# CPSC 440: Advanced Machine Learning
## More DAGs 3

Mark Schmidt

University of British Columbia

Winter 2021

## Last Time: Learning in DAG Models

- Learning in DAGs involves fitting each $p(x_j \mid x_{\mathsf{pa}(j)})$:
    1. Set $\bar{y}^i = x_j^i$ and $\bar{x}^i = x_{\mathsf{pa}(j)}^i$.
    2. Solve a supervised learning problem using $\{\bar{X}, \bar{y}\}$.
        - Gives you a model of $p(x_j \mid x_{\mathsf{pa}(j)})$.
- Combine the $d$ regression/classification models as the density estimator.
    - We've turned unsupervised learning into supervised learning.

- We can use our usual tricks:
    - Linear models, non-linear bases, regularization, kernel trick, random forests, etc.
    - With least squares for continuos $x_j$ it's called a Gaussian belief network.
    - With logistic regression for binary $x_j$ it's called a sigmoid belief networks.
    - Don't need Markov assumptions to tractably fit these models.

# DAGs: Big Picture

- Setting the parameters of a DAG model:
    - Get the graph from an expert, or learn the graph (later).
    - Given the conditional probabilities from an expert, or learn them from data.
        - Counting if you use general discrete distribution for conditionals.
        - Supervised learning for conditions.
        - Combine either of the above with EM if you have hidden/missing values.


- Inference in DAG models:
    - Can use Monte Carlo approximations with ancestral sampling:
        - Sample $x_1$ from $p(x_1)$, $x_2$ from $p(x_2 \mid x_{\mathsf{pa}(2)})$, $x_3$ from $p(x_3 \mid x_{\mathsf{pa}(3)})$,...


    - Can use dynamic programming for exact inference with discrete $x_j$.
        - Also works if all $p(x_j \mid x_{\mathsf{pa}(j)})$ are Gaussian.
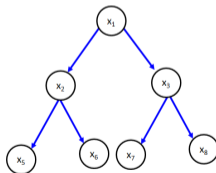        - But dynamic programming may be too expensive (today).

# Inference in Forest DAGs ("Belief Propagation")

- If we try to generalize the CK equations to DAGs we obtain

$$p(x_j = s) = \sum_{x_{\mathsf{pa}(j)}} p(x_j = s, x_{\mathsf{pa}(j)}) = \sum_{x_{\mathsf{pa}(j)}} \underbrace{p(x_j = s \mid x_{\mathsf{pa}(j)})}_{\text{given}} p(x_{\mathsf{pa}(j)}).$$

which works if each node has at most one parent.
  - Such graphs are called trees (connected), or forests (disconnected).
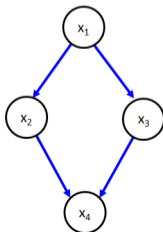    - Also called "singly-connected".



- Forests allow efficient dynamic programming methods as in Markov chains.
  - In particular, decoding and univariate marginals/conditionals in $O(dk^2)$.
  - Forward-backward applied to tree-structured graphs is called belief propagation.

# Inference in General DAGs

- If we try to generalize the CK equations to DAGs we obtain

$$p(x_j = s) = \sum_{x_{\mathsf{pa}(j)}} p(x_j = s, x_{\mathsf{pa}(j)}) = \sum_{x_{\mathsf{pa}(j)}} \underbrace{p(x_j = s \mid x_{\mathsf{pa}(j)})}_{\text{given}} p(x_{\mathsf{pa}(j)}).$$

- What goes wrong if nodes have multiple parents?
  - The expression $p(x_{\mathsf{pa}(j)})$ is a joint distribution depending on multiple variables.

