



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

PROBABILISTIC GRAPHICAL MODEL SESSION # 14: LEARNING

SEETHA PARAMESWARAN
seetha.p@pilani.bits-pilani.ac.in

TABLE OF CONTENTS

- 1 STRUCTURE LEARNING IN BAYESIAN NETWORKS
- 2 SCORE-BASED LEARNING
- 3 BAYESIAN SCORE
- 4 TREE STRUCTURED NETWORK

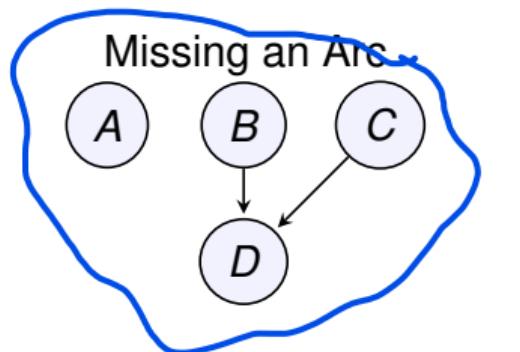
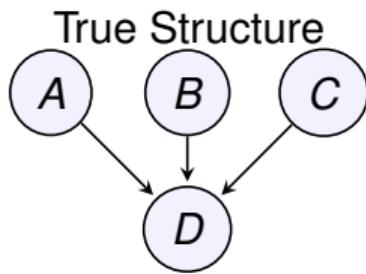
STRUCTURE LEARNING IN BAYESIAN NETWORKS

- We do not know the structure of the Bayesian network in advance.
- The data are generated IID from an underlying distribution $P^*(\mathcal{X})$.
- $P^*(\mathcal{X})$ is induced by some Bayesian network G^* over \mathcal{X} .

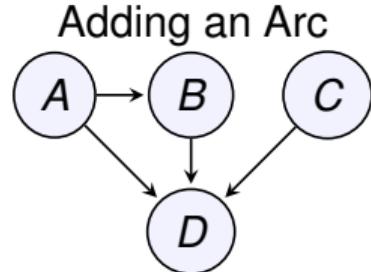
WHY STRUCTURE LEARNING?

- To learn model for new queries when the domain expert is not perfect.
- For structure discovery, when inferring network structure is goal in itself.

STRUCTURE – WHICH ONE?



- Incorrect
independencies
- Correct P^* cannot be
learned
- Could generalize better



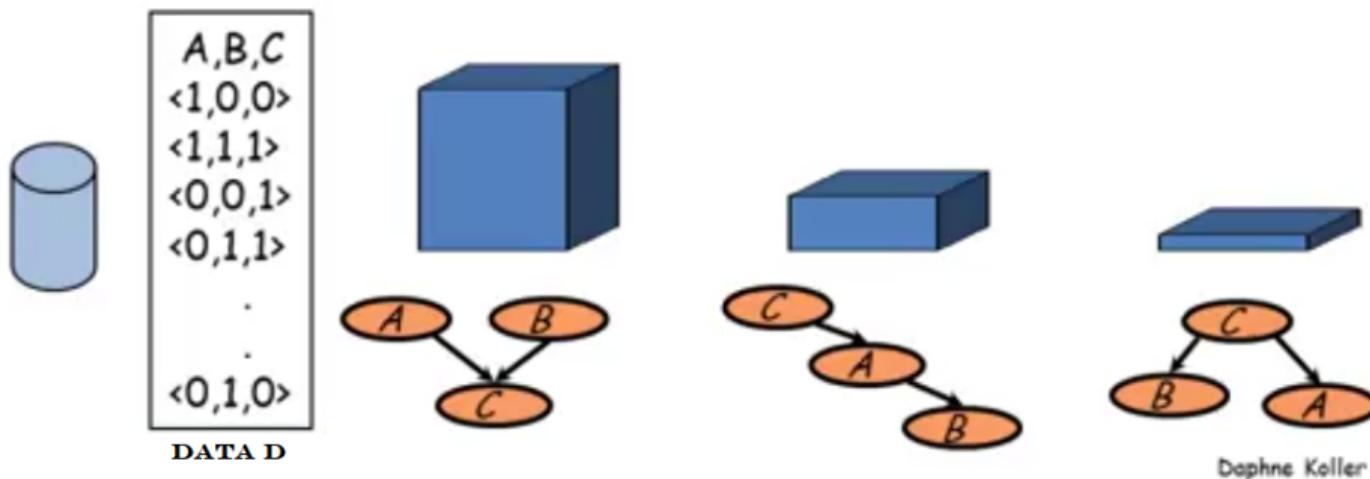
- Spurious dependencies
- Cannot correctly learn P^*
- More number of
parameters are learned
- Worse generalization

