

Learning Graphical Models: Overview

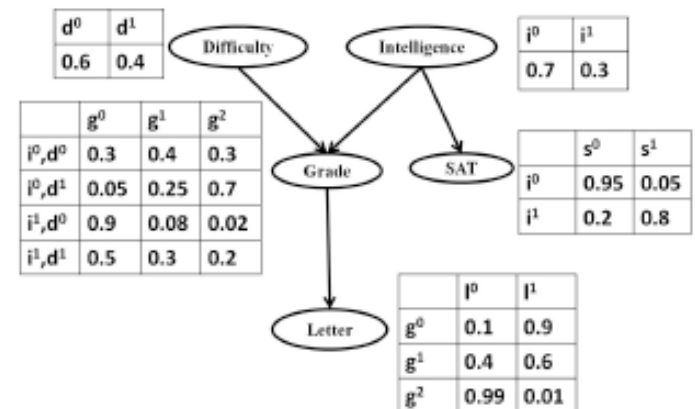
Sargur Srihari

srihari@cedar.buffalo.edu

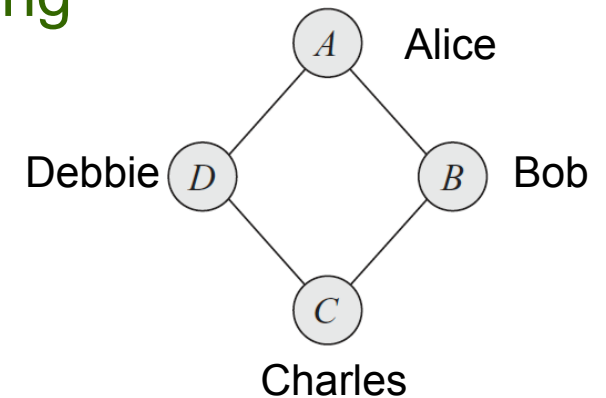
Topics

- Motivation for Learning
- Goals of Learning
 1. Density Estimation
 2. Specific Prediction Tasks
 3. Knowledge Discovery
- Learning as Optimization
 - Empirical Risk and Overfitting
 - Discriminative vs Generative Training
- Learning Tasks
 - Model Constraints
 - Data Observability
 - Taxonomy of Learning Tasks

BN

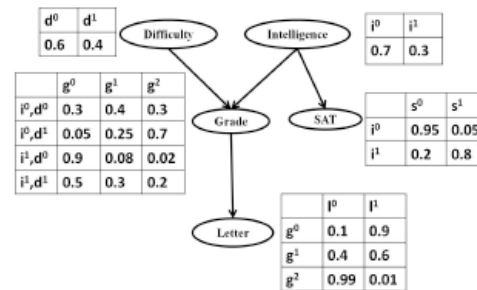


MN



Motivation for Learning

- Usual starting point is a given graphical model
 - Structure and parameters are part of input



- Two approaches to task of acquiring a model
 - Knowledge Engineering
 - Construct a network by hand with expert's help
 - Machine Learning
 - Learn model from a set of instances

Knowledge Engineering vs ML

- Knowledge Engineering Approach
 - Pick variables, pick structure, pick probabilities
 - Too much Effort
 - Simple ones require hours of effort, complex one: months
 - Significant testing of model by evaluating results of typical queries yield plausible answers
- Machine Learning Approach
 - Instances available from distribution we wish to model
 - Easier to get large data sets rather than human expertise

Difficulties with Manual Construction

- In some domains:
 - Amount of knowledge required too large
 - No experts who have sufficient understanding
 - Cost: expert time is valuable
- Properties of distribution change from one site to another
- Change over time
 - Expert cannot redesign every few weeks
- Modeling mistakes have serious impact on quality of answers