- Consider the non-tree graph:

## Inference in General DAGs

- We can compute $p(x_4)$ in this non-tree using:

$$p(x_4) = \sum_{x_3} \sum_{x_2} \sum_{x_1} p(x_1, x_2, x_3, x_4)$$

$$= \sum_{x_3} \sum_{x_2} \sum_{x_1} p(x_4 \mid x_2, x_3) p(x_3 \mid x_1) p(x_2 \mid x_1) p(x_1)$$

$$= \sum_{x_3} \sum_{x_2} p(x_4 \mid x_2, x_3) \underbrace{\sum_{x_1} p(x_3 \mid x_1) p(x_2 \mid x_1) p(x_1)}_{M_{23}(x_2, x_3)}$$

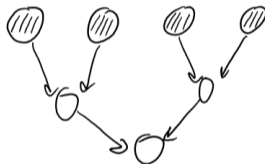- Dependencies between $\{x_1, x_2, x_3\}$ mean our message depends on two variables.

$$p(x_4) = \sum_{x_3} \sum_{x_2} p(x_4 \mid x_2, x_3) M_{23}(x_2, x_3)$$

$$= \sum_{x_3} M_{34}(x_3, x_4),$$

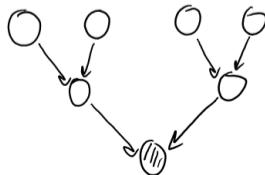# Inference in General DAGs

- With $2$-variable messages, our <span style="color:red">cost increases</span> to $O(dk^3)$.

- If we add the edge $x_1-> x_4$, then the cost is $O(dk^4)$.

  (the same cost as enumerating all possible assignments)

- Unfortunately, cost is <span style="color:red">not as simple as counting number of parents</span>.
  - Even if each node has $2$ parents, we may need huge messages.
  - Decoding is NP-hard and computing marginals is #P-hard in general.

  - We'll see later that maximum message size is "<span style="color:blue">treewidth</span>" of a particular graph.

- On the other hand, <span style="color:green">ancestral sampling is easy</span>:
  - We can obtain Monte Carlo estimates of solutions to these NP-hard problems.

# Conditional Sampling in DAGs

- What about conditional sampling in DAGs?
  - Could be easy or hard depending on what we condition on.
- For example, easy if we condition on the first variables in the order:
  - Just fix these and run ancestral sampling.



- Hard to condition on the last variables in the order:
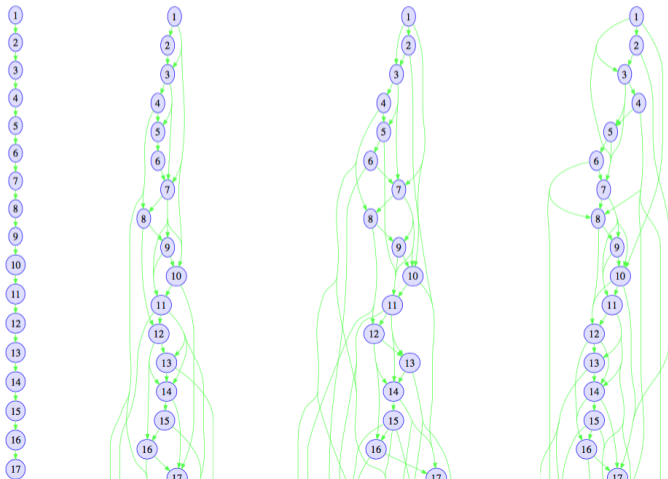  - Conditioning on descendent makes ancestors dependent.

# DAG Structure Learning

- Structure learning is the problem of choosing the graph.
  - Input is data $X$.
  - Output is a graph $G$.

- The "easy" case is when we're given the ordering of the variables.
  - So the parents of $j$ must be chosen from $\{1, 2, \ldots, j-1\}$.

- Given the ordering, structure learning reduces to feature selection:
  - Select features $\{x_1, x_2, \ldots, x_{j-1}\}$ that best predict "label" $x_j$.
  - We can use any feature selection method to solve these $d$ problems.

