

Lecture 8

Bayesian Network Learning

Machine Learning
Ivan Smetannikov

23.04.2021

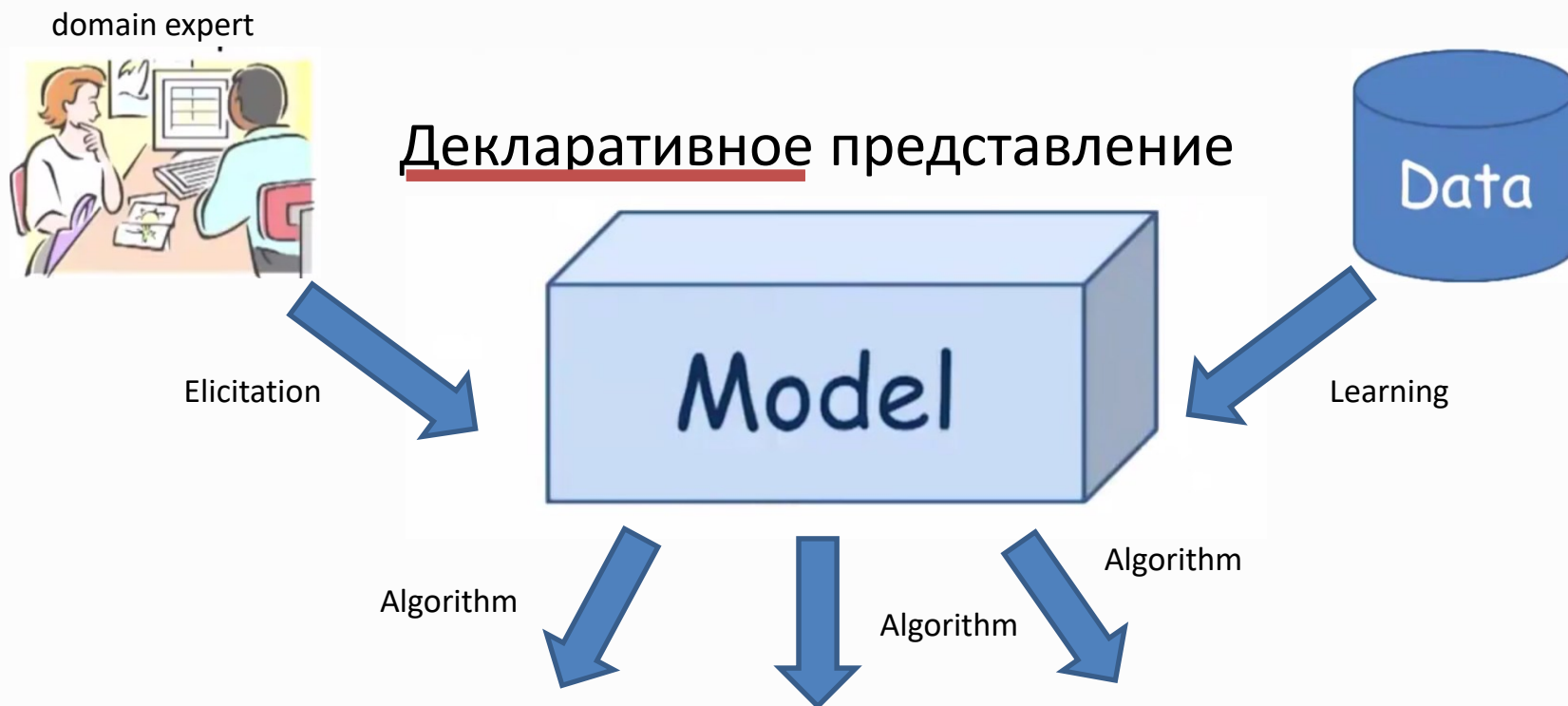
Lecture plan

- Kindly reminder
- Template models brief
- Maximum Likelihood Estimation
- Likelihood, BIC, Bayesian scores
- Learning BN structure