Advantage of ML approach

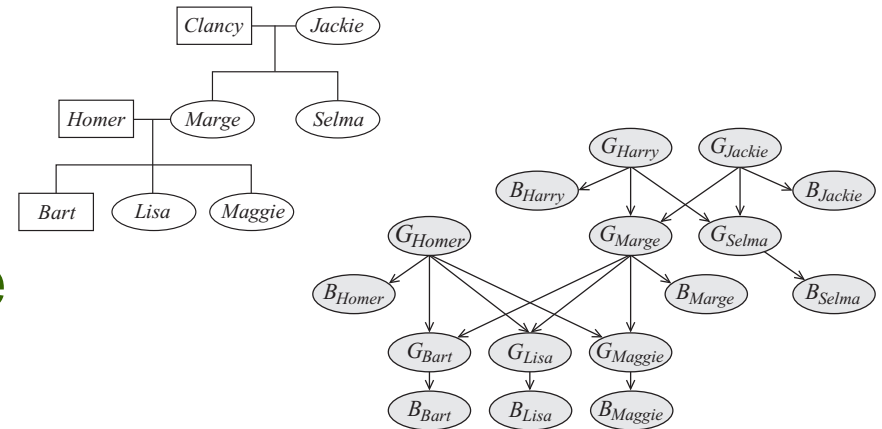
- We are in the Information Age
 - Easier to obtain even large amounts of data in electronic form than to obtain human expertise
- Example Data
 - Medical Diagnosis
 - Patient records
 - Pedigree Analysis (Genetic Inheritance)
 - Family trees for disease transmission
 - Image Segmentation
 - Set of images segmented by a person

Medical Diagnosis Task

- Collection of patient records
 - History:
 - Age, sex, history, medical complications
 - Symptoms
 - Results of tests
 - Diagnosis
 - Treatment
 - Outcome
- Task: Use data to model distribution of patients
 - Pathologist diagnoses disease of lymph nodes (Pathfinder 1992)

Genetic Inheritance: Pedigree

- Set of Family trees
- Task: Learn distribution
 - Breast cancer, Blood Type
- Three types of CPDs:



– Penetrance model: phenotype given genotype

- Probability of a phenotype (say patient's blood type, B) given person's genotype (ordered pair of parent's blood-type, A,B or O): $P(B(c)|G(c))$ $c=person, p=father, m=mother$

– Transmission model: genotype passed to child

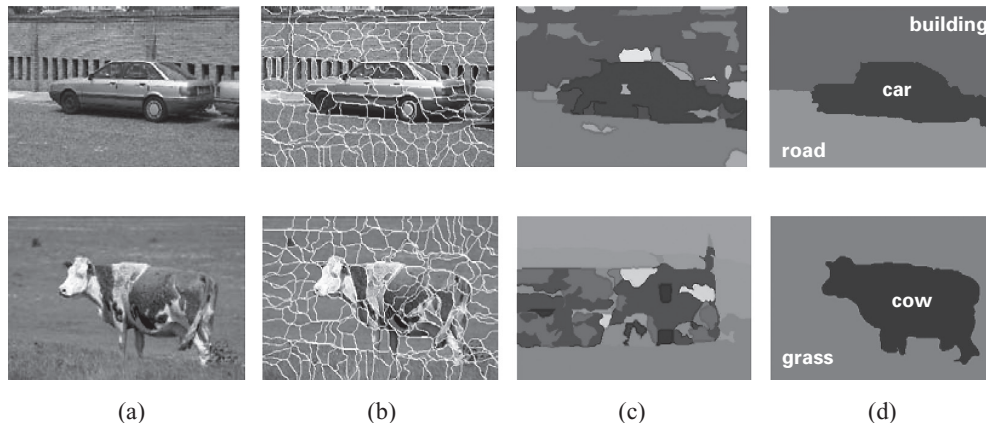
- How often a genotype (locus for disease or blood type) passed from parent to child: $P(G(c)|G(p),G(m))$:

– Genotype Priors: $P(G(c))$

Image Segmentation

- Set of images segmented by a person
 - MRF:
 - Edge potential between neighboring superpixels
 - Penalty λ when $X_i \neq X_j$,
 - Class pairs: Tigers more likely next to vegetation than water
 - Learn parameters of MRF
 - To define characteristics of different regions
 - How strongly neighbor pixels belong to same segment

