# DGM MLE AND POSTERIOR

# ORDER OF PRODUCTS

$$p(x \mid \theta) = \prod_{v=1}^{V} p(x_v \mid x_{\mathrm{pa}(v)}, \theta_v)$$

$$p(\mathscr{D} \mid \theta) = \prod_{n=1}^{N} p(x^n \mid \theta) = \prod_{n=1}^{N} \prod_{v=1}^{V} p(x_v^n \mid x_{\mathrm{pa}(v)}^n, \theta_v)$$

$$= \prod_{v=1}^{V} \prod_{n=1}^{N} p(x_v^n \mid x_{\mathrm{pa}(v)}^n, \theta_v) = \prod_{v=1}^{V} p(\mathscr{D}_v \mid \mathscr{D}_{\mathrm{pa}(v)}, \theta_v)$$

# CATEGORICAL – NOTATION

★ For a $v \in [V]$,

values $\quad k \in S_v$

Cartesian product

combined values $\quad c \in C_v = \displaystyle\prod_{s \in \mathrm{pa}(v)} S_s$

★ Cat CPDs

where $\quad P(x_v | x_{\mathrm{pa}(v)} = c) = \mathrm{Cat}(\boldsymbol{\theta}_{vc})$

and $\quad \theta_{vck} = P(X_v = k | X_{\mathrm{pa}(v)} = c)$

$S_d = \{0,1\}$

$S_i = \{0,1\}$

| $d^0$ | $d^1$ |
|---|---|
| 0.6 | 0.4 |

| $i^0$ | $i^1$ |
|---|---|
| 0.7 | 0.3 |

d

i

*Difficulty*

*Intelligence*

*Grade*

g

| | | $g^1$ | $g^2$ | $g^3$ |
|---|---|---|---|---|
| $\mathrm{Cat}(\boldsymbol{\theta}_{g\langle 0,0\rangle})$ | $i^0, d^0$ | 0.3 | 0.4 | 0.3 |
| $\mathrm{Cat}(\boldsymbol{\theta}_{g\langle 0,1\rangle})$ | $i^0, d^1$ | 0.05 | 0.25 | 0.7 |
| $\mathrm{Cat}(\boldsymbol{\theta}_{g\langle 1,0\rangle})$ | $i^1, d^0$ | 0.9 | 0.08 | 0.02 |
| $\mathrm{Cat}(\boldsymbol{\theta}_{g\langle 1,1\rangle})$ | $i^1, d^1$ | 0.5 | 0.3 | 0.2 |

$\theta_{g\langle 1,1\rangle 2}$

$S_g = \{1,2,3\}$

$$C_g = \prod_{s \in \mathsf{pa}(g)} S_s = S_i \times S_d = \{\langle 0,0\rangle, \langle 0,1\rangle, \langle 1,0\rangle, \langle 1,1\rangle\}$$

# THE DGM LIKELIHOOD FACTORIZES

★ Data

$$\mathcal{D} = \{x^1, \dots, x^N\} \qquad x^n = x_1^n, \dots, x_V^n$$

★ Likelihood

$$p(\mathcal{D} \mid \theta) = \prod_{n=1}^{N} p(x^n \mid \theta) = \prod_{n=1}^{N} \prod_{v=1}^{V} p(x_v^n \mid x_{\text{pa}(v)}^n, \theta_v)$$

★ Notation

$$p(X_v = k \mid X_{\text{pa}} = c) = \theta_{vck}$$

★ Counts

$$N_{vck} = \sum_{n=1}^{N} I(x_v^n = k, x_{\text{pa}(v)}^n = c) \qquad N_{vc} = \sum_{n=1}^{N} I(x_{\text{pa}(v)}^n = c)$$

