$  značí počet vzorků konzistentních s evidencí  $e$ , tj. nabývajících na příslušných veličinách správné hodnoty.
- Vzorky tvoříme úplně stejně, jako dříve, jen ty, co nejsou konzistentní s  $e$  vyškneme, tj.  $\hat{P}(X|e) = \frac{N(X,e)}{N(e)}$
- Problém je v tom, že je-li  $P(e)$  malé, tak většinu vzorků zahazujeme.

# Vážení věrohodností (Likelihood weighting)

- Generuje jen vzorky konzistentní s  $e$ .
- Váhy vzorků jsou různé, podle  $P(e|\text{vzorek})$  (což je věrohodnost  $L(\text{vzorek}|e)$ , odtud likelihood weighting).

Algoritmus **vytvoření váženého vzorku pro**  $(bn, e)$

$w = 1$

v pořadí topologického uspořádání  $bn$ , for  $i = 1$  to  $n$

if  $A_i$  má evidenci  $a_i$  v  $e$

$w = w \cdot P(A_i = a_i | pa(A_i))$

else

$a_i$  vyber podle rozložení  $P(A_i = a_i | pa(A_i))$

return  $(w, \langle a_1, \dots, a_n \rangle)$