# Example: Structure Learning in Rain Data Given Ordering

- Structure learning in rain data using L1-regularized logistic regression.
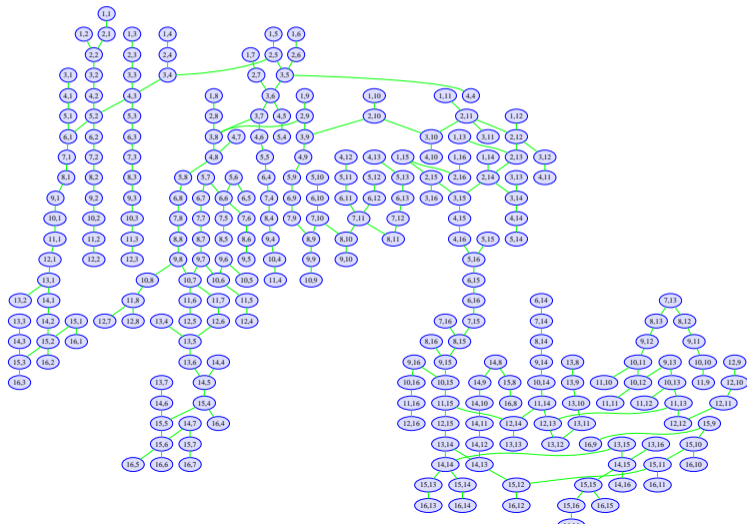  - For different $\lambda$ values, assuming chronological ordering.

# DAG Structure Learning without an Ordering

- Without an ordering, a common approach is "search and score"
  - Define a score for a particular graph structure (like BIC or other L0-regularizers).
  - Search through the space of possible DAGs.
    - "DAG-Search": at each step greedily add, remove, or reverse an edge.

- May have equivalent graphs with the same score (don't trust edge direction).
  - Do not interpret causally a graph learned from data.

- Structure learning is NP-hard in general, but finding the optimal tree is poly-time:
  - For symmetric scores, can be found by minimum spanning tree ("Chow-Liu").
    - Score is symmetric if $\text{score}(x_j \to x_{j'})$ is the same as $\text{score}(x_{j'} \to x_j)$.
  - For asymetric scores, can be found by minimum spanning arborescence.

# Structure Learning on USPS Digits

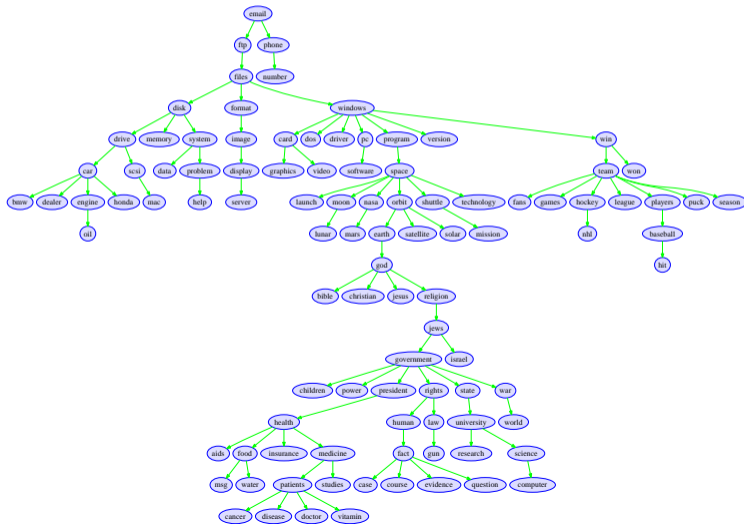An optimal tree on USPS digits (16 by 16 images of digits).

## 20 Newsgroups Data

- Data containing presence of 100 words from newsgroups posts:

| car | drive | files | hockey | mac | league | pc | win |
|-----|-------|-------|--------|-----|--------|----|----|
| 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 |
| 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 |

- Structure learning should give some relationship between word occurrences.