TABLE OF CONTENTS

- ① STRUCTURE LEARNING IN BAYESIAN NETWORKS
- ② SCORE-BASED LEARNING
- ③ BAYESIAN SCORE
- ④ TREE STRUCTURED NETWORK

SCORE-BASED LEARNING

- ① Define **scoring function** that evaluates how well a structure matches the data.



- ② Search for a structure that maximizes the score. This converts the learning problem into an optimization problem.

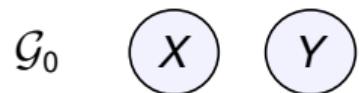
LIKELIHOOD STRUCTURE SCORE

- Simplest score
- Find the graph and parameters that maximize the likelihood of the data.

$$score_L(\mathcal{G} : \mathcal{D}) = \ell((\hat{\theta}, \mathcal{G}) : \mathcal{D})$$

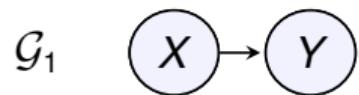
- $\hat{\theta}$ = MLE of parameters given \mathcal{G} and \mathcal{D}

EXAMPLE



$$score_L(\mathcal{G}_0 : \mathcal{D}) = \sum_m (\log \hat{\theta}_{x[m]} + \log \hat{\theta}_{y[m]})$$

\sum $\log \hat{\theta}_{x[m]}$ $\log \hat{\theta}_{y[m]}$



$$score_L(\mathcal{G}_1 : \mathcal{D}) = \sum_m (\log \hat{\theta}_{x[m]} + \log \hat{\theta}_{y[m]|x[m]})$$

EXAMPLE

$$\begin{aligned}
 & \text{score}_L(\mathcal{G}_1 : \mathcal{D}) - \text{score}_L(\mathcal{G}_0 : \mathcal{D}) = \sum_m (\log \hat{\theta}_{y[m] | x[m]}) - \sum_m (\log \hat{\theta}_{y[m]}) \\
 &= \sum_{x,y} M[x,y] \log \hat{\theta}_{y|x} - \sum_y M[y] \log \hat{\theta}_y \quad \text{sufficient statistics} \\
 &= M \sum_{x,y} \hat{P}(x,y) \log \hat{P}(y | x) - M \sum_y \hat{P}(y) \log \hat{P}(y) \quad \text{empirical distribution } \hat{P}(x,y) \\
 &= M \left(\sum_{x,y} \hat{P}(x,y) \log \hat{P}(y | x) - \sum_{x,y} \hat{P}(x,y) \log \hat{P}(y) \right) \quad \sum_x \hat{P}(x,y) = \hat{P}(y) \\
 &= M \sum_{x,y} \hat{P}(x,y) \log \frac{\hat{P}(x,y)}{\hat{P}(x)\hat{P}(y)} \quad \hat{P}(y | x) = \frac{\hat{P}(x,y)}{\hat{P}(x)} \\
 &= M \cdot I_{\hat{P}}(X; Y) \quad \# \text{ samples times Mutual Information}
 \end{aligned}$$

LIKELIHOOD STRUCTURE SCORE

- In general, the Likelihood Structure Score decomposes as number of samples M times the difference between the mutual information and the sum of entropies of variables X .

$$score_L(\mathcal{G} : \mathcal{D}) = M \sum_{i=1}^n I_{\hat{P}}(X_i; Pa_{X_i}^{\mathcal{G}}) - M \sum_i H_{\hat{P}}(X_i)$$

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

$$H_{\hat{P}}(X) = - \sum_x P(x) \log P(x)$$

- The score is higher if X_i is correlated with its parents. Helps in identifying the Parent-Child relationship and place the nodes accordingly.