# THE DGM LIKELIHOOD FACTORIZES

★ Complete data

$$\mathscr{D} = \{x^1, \ldots, x^N\}$$

$$x^n = x^n_1, \ldots, x^n_V$$

v's CPD

★ Likelihood

$$p(\mathscr{D} \,|\, \theta) = \prod_{n=1}^{N} p(x^n \,|\, \theta) = \prod_{n=1}^{N} \prod_{v=1}^{V} p(x^n_v \,|\, x^n_{\mathrm{pa}(v)}, \theta_v)$$

# THE DGM LIKELIHOOD FACTORIZES

★ Complete data

$$\mathscr{D} = \{x^1, \ldots, x^N\}$$

$$x^n = x_1^n, \ldots, x_V^n$$

★ Likelihood

$$p(\mathscr{D} \,|\, \theta) = \prod_{n=1}^{N} p(x^n \,|\, \theta) = \prod_{n=1}^{N} \prod_{v=1}^{V} p(x_v^n \,|\, x_{\mathrm{pa}(v)}^n, \theta_v)$$

$$= \prod_{v=1}^{V} \prod_{n=1}^{N} p(x_v^n \,|\, x_{\mathrm{pa}(v)}^n, \theta_v)$$

# THE DGM LIKELIHOOD FACTORIZES

★ Complete data

$$\mathscr{D} = \{x^1, \ldots, x^N\}$$

$$x^n = x_1^n, \ldots, x_V^n$$

★ Likelihood

$$p(\mathscr{D} \mid \theta) = \prod_{n=1}^{N} p(x^n \mid \theta) = \prod_{n=1}^{N} \prod_{v=1}^{V} p(x_v^n \mid x_{\mathrm{pa}(v)}^n, \theta_v)$$

$$= \prod_{v=1}^{V} \prod_{n=1}^{N} p(x_v^n \mid x_{\mathrm{pa}(v)}^n, \theta_v) = \prod_{v=1}^{V} p(\mathscr{D}_v \mid \mathscr{D}_{\mathrm{pa}(v)}, \theta_v)$$

values of v        values of v's parents

Called: decomposable likelihood (factorizes into family-factors)

# PARAMETER AND COUNTS

$$N_{vc} = \sum_{n=1}^{N} I(x_{\text{pa}(v)}^n = c)$$

$$N_{vck} = \sum_{n=1}^{N} I(x_v^n = k, x_{\text{pa}(v)}^n = c)$$

★ Complete data

$$\mathcal{D} = \{x^1, \ldots, x^N\}$$

$$x^n = x_1^n, \ldots, x_V^n$$

★ Likelihood

$$p(\mathcal{D} \mid \theta) = \prod_{n=1}^{N} p(x^n \mid \theta) = \prod_{n=1}^{N} \prod_{v=1}^{V} p(x_v^n \mid x_{\text{pa}(v)}^n, \theta_v)$$

$$= \prod_{v=1}^{V} \prod_{n=1}^{N} p(x_v^n \mid x_{\text{pa}(v)}^n, \theta_v) = \prod_{v=1}^{V} p(\mathcal{D}_v \mid \mathcal{D}_{\text{pa}(v)}, \theta_v)$$

$$= \prod_{v=1}^{V} \prod_{c \in C_v} \prod_{k \in S_v} \theta_{vck}^{N_{vck}}$$

# THE LIKELIHOOD FACTORIZES

$$N_{vc} = \sum_{n=1}^{N} I(x_{\text{pa}(v)}^n = c)$$

$$N_{vck} = \sum_{n=1}^{N} I(x_v^n = k, x_{\text{pa}(v)}^n = c)$$

★ Complete data

$$\mathcal{D} = \{x^1, \ldots, x^N\}$$

$$x^n = x_1^n, \ldots, x_V^n$$

★ Likelihood

$$p(\mathcal{D}\,|\,\theta) = \prod_{n=1}^{N} p(x^n\,|\,\theta) = \prod_{n=1}^{N}\prod_{v=1}^{V} p(x_v^n\,|\,x_{\text{pa}(v)}^n, \theta_v)$$