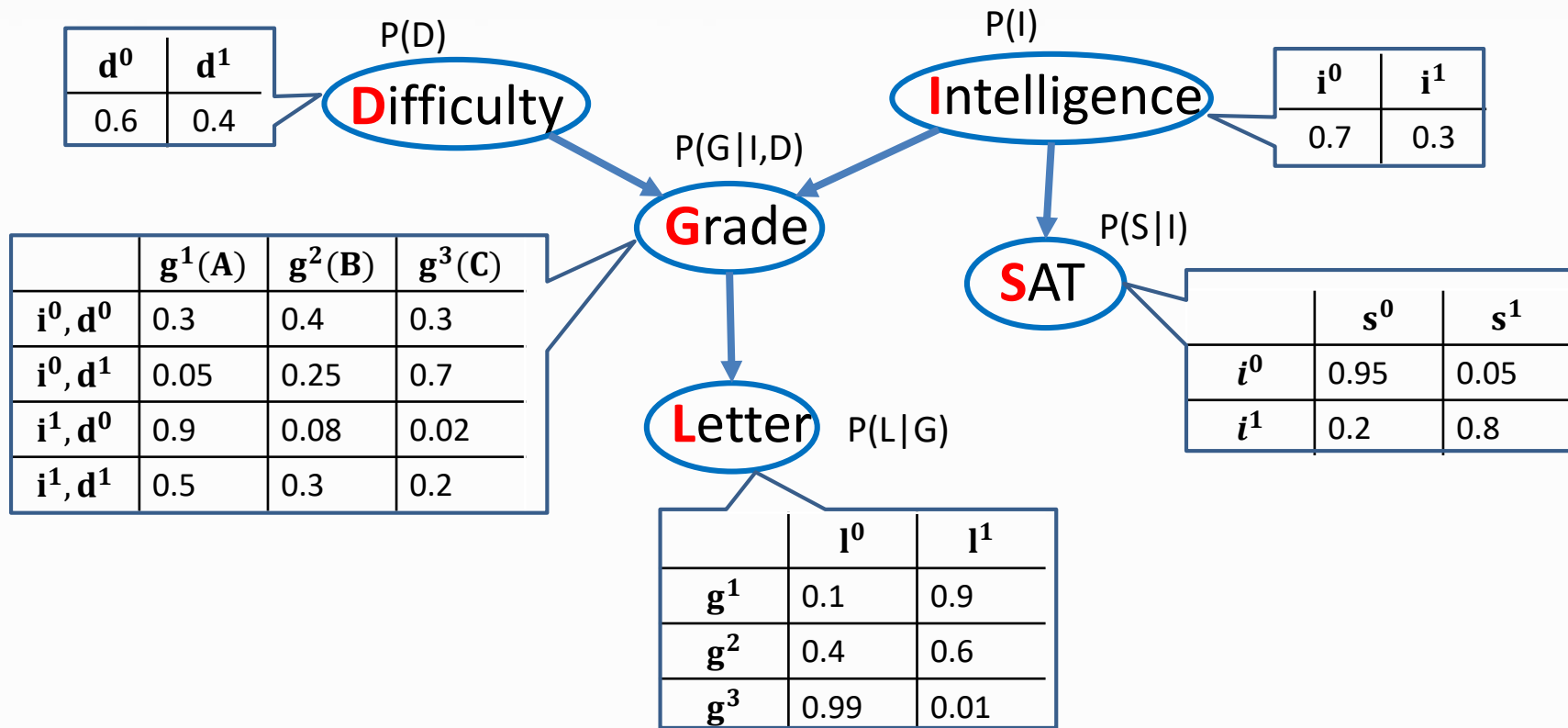
Lecture plan

- Kindly reminder
- Template models brief
- Maximum Likelihood Estimation
- Likelihood, BIC, Bayesian scores
- Learning BN structure

Модели



CPD= Условное вероятностное распределение



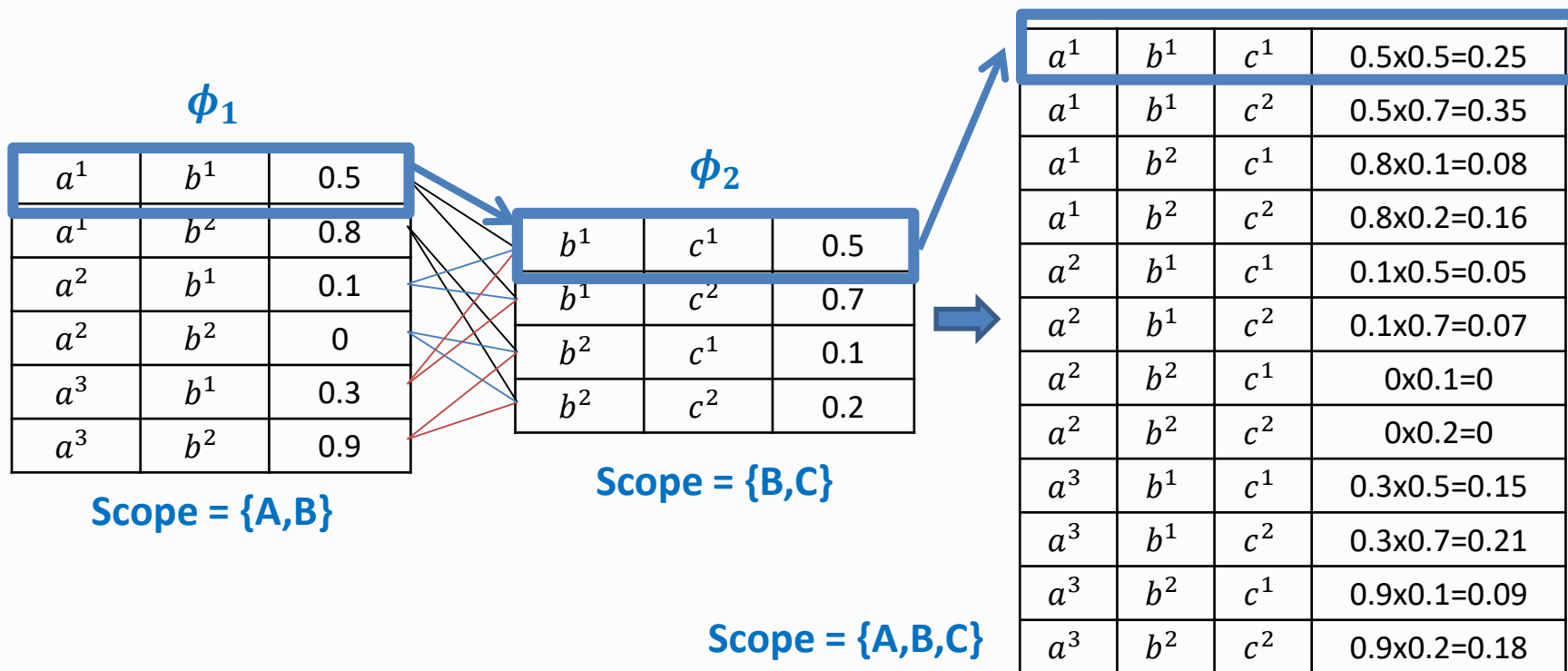
Р факторизуется над G

- Пусть G граф на случайных величинах X_1, \dots, X_n .
- Р факторизуется над G если