LIMITATIONS OF LIKELIHOOD SCORE

$$score_L(\mathcal{G}_1 : \mathcal{D}) - score_L(\mathcal{G}_0 : \mathcal{D}) = M \cdot I_{\hat{P}}(X; Y)$$

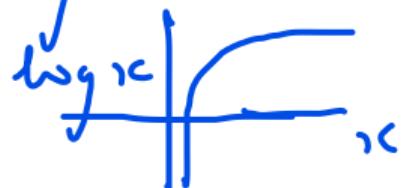
- Positive difference suggests that \mathcal{G}_1 has to be chosen.
 - Negative difference implies that \mathcal{G}_0 has to be chosen.
 - Mutual Information is always positive. $I_{\hat{P}}(X; Y) \geq 0$
 - Mutual Information $I_{\hat{P}}(X; Y) = 0$ iff X and Y are independent.
 - In empirical distribution \hat{P} , due to statistical fluctuations while taking samples, $I_{\hat{P}}(X; Y) > 0$ almost always.
 - Score is maximized for fully connected network.
 - Over-fits the data
- why is this always true?*

why is mutual information always non-negative?

$$I(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

\log is a concave function since

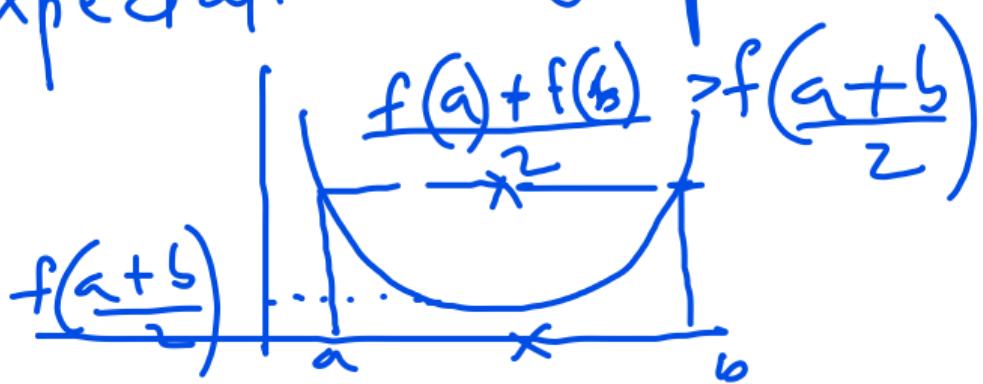
$$\frac{d^2}{dx^2} \log x = -\frac{1}{x^2} < 0$$



- \log is a convex function

We can apply Jensen's inequality on
 $-\log$, which is a convex function

$E(F(X)) \geq F(E(X))$ where X is
a random variable, and F is a convex
function, E is expectation over $\{x\}$ of X



$$\begin{aligned}
 I(x, y) &\geq -\log \sum_{x \in X} \sum_{y \in Y} p(x, y) \frac{p(x)p(y)}{p(x, y)} \\
 &\geq -\log \sum_{x \in X} \sum_{y \in Y} p(x)p(y) \\
 &\geq -\log \left(\left(\sum_{x \in X} p(x) \right) \left(\sum_{y \in Y} p(y) \right) \right) \\
 &\geq -\log 1 \geq 0
 \end{aligned}$$

AVOID OVER-FITTING

- Restrict the hypothesis space
 - ▶ Restrict # parents – common strategy
 - ▶ Restrict # parameters
- Score that penalize complexity
 - ▶ Explicitly penalize model complexity
 - ▶ Use Bayesian score that averages over all possible parameter values.

TABLE OF CONTENTS

- ① STRUCTURE LEARNING IN BAYESIAN NETWORKS
- ② SCORE-BASED LEARNING
- ③ BAYESIAN SCORE
- ④ TREE STRUCTURED NETWORK

BAYESIAN SCORE

The image shows two handwritten mathematical expressions. The first expression, enclosed in a blue hand-drawn oval, is $\ell(\hat{\theta} : \mathcal{G})$. The second expression, to its right, is $\Pr(\hat{\theta})$.