Original Image | Super Pixels | Node Pots. Only | Pairwise MN



Model Learning

- Construct model for underlying distribution either in isolation or with prior knowledge from human expert
- There are several variants of this task

Model Learning Task

- Domain is governed by distribution P^*
- Induced by a directed/undirected network model $\mathcal{M}^*=(K^*,\theta^*)$: graph structure and parameters
- Given data set $D=\{d[1],\dots,d[M]\}$ of M samples from P^* sampled independently
- Given a family of models, task is to learn some model \tilde{M} from this family that defines a distribution $P_{\tilde{M}}$

Model Learning Variants

1. Given a family of models, learn some model \tilde{M} from this family that defines a distribution $P_{\tilde{M}}$
2. Learn only model parameters for a fixed structure
3. Some or all of the *structure* of the model
4. A probability distribution over models
5. Estimate of our confidence in model learned

Goal of Learning a Probabilistic Model

- Ideal goal of returning a model that completely captures the probability model \mathcal{M} is not achievable
 - Computational complexity
 - Limited data set provides rough approximation of true distribution
- We have to select \tilde{M} that is a best approximation to \mathcal{M}^*
- Notion of best depends on our goals
- Need to define goals and metrics for evaluation

Three Goals of Learning

1. Density Estimation

- Most common reason for learning network model is to use it for probabilistic inference
 - Evaluate probability of a full instance ξ

2. Specific Prediction Tasks

- We may only be interested in special case $P(Y|X)$

3. Knowledge Discovery

- A tool for discovering knowledge about P^*
 - What are the direct/indirect independencies?
 - Nature of dependencies
 - E.g., positive or negative correlation
 - In medical domain, which factors lead to a disease¹⁴

1. Density Estimation

- Construct model \tilde{M} whose distribution \tilde{P} is close to the generating distribution P^*
- How to evaluate quality of approximation \tilde{M} ?
 - Common measure is *Relative Entropy Distance Measure* or *K-L Divergence* defined as
 - If P and Q are distributions over variables $\{X_1, \dots, X_n\}$

$$D(P(X_1, \dots, X_n) \parallel Q(X_1, \dots, X_n)) = E_P \left[\log \frac{P(X_1, \dots, X_n)}{Q(X_1, \dots, X_n)} \right]$$

- Denoting variable settings by ξ

$$D(P^* \parallel \tilde{P}) = E_{\xi \sim P^*} \left[\log \left(\frac{P^*(\xi)}{\tilde{P}(\xi)} \right) \right]$$

- This measure is zero when $\tilde{P} = P^*$ and positive otherwise
- It measures the *compression loss* in bits using \tilde{P} rather than P^*
- Wish to find \tilde{P} for which this metric is low

Evaluation of K-L Divergence

- To evaluate the relative-entropy distance measure

$$D(P^* \parallel P) = E_{\xi \sim P^*} \left[\log \left(\frac{P^*(\xi)}{P(\xi)} \right) \right]$$

$$\xi = [X_1, \dots, X_n]$$

- Need to know P^*
- If learning algorithm is being evaluated over synthetic data we know P^*
- In real world application P^* is not known
- We can simplify this metric to one easier to evaluate (derivation given next)

Simplifying Relative Entropy

- For any distributions P, P' over χ we have

Proposition :

$$D(P \parallel P') = -H_P(X) - E_{\xi \sim P} [\log P'(\xi)]$$

Proof:

$$\begin{aligned} D(P \parallel P') &= E_{\xi \sim P} \left[\log \left(\frac{P(\xi)}{P'(\xi)} \right) \right] \\ &= E_{\xi \sim P} [\log P(\xi) - \log P'(\xi)] \\ &= E_{\xi \sim P} [\log P(\xi)] - E_{\xi \sim P} [\log P'(\xi)] \\ &= -H_P(\chi) - E_{\xi \sim P} [\log P'(\xi)] \end{aligned}$$