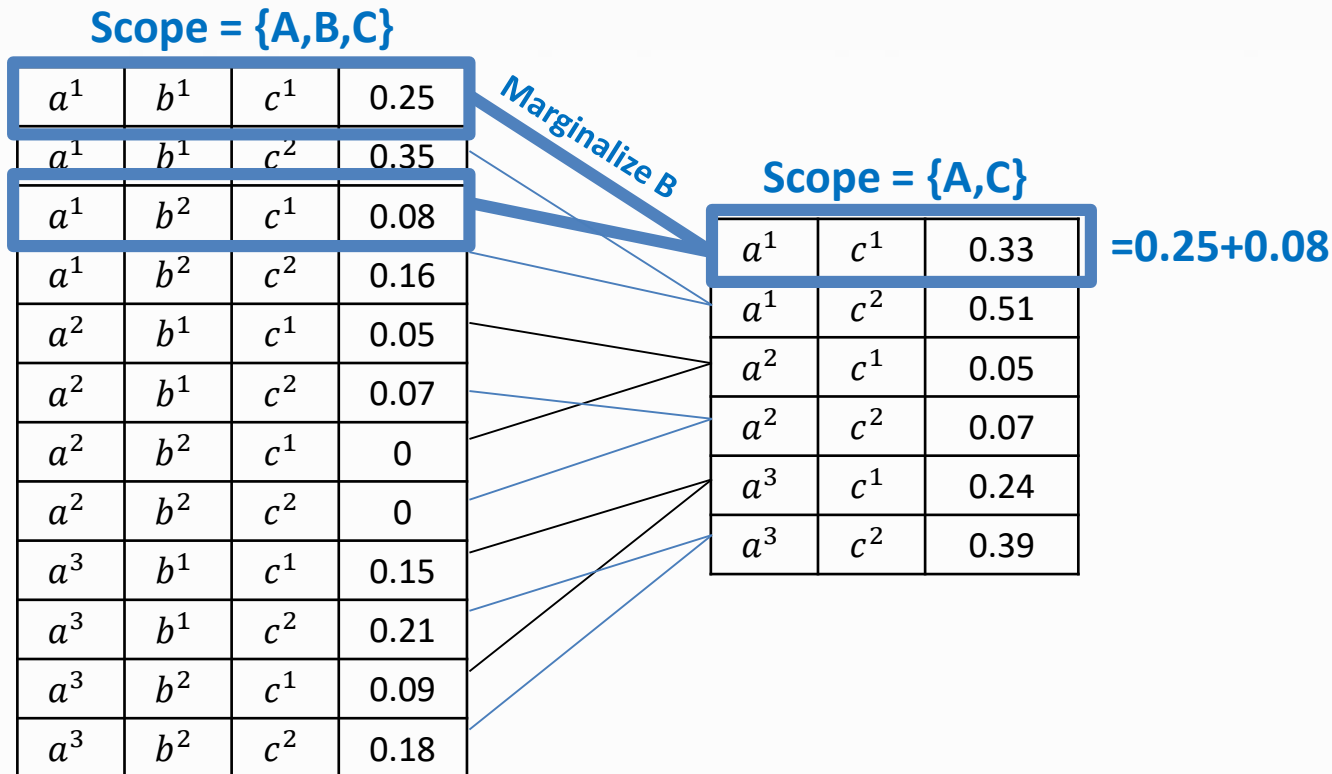
$$P(X_1, \dots, X_n) = \prod_i P(X_i | \text{Par}_G(X_i))$$

chain rule (цепное правило)

Произведение факторов



Маргинализация факторов



Редукция факторов

a^1	b^1	c^1	0.25
a^1	b^1	c^2	0.35
a^1	b^2	c^1	0.08
a^1	b^2	c^2	0.16
a^2	b^1	c^1	0.05
a^2	b^1	c^2	0.07
a^2	b^2	c^1	0
a^2	b^2	c^2	0
a^3	b^1	c^1	0.15
a^3	b^1	c^2	0.21
a^3	b^2	c^1	0.09
a^3	b^2	c^2	0.18

Reduce to the context c^1



a^1	b^1	c^1	0.25
a^1	b^2	c^1	0.08
a^2	b^1	c^1	0.05
a^2	b^2	c^1	0
a^3	b^1	c^1	0.15
a^3	b^2	c^1	0.09

Scope = {A,B}

Цепное правило это Legal Distribution: $\sum P = 1$

$$\begin{aligned}
 \sum_{D,I,G,S,L} P(D, I, G, S, L) &= \sum_{D,I,G,S,L} \underbrace{P(D)P(I)P(G|I,D)P(S|I)P(L|G)}_{\text{chain rule}} \\
 &= \sum_{D,I,G,S} P(D)P(I)P(G|I,D)P(S|I) \sum_L P(L|G) \quad \text{=1} \\
 &= \sum_{D,I,G,S} P(D)P(I)P(G|I,D)P(S|I) \\
 &= \sum_{D,I,G} P(D)P(I)P(G|I,D) \sum_S P(S|I) = \sum_{D,I} P(D)P(I) \sum_G P(G|I,D) \\
 &= \dots \quad \text{=1}
 \end{aligned}$$