- Bayesian score is derived from Bayesian paradigm – Any uncertainty can be represented by using a probability distribution for it.
- Uncertainty over graph can be represented using a prior over graph.
- Uncertainty over parameters can be represented using a prior over parameters.

BAYESIAN SCORE

- Optimization problem – find a graph that maximizes a probability of \mathcal{G} over \mathcal{D}

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

$P(\mathcal{D} | \mathcal{G})$ – Marginal likelihood

$P(\mathcal{G})$ – Prior over structures

$P(\mathcal{D})$ – Marginal probability of Data – normalizing constant (ignored)

BAYESIAN SCORE

- Taking log on both sides, Bayesian Score

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

$\text{score}_B(\mathcal{G} : \mathcal{D})$ – Bayesian score

$\log P(\mathcal{D} | \mathcal{G})$ – Marginal likelihood of Data given graph \mathcal{G}

$\log P(\mathcal{G})$ – Prior over structures – less significant

MARGINAL LIKELIHOOD OF DATA GIVEN GRAPH

- As $M \rightarrow \infty$ a network with Dirichlet prior satisfies

$$\log P(\mathcal{D} | \mathcal{G}) = \ell(\hat{\theta}_{\mathcal{G}} : \mathcal{D}) - \frac{\log M}{2} \text{score}_{\text{BIC}}(\mathcal{G} : \mathcal{D}) + \mathcal{O}(1)$$

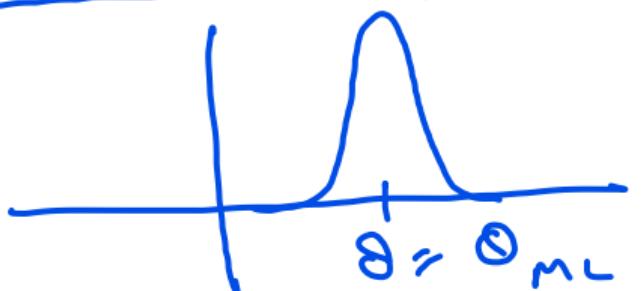
score_{BIC}(G : D)

- Marginal likelihood is the Bayesian Information Criterion (BIC) score.
- $\ell(\hat{\theta}_{\mathcal{G}} : \mathcal{D})$ – Log likelihood score with Maximum likelihood parameters $\hat{\theta}_{\mathcal{G}}$
- M – # training data instances
- $\text{Dim}[\mathcal{G}]$ – # independent parameters = # entries in the distribution - 1

$$\text{Pr}(\theta|G) = \int \text{Pr}(P|\theta, G) \text{Pr}(\theta|G) d\theta$$

$\int \text{Pr}(\theta|G) d\theta$
 $\text{Pr}(\theta|G)$

$\text{Pr}(P|\theta, G) \text{Pr}(\theta|G)$ can be approximated as
 a function that peaks at $\theta = \hat{\theta}_{ML} =$
max likelihood parameters



We can use second-order approximation
(Taylor's series) and write

$$\log(\Pr(\theta|\phi_1, g)\Pr(\theta|g)) = \log\left(\Pr\left(\theta|\theta_{MC}, g\right)\Pr\left(\frac{\theta_{MC}}{g}\right)\right)$$
$$-\frac{1}{2}(\theta - \theta_{MC})^T H (\theta - \theta_{MC})$$
$$g(\theta) = g(\theta_{MC}) e^{-\frac{1}{2}(\theta - \theta_{MC})^T H (\theta - \theta_{MC})}$$

where $H = \begin{bmatrix} \frac{\partial^2 \ln g(\theta)}{\partial \theta \partial \theta} & \\ & 0 \end{bmatrix}$ a $d \times d$ matrix

$P(\theta|G) = \int g(\theta) d\theta$ is a Gaussian integral can be approximated as

$$\frac{(2\pi)^{d/2}}{|H|^{\frac{1}{2}}} P_v(\theta|\theta_{MC}) \text{ or } (\theta_{MC}/g)$$

Taking $\log P(\theta|G)$ and approximate $|H| = M^d$ + get the result.

