

Machine Learning (6.867)

David Sontag

Massachusetts Institute of Technology

Lecture 15, Nov. 2, 2017

Course announcements

Project milestone 4 (*produce initial results*) due Tuesday

Where are we in the course?

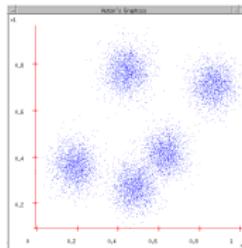
- ➊ Until now, we focused on *supervised learning*
 - Regression, classification
 - Loss functions, regularization
 - Linear, SVMs, nearest neighbor, neural networks, decision trees

- ➋ Today's lecture: *probabilistic modeling*
- ➌ 11/7 – 11/14: using probabilistic models for *unsupervised learning*
- ➍ 11/16: *approximate inference* in probabilistic models (Gibbs sampling)

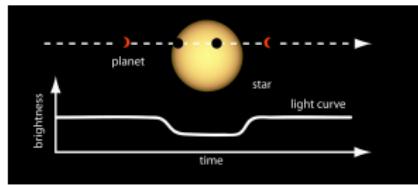
- ➎ 11/21 and 12/5: *dimensionality reduction* with matrix factorization
- ➏ 11/28: *recommendation systems* (not quite (un)supervised learning)
- ➐ 12/7 and 12/12: *reinforcement learning*

How do we tackle problems like these?

- How many clusters are there in this data?



- Does there exist an exoplanet here?



- How do we learn to predict y in the presence of missing data?

x	y
0 1 ?? 1 1 0 ?	1
1 ? 1 0 ?? 0 0	0
1 0 0 0 0 0 1 0	?
1 ?? ??? ? 0	0

Probabilistic modeling

① **Represent** the world as a collection of random variables

- Observed variables X_1, \dots, X_n
- Latent variables Z_1, \dots, Z_m (*optional*)

with joint distribution $p(X_1, \dots, X_n, Z_1, \dots, Z_m)$.

② **Learn** the distribution from data

③ Either:

- Obtain insight from the distribution itself (e.g. the learned parameters)
- Use learned model for **prediction** by performing posterior **inference**,

$$f(x) = p(Z, X_{\text{query}} \mid X_{\text{evidence}} = x)$$

where $f(x)$ is a vector-valued function which returns the posterior distribution given observed evidence

Applications: Modeling sequential data (*today*)

Machine translation, speech recognition, time series, ...

The screenshot shows the Google Translate interface. At the top, it says "Translate" and has dropdown menus for "From: English" and "To: Spanish". A "Translate" button is also present. Below this, there are two text boxes. The left text box contains the following English text:

The top U.S. general, visiting Israel at a delicate and dangerous moment in the global standoff with Tehran, is expected to press for restraint amid fears that the Jewish state is nearing a decision to attack Iran's nuclear program.

The right text box contains the corresponding Spanish translation:

El máximo general de EE.UU., de visita en Israel en un momento delicado y peligroso en el enfrentamiento global con Teherán, se espera que presione a la moderación en medio de temores de que el estado judío se acerca a una decisión de atacar el programa nuclear de Irán.

Below the text boxes, there is a note: "New! Hold down the shift key, click, and drag the words above to reorder. Dismiss". At the bottom of the interface, there are links for "Turn off instant translation", "About Google Translate", "Mobile", "Privacy", "Help", and "Send feedback".

Recurrent neural networks, hidden Markov models, Gaussian processes

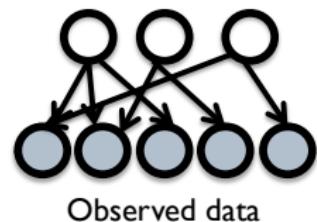
Applications: Modeling survey data (*next week*, 11/2)

Factor analysis

1 0 0 1 1 0 0 0 1
1 1 0 1 1 0 0 1 0
0 1 0 1 0 0 1 0 0
1 0 0 0 1 0 0 1 1
0 1 0 1 0 0 1 1 0



Factors



- Social surveys with questions such as:
 - Should the tax rate be progressive?
 - Do you consider Edward Snowden a hero or a criminal?
 - What is your annual salary?
 - Should affirmative action be used in college admissions?
- Automatically *discover* the relevant hidden variables or factors, e.g.
 - Socioeconomic status
 - Health
 - Political attitudes
 - Family values
- *Infer* the underlying beliefs of each respondent

Applications: Modeling text data (next week, 11/9)

Topic models

Poisoning by ice cream.

No chemist certainly would suppose that the same poison exists in all samples of ice-cream which have produced untoward symptoms in man. Mineral poisons, copper, lead, arsenic, and mercury, have all been found in ice cream. In some instances these have been used with criminal intent; in other cases they have simply been accidentally. Likewise, that vanilla is sometimes the bearer, at least, of the poison, is well known to all chemists. Mr. Bartery's idea that the poisonous properties of the cream which he examined were due to gelatin is not only a natural theory. The poison's origin might in this case arise from the decomposition of the gelatine; or with the gelatine there may be introduced into the milk a ferment, by the growth of which a poison is produced.