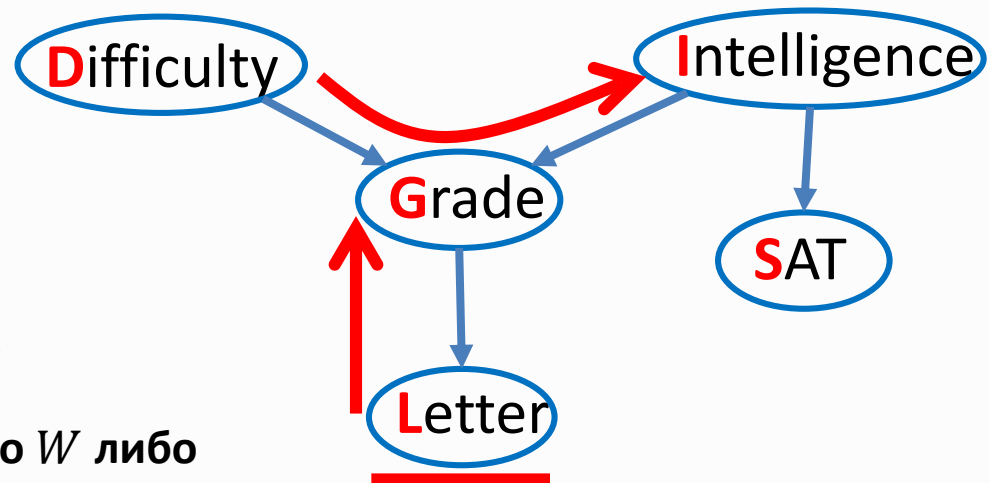
Поток вероятностей: когда X влияет на Y при наличии наблюдения Z ?

- $X \rightarrow Y$ ✓
- $X \leftarrow Y$ ✓
- $X \rightarrow W \rightarrow Y$ ✓ $W \notin Z$
- $X \leftarrow W \leftarrow Y$ ✓
- $X \leftarrow W \rightarrow Y$ ✓
- $X \rightarrow W \leftarrow Y$ ✗ if W и все его потомки не в Z

$W \in Z$

✗
✗
✗
✓

либо W либо
один из его
потомков в Z



Активные пути

Путь $X_1 - \dots - X_k$ называется активным для данного Z если:

- для любой v-structure $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$ получается что X_i или один из его потомков $\in Z$ активирует v-structure
- никакой другой X_i не находится в Z не формирует v-structure

Lecture plan

- Kindly reminder
- **Template models brief**
- Maximum Likelihood Estimation
- Likelihood, BIC, Bayesian scores
- Learning BN structure

Марковское предположение

$$\underline{P(X^{(0:T)})} = \underline{P(X^{(0)})} \prod_{t=0}^{T-1} P(\underline{X^{(t+1)}} | X^{(0:t)})$$