$$= \prod_{v=1}^{V}\prod_{n=1}^{N} p(x_v^n\,|\,x_{\text{pa}(v)}^n, \theta_v) = \prod_{v=1}^{V} p(\mathcal{D}_v\,|\,\mathcal{D}_{\text{pa}(v)}, \theta_v)$$

$$= \prod_{v=1}^{V}\prod_{c\in C_v}\prod_{k\in S_v} \theta_{vck}^{N_{vck}}$$

So ML estimate

$$\theta_{vck} = N_{vck}/N_{vc}$$

# BAYESIAN PARAMETER LEARNING

★ Decomposable prior

$$p(\boldsymbol{\theta}) = \prod_{v=1}^{V} p(\boldsymbol{\theta}_v) = \prod_{\substack{v \in [V] \\ c \in S_{\mathrm{pa}(v)}}} \mathrm{Dir}(\boldsymbol{\theta}_{vc} | \boldsymbol{\alpha}_{vc})$$

★ Gives decomposable posterior

$$p(\theta | \mathcal{D}) \propto p(\mathcal{D} | \theta) p(\theta)$$

$$= \prod_{v=1}^{V} p(\mathcal{D}_v | \mathcal{D}_{\mathrm{pa}(v)}, \theta_v) p(\theta_v)$$

# BAYESIAN PARAMETER LEARNING

★ Decomposable prior

$$p(\boldsymbol{\theta}) = \prod_{v=1}^{V} p(\boldsymbol{\theta}_v) \quad = \prod_{\substack{v \in [V] \\ c \in S_{\mathrm{pa}(v)}}} \mathrm{Dir}(\boldsymbol{\theta}_{vc}|\boldsymbol{\alpha}_{vc})$$

★ Gives decomposable posterior

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

$$= \prod_{v=1}^{V} p(\mathcal{D}_v|\mathcal{D}_{\mathrm{pa}(v)}, \theta_v)p(\theta_v)$$

$$= \prod_{\substack{v \in [V] \\ c \in S_{\mathrm{pa}(v)}}} \mathrm{Cat}(N_{\boldsymbol{vc}1}, \ldots, N_{\boldsymbol{vc}|S_v|}|\boldsymbol{\theta}_{vc})\mathrm{Dir}(\boldsymbol{\theta}_{vc}|\boldsymbol{\alpha}_{vc})$$

$$\propto \prod_{\substack{v \in [V] \\ c \in S_{\mathrm{pa}(v)}}} \mathrm{Dir}(\boldsymbol{\theta}_{vc}|N_{\boldsymbol{vc}1} + \alpha_{\boldsymbol{vc}1}, \ldots, N_{\boldsymbol{vc}|S_v|} + \alpha_{\boldsymbol{vc}|S_v|})$$

# BAYESIAN PARAMETER LEARNING

★ Decomposable prior

$$p(\boldsymbol{\theta}) = \prod_{v=1}^{V} p(\boldsymbol{\theta}_v) \quad = \prod_{\substack{v \in [V] \\ c \in S_{\mathrm{pa}(v)}}} \mathrm{Dir}(\boldsymbol{\theta}_{vc}|\boldsymbol{\alpha}_{vc})$$

★ Gives decomposable posterior

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

$$= \prod_{v=1}^{V} p(\mathcal{D}_v|\mathcal{D}_{\mathrm{pa}(v)}, \theta_v)p(\theta_v)$$

$$= \prod_{\substack{v \in [V] \\ c \in S_{\mathrm{pa}(v)}}} \mathrm{Cat}(N_{\boldsymbol{vc}1}, \ldots, N_{\boldsymbol{vc}|S_v|}|\boldsymbol{\theta}_{vc})\mathrm{Dir}(\boldsymbol{\theta}_{vc}|\boldsymbol{\alpha}_{vc})$$