More specifically, $\int_{-\infty}^{\infty} e^{-\frac{1}{2}x^T A x} dx$

$$= \frac{(2\pi)^{n/2}}{|A|^{1/2}}$$

Here $A = H$ and $n=d$, and we have a constant $\Pr(D|\theta_{MC}, G) \Pr(\theta_{MC}|G)$ in the integrand leading to the integral being equal to $\frac{(2\pi)^{d/2}}{|H|^{1/2}} \Pr(D|\theta_{MC}, G) \Pr(\theta_{MC}|G)$

$$\Pr(D|G) = \frac{(2\pi)^{d/2}}{|H|^{\frac{d}{2}}} \Pr(D|O_{MC}, G) \Pr(O_{MC}|G)$$

$$\begin{aligned} \log \Pr(D|G) &= \frac{d}{2} \log 2\pi - \frac{1}{2} \log |H| + \\ &\quad \boxed{\log (\Pr(D|O_{MC}, G) \Pr(O_{MC}|G))} \\ &= O(1) - \frac{1}{2} \log |H| + \ell(O_{MC}; G) \end{aligned}$$

We can take $|H| \approx M^d$ + const $H = \begin{bmatrix} M & M & \dots & M \end{bmatrix}_{d \times d}$

$$\log(\Pr(D|G)) = \ell(O_{MC}; G) - \frac{d}{2} \log M + O(1)$$

BAYESIAN SCORE

- As $M \rightarrow \infty$ Marginal Likelihood is equivalent to BIC score.
- That implies Bayesian score is equivalent to BIC score.
- The graph or its I-equivalent graph will have the highest score among all possible graphs.

A closer look at the formula:

$$\text{score}_{\text{BIC}}(G:D) = l(\hat{\theta}:G) - \frac{-\log M}{2} \dim(G)$$

We can rewrite $l(\hat{\theta}:G)$ as follows

$$l(\hat{\theta}:G) = M \sum_{i=1}^n I_p(x_i, p_{x_i}) - M \sum_{i=1}^n H_p(x_i) \\ - \frac{-\log M}{2} \dim(G)$$

$$score_{\text{BIC}}(G:D) = M \sum_{i=1}^n I_p^{\hat{P}}(X_i, \text{par}_i) - M \sum_{i=1}^n H_p^{\hat{P}}(X_i) - \frac{\log M \cdot \dim(G)}{2}$$

- Entropy terms do not depend on the graph and can be ignored
- The stronger the dependence between a variable and its parents, the higher the score
- The more complex the network, the lower the score
- Mutual information term grows linearly in M , complexity term grows logarithmically \Rightarrow more the data, more emphasis given to fitting data

TABLE OF CONTENTS

- 1 STRUCTURE LEARNING IN BAYESIAN NETWORKS
- 2 SCORE-BASED LEARNING
- 3 BAYESIAN SCORE
- 4 TREE STRUCTURED NETWORK

OPTIMIZATION PROBLEM

- Input
 - ▶ Training data
 - ▶ Scoring function
 - ▶ Set of possible structures
- Output
 - ▶ Network that maximizes the score
- Key property for computation efficiency
 - ▶ Decomposability of score

$$score(\mathcal{G} : \mathcal{D}) = \sum_i score(X_i \mid Pa_{X_i}^{\mathcal{G}} : \mathcal{D})$$

LEARNING TREES / FORESTS

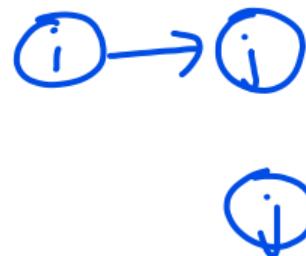
$$\begin{aligned}
 \underline{\text{score}(\mathcal{G} : \mathcal{D})} &= \sum_i \text{score}(X_i | \text{Pa}_{X_i}^{\mathcal{G}} : \mathcal{D}) \\
 &= \left[\sum_{i:p(i)>0} \text{score}(X_i | X_{p(i)} : \mathcal{D}) - \text{score}(X_i : \mathcal{D}) \right] + \sum_{i=1}^n \text{score}(X_i : \mathcal{D})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


- Score = sum of edge scores + constant.