цепное правило для вероятностей

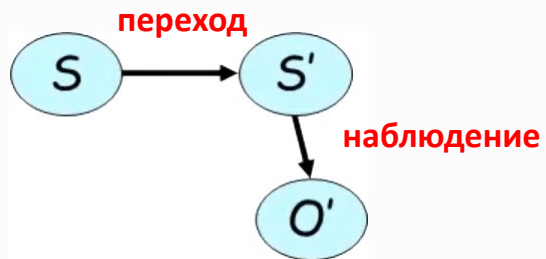
время идет вперед

состояние t+1 состояния 0...t

$$(\underline{X^{(t+1)}} \perp \underline{X^{(0:t-1)}} | \underline{X^{(t)}})$$

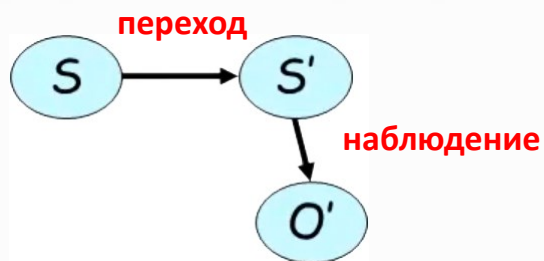
следующий шаг прошлое настоящее

Скрытые марковские модели

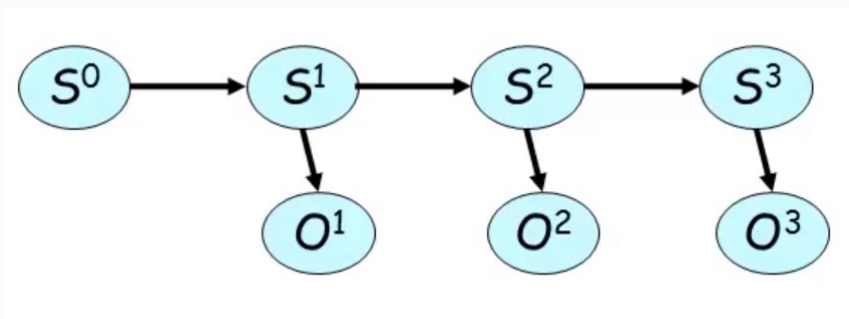


Скрытые марковские модели

2TVN



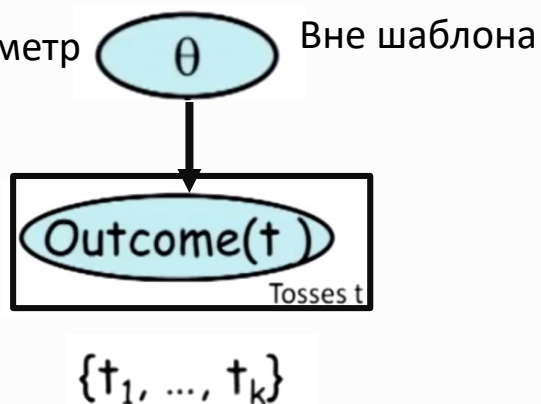
Развернутая сеть



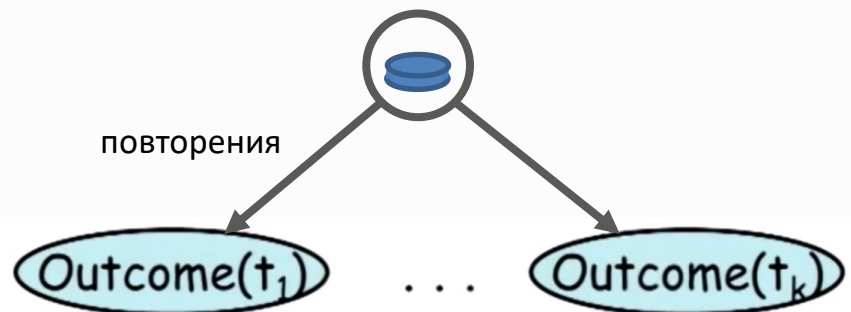
Повторение моделей

около t нет индекса = одинаковое для любого t

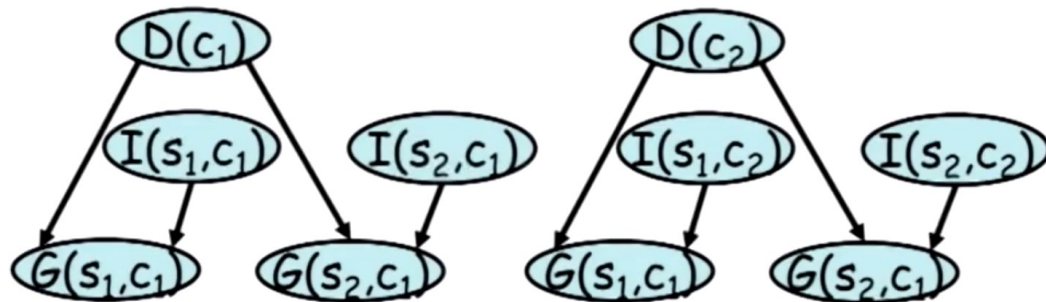
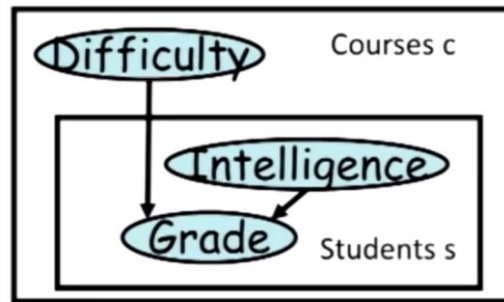
Параметр
(CPD)



повторения



Nested Plates



Lecture plan

- Kindly reminder
- Template models brief
- **Maximum Likelihood Estimation**
- Likelihood, BIC, Bayesian scores
- Learning BN structure

Пример со смещенной монетой

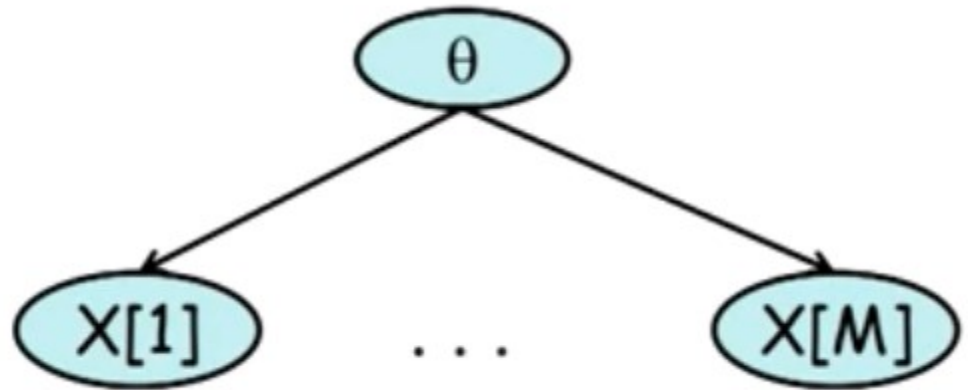
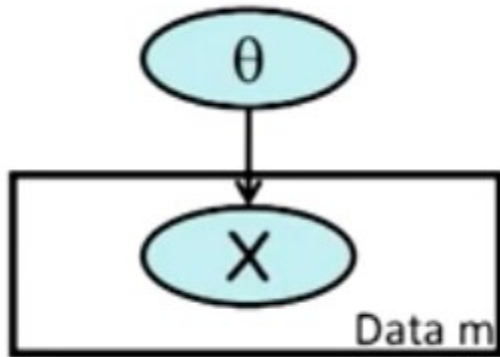
G это распределение Бернулли:

$$P(X = 1) = \theta, P(X = 0) = 1 - \theta$$

$D = \{x[1], \dots, x[M]\}$ это IID (independent identically distributed) семплы из P :

