

## Bayesovská sieť

- Orientovaný acyklický graf
- Vrcholy: náhodné premenné  $X_1 \dots X_n$
- Nech  $Z_i$  sú predchodcovia premennej  $X_i$
- Pre každé  $X_i$  tabuľka  $\Pr(X_i|Z_i)$
- Sieť definuje celkové rozdelenie

$$\Pr(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n \underbrace{\Pr(X_i = x_i | Z_i = x_{Z_i})}_{\Theta_{i, x_i, x_{Z_i}}}$$

## Trénovanie s úplnými dátami

$$\mathbf{x}_1 = (x_{1,1}, x_{1,2}, x_{1,3}, \dots, x_{1,n})$$

$$\mathbf{x}_2 = (x_{2,1}, x_{2,2}, x_{2,3}, \dots, x_{2,n})$$

...

$$\mathbf{x}_t = (x_{t,1}, x_{t,2}, x_{t,3}, \dots, x_{t,n})$$

Hľadáme najvierohodnejšie parametre:

$$\Theta^* = \arg \max_{\Theta} \sum_{i=1}^t \log \Pr(x_{i,1}, \dots, x_{i,n} | \Theta)$$

Nastavíme na pozorované frekvencie:  $\Theta_{j,u,v}^* = \frac{\#_{j,u,v}}{n_{j,v}}$

- $\#_{j,u,v}$  je počet vzoriek, kde  $X_j = u$  a  $Z_j = v$
- $n_{j,v}$  je počet vzoriek, kde  $Z_j = v$

## Trénovanie s neúplnými dátami

- $\chi^{(i)}$  - pozorované dáta vo vzorke  $i$
- $X^{(i)}$  - pozorované premenné
- $y^{(i)}$  - chýbajúce dáta vo vzorke  $i$
- $Y^{(i)}$  - chýbajúce premenné
- Cieľ: maximalizovať vierohodnosť

$$\Theta^* = \arg \max_{\Theta} \prod_{i=1}^t \Pr(X^{(i)} = \chi^{(i)} | \Theta)$$

## EM algoritmus (Expectation maximization)

- Zvoľ počiatkové  $\Theta^{(0)}$
- Iteruj kým nenastane konvergencia ( $k = 0, 1, \dots$ )
  - E-krok: spočítaj pre každé  $i, y^{(i)}$ :

$$Q_i^{(k)}(y^{(i)}) = \Pr(Y^{(i)} = y^{(i)} | X^{(i)} = x^{(i)}, \Theta^{(k)})$$

- M-krok: spočítaj

$$\Theta^{(k+1)} = \arg \max_{\Theta} \sum_{i=1}^t \sum_{y^{(i)}} Q_i^{(k)}(y^{(i)}) \log \Pr(X^{(i)} = x^{(i)}, Y^{(i)} = y^{(i)} | \Theta)$$

Maximalizujeme vierohodnosť dát doplnených vš. spôsobmi o chýbajúce hodnoty, váhované cez  $Q^{(k)}$

## M-krok EM algoritmu

Spočítaj

$$\Theta^{(k+1)} = \arg \max_{\Theta} \sum_{i=1}^t \sum_{\mathbf{y}^{(i)}} Q_i^{(k)}(\mathbf{y}^{(i)}) \log \Pr(X^{(i)} = \mathbf{x}^{(i)}, Y^{(i)} = \mathbf{y}^{(i)} | \Theta)$$

Maximalizujeme vierohodnosť dát doplnených vš. spôsobmi o chýbajúce hodnoty, váhované cez  $Q^{(k)}$

Maximum pre pozorované frekvencie s váhovanými dátami:

$$\Theta_{j,u,v}^{(k+1)} = \frac{\#_{j,u,v}}{n_{j,v}}$$

- $\#_{j,u,v}$  je súčet váh doplnených vzoriek, v ktorých  $X_j = u$  a  $Z_j = v$
- $n_{j,v}$  je súčet váh doplnených vzoriek, v ktorých  $Z_j = v$

## EM algoritmus (Expectation maximization)

Označenie:  $L(\Theta|X) = \prod_{i=1}^t \Pr(X^{(i)} = x^{(i)}|\Theta)$

Lema:  $L(\Theta^{(k+1)}|X) \geq L(\Theta^{(k)}|X)$ , pričom rovnosť nastáva iba v lokálnom maxime alebo inflexnom bode  $L(\Theta|X)$ .

- Iterácie postupne zvyšujú vierohodnosť
- Môžu skončiť v lokálnom maxime