ALGORITHM FOR LEARNING TREES / FORESTS

Algorithm

- Define unirected graph with nodes $\{1, \dots, n\}$
- Set $w(i \rightarrow j) = \max[score(X_j | X_i) - score(X_j), 0]$
- Find forest with maximal weights
 - ▶ Use MST algorithms like Prim's or Kruskal's
 - ▶ Remove edges with weight 0.
- Generates tree if likelihood score is used and a forest if BIC score is used.
- Network with only one parent.



GENERAL NETWORK

- For general trees, with # parents ≥ 2 , use
 - ▶ Heuristic hill climbing algorithm
 - ▶ Best first search
 - ▶ Simulated annealing

GREEDY HILL CLIMBING ALGORITHM

- Start with a given network
 - ▶ Empty network
 - ▶ Best tree
 - ▶ Random network
 - ▶ Prior knowledge
- At each iteration
 - ▶ Consider score of all possible changes like addition, removal and reversal of edges.
 - ▶ Apply change that most improves the score (greedy approach)
- Stop when no modification improves score. This gives local maxima.

GREEDY HILL CLIMBING ALGORITHM

- Pitfalls
 - ▶ Local maxima
 - ▶ Plateau – often in BN due to I-equivalent structures.
- To remove pitfalls
 - ▶ Use edge reversals
 - ▶ Random restarts
 - ★ Take some random steps and then start climbing again
 - ▶ Tabu lists
 - ★ Keep a list of most recent K steps and search cannot reverse any of these steps.

CHOW LIU ALGORITHM

Chow-Liu Algorithm

Step 1: For each pair of variables A,B, use data to estimate $P(A,B)$, $P(A)$, $P(B)$

Step 2: For each pair of variables A,B, calculate mutual information

$$I(A, B) = \sum_a \sum_b P(a, b) \log \frac{P(a, b)}{P(a)P(b)}$$

Step 3: Calculate the maximum spanning tree over the set of variables, using edge weights $I(A,B)$ (given N variables, this costs only $O(N^2)$ time)

Step 4: Add arrows to edges to form a directed-acyclic graph by picking an arbitrary node as root and directing edges outward from the root

Step 5: Learn the CPD's for this graph

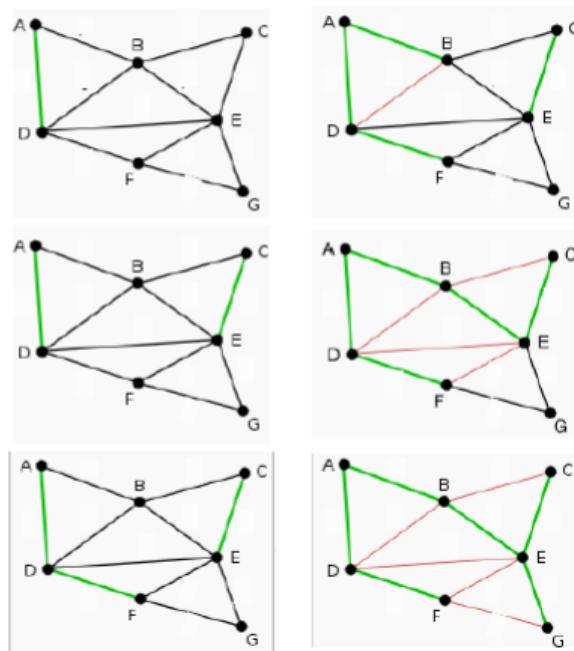
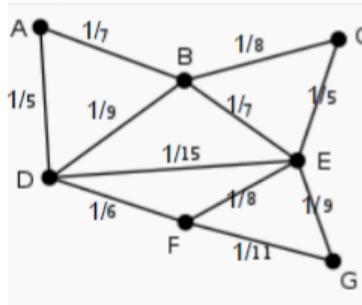
KRUSKAL'S ALGORITHM

Maximum Spanning Tree Algorithm

- Kruskal's Algorithm
 - Start with the empty graph and add edges one by one
 - As the next edge to add, choose one that
 - Is not in graph yet
 - Does not introduce a cycle. Has the maximum weight

CHOW LIU ALGORITHM

Chow-Liu algorithm example Greedy Algorithm to find Max-Spanning Tree



[courtesy A. Singh, C. Guestrin]