- Броски независимы
- Броски из одного распределения

IID в виде PGM (Вероятностной графической модели)



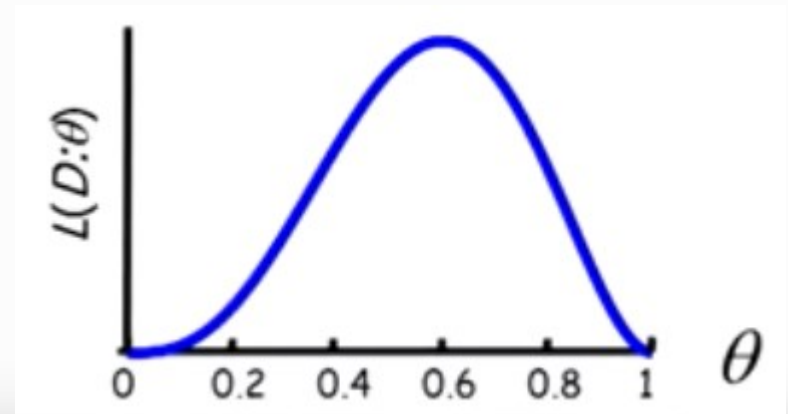
Оценка максимального правдоподобия (MLE)

Цель: найти $\theta \in [0,1]$ которая хорошо предсказывает D

Качество предсказания: правдоподобие D при известном θ

$$L(\theta: D) = P(D|\theta) = \prod_{m=1}^M P(x[m]|\theta)$$

$$L(\theta: \langle H, T, T, H, H \rangle)$$



Maximum Likelihood Estimator

- Наблюдения: M_H орлов и M_T решек
- Ищем θ максимизируя правдоподобие
$$L(\theta: M_H, M_T) = \theta^{M_H} (1 - \theta)^{M_T}$$
- Сводим к максимизации log-likelihood
$$l(\theta: M_H, M_T) = M_H \log \theta + M_T \log(1 - \theta)$$
- Дифференцируем log-likelihood и решаем для заданного θ :

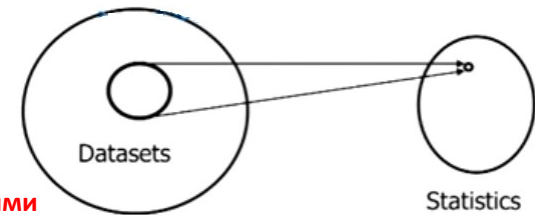
$$\hat{\theta} = \frac{M_H}{M_H + M_T}$$

Sufficient Statistics

- В примере с монетой для вычислений нужны только M_H и M_T
- Функция $s(D)$ называется **sufficient statistic** из экземпляров в вектор R^k если для любых двух наборов данных D и D' и любой $\theta \in \Theta$ выполняется условие:

Если $\sum_{x[i] \in D} s(x[i]) = \sum_{x[i] \in D'} s(x[i])$ то
 $L(\theta: D) = L(\theta: D')$

Набор данных с
несколькими
последовательностями
бросков



Sufficient Statistic для Мультинома

- Для набора данных D на переменной X с k возможными значениями, sufficient statistics это просто счетчики $\langle M_1, \dots, M_k \rangle$, где M_i число раз, когда $X[m] = x^i$

- Sufficient statistic $s(x)$ это k -мерный tuple где

$$s(x^i) = (0, \dots, 0, \underset{i}{1}, 0, \dots, 0), \quad \sum s(x[m]) = \langle M_1, \dots, M_k \rangle$$

$$L(\theta: D) = \prod_{i=1}^k \theta_i^{M_i}$$

Maximum Likelihood Estimation

Принцип MLE в общем случае: выбираем θ
для максимизации $L(D: \theta)$

MLE для байесовских сетей

- Параметры: $\{\theta_x: x \in \text{val}(X)\}, \{\theta_{y|x}: x \in \text{val}(X), y \in \text{val}(Y)\}$

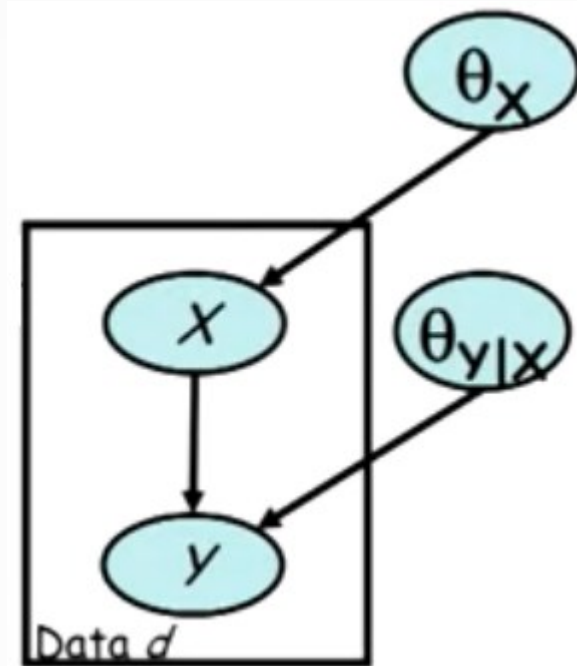
$$L(\Theta: D) = \prod_{m=1}^M P(x[m], y[m]: \theta)$$

$$= \prod_{m=1}^M P(x[m]: \theta) P(y[m]|x[m]: \theta)$$

$$= \left(\prod_{m=1}^M P(x[m]: \theta) \right) \left(\prod_{m=1}^M P(y[m]|x[m]: \theta) \right)$$

$$= \left(\prod_{m=1}^M P(x[m]: \theta_X) \right) \left(\prod_{m=1}^M P(y[m]|x[m]: \theta_{Y|X}) \right)$$