# Structure Learning on News Words
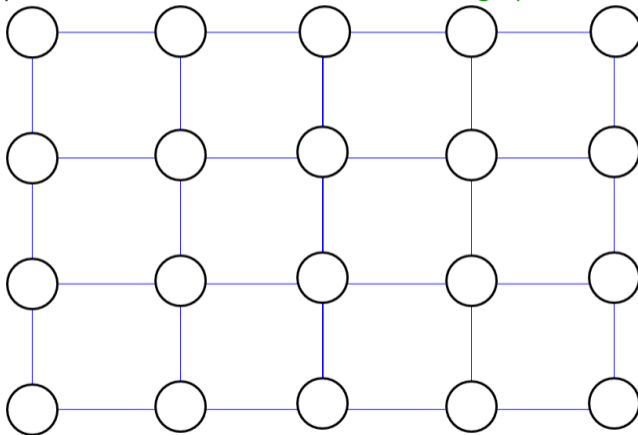
Optimal tree on newsgroups data:

# Outline

# Directed vs. Undirected Models

- In some applications we have a natural ordering of the $x_j$.
  - In the "rain" data, the past affects the future.

- In some applications we don't have a natural order.
  - E.g., pixels in an image.

- In these settings we often use undirected graphical models (UGMs).
  - Also known as Markov random fields (MRFs) and originally from statistical physics.
    - Another name is "Markov networks".

# Directed vs. Undirected Models

- Undirected graphical models are based on undirected graphs:



- They are a classic way to model dependencies in images:
  - Can capture dependencies between neighbours without imposing an ordering.

# Multi-Label Classification

- Consider multi-label classification:



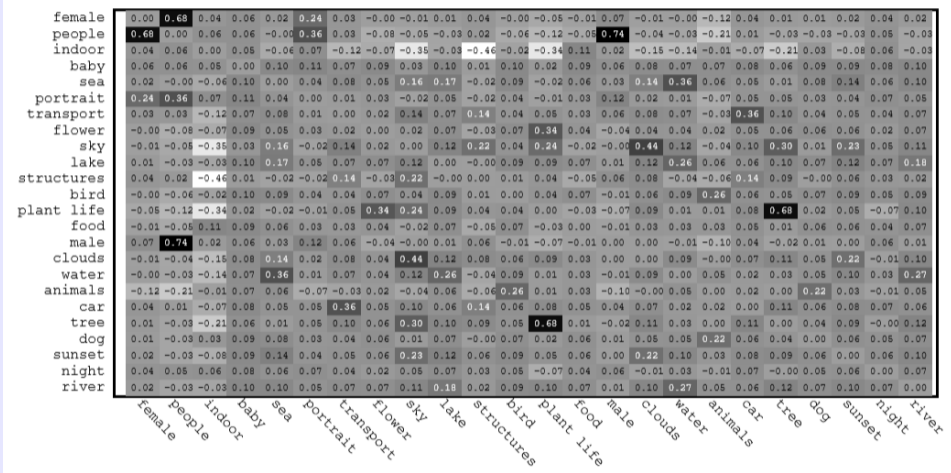| female/indoor/portrait | sky/plant life/tree | water/animals/sea | animals/dog/indoor | indoor/flower/plant life |

http://proceedings.mlr.press/v37/chenb15.pdf

- Flickr dataset: each image can have multiple labels (out of 38 possibilities).

- Use neural networks to generate "factors" in an undirected model.
  - Decoding undirected model makes predictions accounting for label correlations.
  - We'll discuss how neural networks and density models fit together later.

# Multi-Label Classification

- Learned correlation matrix:

# Summary

- Inference in DAGs:
    - Ancestral sampling and Monte Carlo methods work as before.
    - Message-passing message sizes depend on graph structure.

- Structure learning is the problem of learning the graph structure.
    - Hard in general, but easy for trees and L1-regularization gives fast heuristic.

- Undirected graphical models do not require an ordering of the variables.

- Next time: easy conditional dependence and hard "everything else" in UGMs.

# "Constraint-Based" DAG Structure Learning

- Another common structure learning approach is "constraint-based":
  - Based on performing a sequence of conditional independence tests.
  - Prune edge between $x_i$ and $x_j$ if you find variables $S$ making them independent,

$$x_i \perp x_j \mid x_S.$$

  - Challenge is considering exponential number of sets $x_S$ (heuristic: "PC algorithm").
  - Assumes "faithfulness" (all independences are reflected in graph).
    - Otherwise it's weird (a duplicated feature would be disconnected from everything.)