- Apply this derivation to P^*, \tilde{P}

Expected Log-likelihood Metric

Relative Entropy Distance Measure has two terms:

$$D(P^* \parallel \tilde{P}) = \underbrace{-H_{P^*}(X)}_{\text{Negative entropy of } P^*} - \underbrace{E_{\xi \sim P^*} [\log \tilde{P}(\xi)]}_{\text{Focus on making this term large}}$$

Negative entropy of P^*
Since it does not depend on P
it can be dropped

Focus on making this term large
Encodes our preference for models that assign
High probability to instances sampled from P^*

Called *Expected log-likelihood*
Whose negative is called *Expected log-loss*
which we want to minimize

Useful to compare two learned models but cannot determine
closeness to optimum
(since first term is dropped)

Log Loss

- In machine learning, we are interested in the *likelihood of the data* given a model M

$$P(D:M)$$

- We replace
- For convenience the *log-likelihood*

$$l(D:M) = \log P(D:M)$$

- Negated form is called log-loss
 - Note that log-likelihood is itself a negative quantity since log-probabilities are negative

Relationship of Log-loss to Likelihood

Loss that model makes on instance ξ is $loss(\xi:M)$

Expected loss(or risk) over D instances is

$$E_D[loss(\xi : M)] = \frac{1}{|D|} \sum_{\xi \in D} loss(\xi : M)$$

which we want to minimize

Relationship to log-likelihood:

Given data set $\mathcal{D} = \{\xi[1], \dots, \xi[M]\}$ Probability that \mathcal{M} ascribes to \mathcal{D} is

$$P(\mathcal{D} : \mathcal{M}) = \prod_{m=1}^M P(\xi[m] : \mathcal{M})$$

$$\log P(\mathcal{D} : \mathcal{M}) = \sum_{m=1}^M \log P(\xi[m] : \mathcal{M}) \quad \text{which we want to maximize}$$

Compare with AIC: $2k - 2\ln L$

k = no of parameters,

L = maximum likelihood

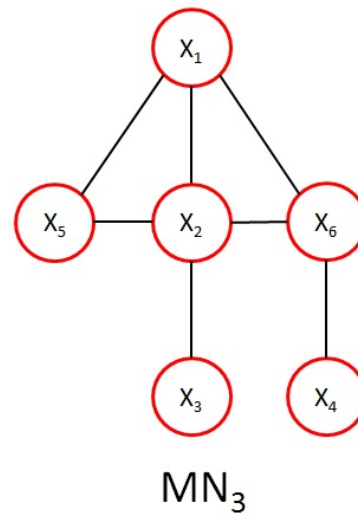
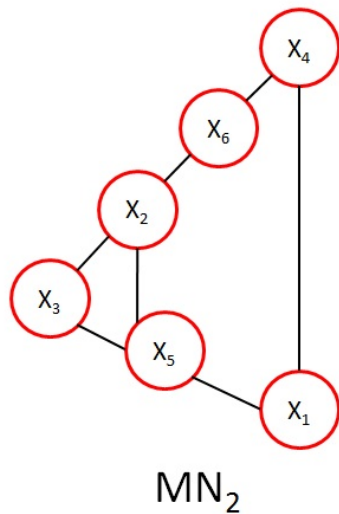
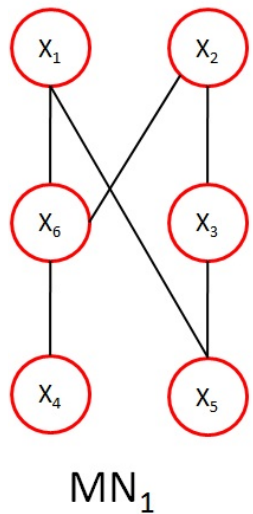
M = no of samples

BIC: $-2\ln L + k\ln(M)$

Expected (empirical) log-loss:

$$\frac{-1}{|D|} \sum_{m=1}^M \log P(\xi[m] : \mathcal{M})$$

Example of log-likelihood metric



MRF	average log-likelihood
MN_1	-2803.1
MN_2	-2842.7
MN_3	-2792.1

We want to make log-likelihood large. So MN₃ is best.