Local likelihood



MLE для байесовских сетей

- Likelihood для Байесовской сети более общего вида

$$\begin{aligned} L(\Theta: D) &= \prod_{m=1}^M P(x[m]: \Theta) \\ &= \prod_m \prod_{\substack{\text{parents of } x_i \\ \text{chain rule}}} P(x_i[m] | U_i[m]: \Theta_i) \\ &= \prod_m \prod_i^m P(x_i[m] | U_i[m]: \Theta_i) = \prod_i L_i(D: \Theta_i) \\ &\quad \text{Local likelihood} \end{aligned}$$

MLE для табличных CPDs

$$\begin{aligned} \prod_{m=1}^M P(x[m]|u[m]: \theta) &= \prod_{m=1}^M P(x[m]|u[m]: \theta_{X|U}) \\ &= \prod_{x,u} \left(\prod_{m: x[m]=x, u[m]=u} P(x[m]|u[m]: \theta_{X|U}) \right) \\ &= \prod_{x,u} \left(\prod_{m: x[m]=x, u[m]=u} \theta_{x|u} \right) \\ &= \prod_{x,u} \theta_{x|u}^{M[x,u]} \\ \theta_{x|u} &= \frac{M[x, u]}{\sum_{x'} M[x', u]} = \frac{M[x, u]}{M[u]} \end{aligned}$$

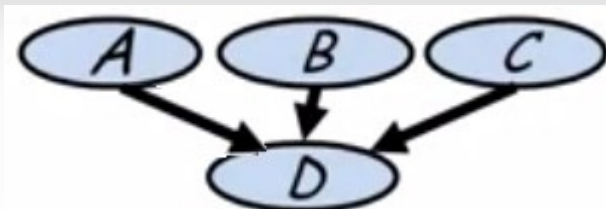
Итого

- Для Байесовской Сети с **disjoint параметрами** в CPD, правдоподобие **декомпозируется** как произведение **локальных** функций правдоподобия, по одной на каждую переменную
- Для табличных CPD, локальные правдоподобия можно дальше **декомпонировать** как произведение правдоподобий для мультиномов, одно для каждой комбинации родителей

Lecture plan

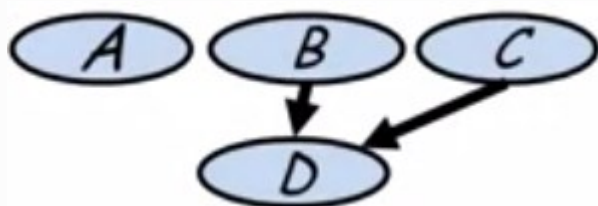
- Kindly reminder
- Template models brief
- Maximum Likelihood Estimation
- **Likelihood, BIC, Bayesian scores**
- Learning BN structure

Важность точного построения



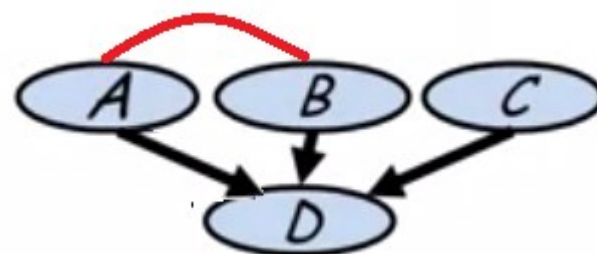
Оригинальный граф

Пропустили ребро



- Некорректные независимости
- Нельзя выучить правильное распределение P^*
- Лучше обобщающая способность

Лишнее ребро



- Ложные зависимости
- Можем выучить правильное распределение P^*
- Увеличенное число параметров
- Хуже обобщающая способность

Что нужно для построения?

- Метрика
- Оптимизатор

Пример

- Найти (G, θ) которая максимизирует правдоподобие

$$score_L(G: D) = l\left((\hat{\theta}, G): D\right)$$

Пример



$$score_L(G_0: D) = \sum_m (\log \hat{\theta}_{x[m]} + \log \hat{\theta}_{y[m]})$$



$$score_L(G_1: D) = \sum_m (\log \hat{\theta}_{x[m]} + \log \hat{\theta}_{y[m]|x[m]})$$

$$score_L(G_1: D) - score_L(G_0: D) = \sum_m (\log \hat{\theta}_{y[m]|x[m]} - \log \hat{\theta}_{y[m]})$$

$$= \sum_{x,y} M[x,y] \log \hat{\theta}_{y|x} - \sum_y M[y] \log \hat{\theta}_y$$

$$= M \sum_{x,y} \hat{P}(x,y) \log \hat{P}(y|x) - M \sum_y P(y) \log \hat{P}(y)$$

$$= M \left(\sum_{x,y} \hat{P}(x,y) \log \hat{P}(y|x) - \sum_{x,y} \hat{P}(x,y) \log \hat{P}(y) \right)$$

$$= M \left(\sum_{x,y} \hat{P}(x,y) \log \frac{\hat{P}(y,x)}{\hat{P}(x)\hat{P}(y)} \right) = M \cdot I_{\hat{P}}(X; Y)$$

Избегаем переобучения

- Ограничиваем пространство
 - Число родителей или число параметров
- Дополнительные штрафы на сложность
 - Явные
 - Bayesian score усредняется по всем возможным значениям параметров

BIC score

- Штрафуем сложность

$$score_{BIC}(G:D) = score_L(G:D) - \frac{\log M}{2} Dim[G]$$

$Dim[G] = 2$ в степени числа ребер

Bayesian score

Marginal likelihood

Prior over structures

$$P(G:D) = \frac{P(D|G)P(G)}{P(D)}$$

Marginal probability of Data

$$\text{score}_B(G:D) = \log P(D|G) + \log P(G)$$

Marginal Likelihood

$$P(\mathcal{D} \mid \mathcal{G}) = \prod_i \prod_{\mathbf{u}_i \in \text{Val}(\mathbf{Pa}_{X_i}^{\mathcal{G}})} \frac{\Gamma(\alpha_{X_i|\mathbf{u}_i})}{\Gamma(\alpha_{X_i|\mathbf{u}_i} + M[\mathbf{u}_i])} \prod_{x_i^j \in \text{Val}(X_i)} \left[\frac{\Gamma(\alpha_{x_i^j|\mathbf{u}_i} + M[x_i^j, \mathbf{u}_i])}{\Gamma(\alpha_{x_i^j|\mathbf{u}_i})} \right]$$

$$\Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt$$

$$\Gamma(x) = x \cdot \Gamma(x-1)$$

Lecture plan

- Kindly reminder
- Template models brief
- Maximum Likelihood Estimation
- Likelihood, BIC, Bayesian scores
- Learning BN structure

Строим деревья/леса

- Леса
 - Максимум один родитель для каждой переменной
- Почему деревья?
 - Математика
 - Эффективная оптимизация
 - Разреженная параметризация

Обучение

- $p(i)$ = родитель X_i или 0 если у X_i нет родителей

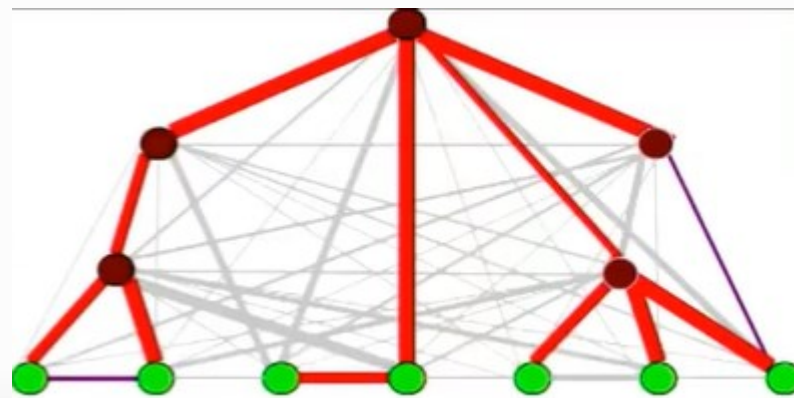
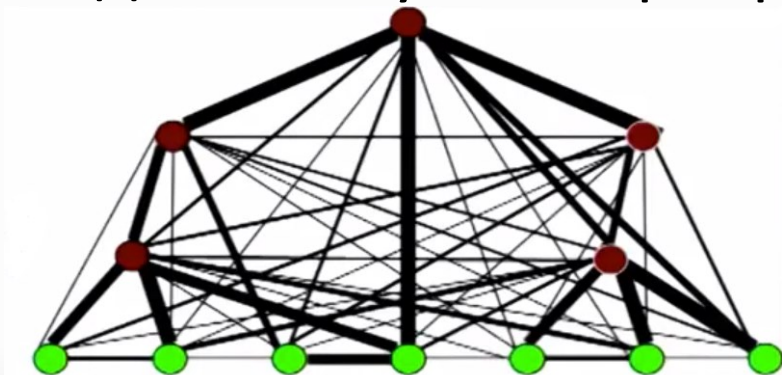
$$\begin{aligned}\text{score}(\mathcal{G} : \mathcal{D}) &= \sum_i \text{score}(X_i \mid \text{Pa}_{X_i}^{\mathcal{G}} : \mathcal{D}) \\ &= \sum_{i:p(i)>0} \text{score}(X_i \mid X_{p(i)} : \mathcal{D}) + \sum_{i:p(i)=0} \text{score}(X_i : \mathcal{D}) \\ &= \sum_{i:p(i)>0} (\text{score}(X_i \mid X_{p(i)} : \mathcal{D}) - \text{score}(X_i : \mathcal{D})) + \sum_{i=1}^n \text{score}(X_i : \mathcal{D})\end{aligned}$$

Улучшение по сравнению
с пустой сетью

Score of empty network

Обучение

- Определить неориентированный граф с вершинами $\{1, \dots, n\}$
- Назначаем $w = \max(\text{score}, 0)$
- Строим остовное дерево Крускалом или еще чем-то
- Удаляем нулевые ребра



Что-то кроме деревьев?

- В общем случае всё не так очевидно
 - Пример: Если позволить наличие двух родителей, то жадный алгоритм уже не гарантирует оптимальное множество
- Теорема
 - Поиск сети со структурой имеющий максимальный score для ситуации когда для каждой переменной разрешено не более k родителей является NP-трудной для $k > 1$

Эвристический поиск

- Операторы:
 - Пошаговые: добавление, удаление, инверсия ребер
 - Глобальные шаги
- Техники:
 - Greedy hill-climbing
 - Best first search
 - Simulated Annealing
 - ...