2. Specific Prediction Tasks

- In density estimation goal was to use learning model for probabilistic inference
 - Jumped to conclusion we wish to fit P^* well
 - To evaluate joint probability of ξ
- We are only interested in answering queries of the form $P(Y|X)$
- Example: classification problems
 - Topic categorization (Single decision)
 - X consists of words and other features, Y is topic
 - Image segmentation (Multiple decisions)
 - Predict labels for all pixels (Y) given image features (X)

Model Trained for Prediction Task

- For any instance x produce the probability distribution $P(Y|x)$
- Or MAP assignment of this conditional distribution to produce a specific prediction

$$h_p(x) = \arg \max_y P(y|x)$$

Value of y that has the highest probability

- What loss functions to use?
 - Classification error
 - Hamming loss

Classification Error

- Also called 0/1 loss

$$E_{(x,y) \sim P}[I\{h_P(x) \neq y\}]$$

- Expected value of I is its average value
 - Determined by adding all values of I over different (x,y) and dividing by the total number
- Probability that over all pairs (x,y) sampled from P our classifier selects the wrong label
- Suitable for labelling a single sample, not suited for many simultaneous decisions

Hamming Loss

- Simultaneously provide labels to a large number of inputs, e.g., image segmentation
- Do not wish to penalize the entire prediction if classification is wrong on few
- Instead of using $I(h_P(x) \neq y)$
 - count number of variables Y in which $h_P(x)$ differs from ground-truth y
 - Or take confidence of prediction into account
 - Conditional log-likelihood

Conditional log-likelihood Criterion

- Taking into account prediction confidence

$$E_{(x,y) \sim P^*} [\log P(y | x)]$$

- Similar to expected log-loss but use conditional distribution instead
- In general, if the model is never used to predict X , we want to design our training to optimize the quality of answers Y

3. Knowledge Discovery

- Very different motivation to learn model for P^*
- A tool for discovering knowledge about P^*
 - What are the direct/indirect independencies?
 - What characterizes the nature of dependencies?
 - E.g., positive or negative correlation
 - In genetic inheritance domain
 - Discover parameter for inheritance of certain property
 - Parameter gives insight regarding disease alleles
 - In medical domain
 - Which factors lead to a disease
 - Which symptoms associated with different diseases

Statistical Correlation vs Learned Network

- Simpler statistical methods can be used
 - Highlight most significant correlations between pairs of variables
 - Variables being uncorrelated can be determined using chi-square test using $P(x,y)-P(x)P(y)$
- Learned network models
 - Provide parameters that have direct causal interpretation
 - Reveal much finer structure
 - Distinguish between direct and indirect independencies
 - Both lead to correlations in resulting distribution

Evaluation of Knowledge Discovery

- Different compromises from prediction task
- We do care about the model \mathcal{M}^* rather than a model \mathcal{M} that induces a distribution similar to that of \mathcal{M}^*
- In contrast to density estimation where metric is on distribution defined by model. i.e., $D(P^*||P)$
 - Here measure of success is on the model
 - Differences between \mathcal{M}^* and \mathcal{M}
- Goal often not achievable

Difficulties in Knowledge Discovery

- Identifiability

- Even with large amounts of data

- A Bayesian network has several I-equivalent structures
 - Best hope to recover \mathcal{K}^* is to get I-equivalent structure

- With limited data

- If X, Y directly related in \mathcal{K}^* but parameters induce weak relationship, cannot distinguish from random fluctuations
 - Fewer repercussions in density estimation

Significance of Knowledge Discovery

- High probability of model identification errors
- Expensive wet-lab experiments
 - To discover relationships between genes
- In knowledge discovery, more critical to assess confidence in a prediction