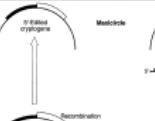
But the cream which I examined, none of the above sources of the poisoning existed. There were no mineral poisons present. No gelatine of any kind had been used in making the cream. The vanilla used was shown to be not poisonous. The slugs were made up of two chemical catalysts which might not have been conclusive, but Mr. Novis and I drank of the vanilla extract which was used, and no ill results followed. Still, from this cream we isolated the same poison which I had before found in poisonous cheese (*Zeitschrift für physiologische Chemie*, x,

RNA Editing and the Evolution of Parasites

Larry Simpson and Dmitri A. Maslov

The *Leucaspis* gallflies, together with their close group of endoparasitoids, represent the earliest known lineage of eukaryotes to edit their own RNA. Although these insects are common, there are two major groups, the polytrophic holometabolous species and the parasitic species, and the latter are nonhomoplasmic, which are obligate parasites [1]. Perhaps because of the unique nature of the transcript editing process, the parasite cell can possess several distinct genetic lineages. We report here accompanying phylogenetic analysis of a 134 kb segment of the *Leucaspis* genome in which RNA editing of all 25 tRNAs is required for normal protein synthesis. The RNA editing function [3–7] causes open reading frames for optional deletion of amino acids at specific sites within the coding regions of mRNA. This pattern of RNA editing is unique and specific among insect genomes thus far described (RNA “splicing”). The

rate is thus a diagnostic of the status of the primary parasitic host. The “virulence” (or) model [1,2] states that virulence is the result of a trade-off between host and parasite, such that the parasite will have to “trade off” its own benefit for that of the host [8]. This theory, together with the lack of evidence for homoplasmic species, suggests that, in this theory, digenetic life cycles (allowing both protein and viral evolution) could result in complex dynamics of behavior. Natural selection for virulence reduction could occur via the acquisition by some hermaphrodites of the ability to feed on the blood



Chaotic Beetles

Charles Godfray and Michael Hassell

Epidemiologists have been the pioneer work of May in the mid-1970s [1,2] that drug population dynamics of arthropod and plant diseases in field populations can exhibit chaotic attractors from two sources. The first is a population of organisms that constitutes any natural community; the second is the parasite, virus or pathogen for species or insects, both able to affect disease dynamics in field populations. The nonlinear feedback processes that result in complex dynamics behavior. Natural population can share nonlinear oscillations dynamics that are induced by the combined effects by parasite sensitivity to initial conditions with the time delay, when this would have important role in establishing the population dynamics of disease [3,4]. Ode page 209 of this issue, Gottsche et al. [2,3] provide the most

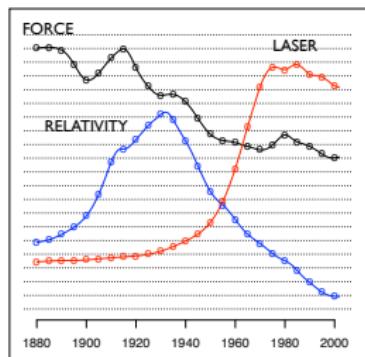
interesting dynamics induced by complex dynamics and chaos as a biological population of different components in a field population of disease (see figure).

It has been estimated that the number of chaotic attractors in epidemics in the field population may be up to 100 million [1,2], which fluctuating populations in the field. It is also a possibility that a single population may apparently remain stable, while in fact it fluctuates randomly, all induced by the normal nonlinearities of the system of equations of the model. Given a long enough time to record the data of the system, then nonlinear mathematics can be used to obtain a better understanding of the system. Given the number of chaotic attractors, one can make predictions of disease. For example, chaotic trajectories coming from different initial conditions are in different geometric shapes with different periods. One can predict the behavior of the disease if the initial conditions are allowed in a mathematical model.



THE AUTHOR AND CO-DIRECTOR OF THE INSTITUTE OF CHAOS, READING COLLEGE, INSTITUTE PARK, ABERDEEN, SCOTLAND, UK; E-MAIL: M.HASSELL@AOL.COM

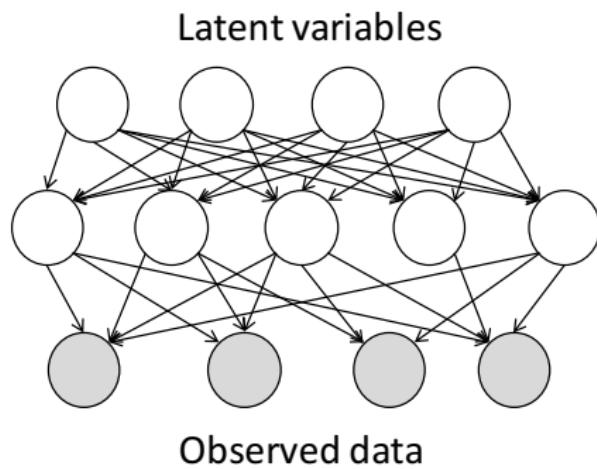
"Theoretical Physics"



human genome	evolutionary	disease	computer models
dna	species	host	information
genetic	organisms	bacteria	data
genes	life	diseases	computers
sequence	origin	resistance	system
gene	biology	bacterial	network
molecular	groups	new	systems
sequencing	phylogenetic	strains	model
map	living	control	parallel
information	diversity	infectious	methods
genetics	group	malaria	networks
mapping	new	parasite	software
project	two	parasites	new
sequences	common	united	simulations
		tuberculosis	

Applications: Modeling images and high-dimensional data (in two weeks, 11/14)

Deep generative models and variational autoencoders, adversarial nets



Samples from the model:

6	6	9	0	5	7	0	8	0	5
1	5	7	6	6	9	5	1	3	0
4	1	4	5	5	4	0	6	4	9
8	9	7	7	2	3	0	7	4	8
6	5	4	0	0	9	9	2	2	3
0	9	8	6	1	5	0	7	7	6
6	6	2	9	7	6	9	4	0	9
2	0	1	3	4	1	5	4	4	0
1	2	4	7	6	1	9	5	3	7
6	4	3	8	7	9	0	9	4	5

Key challenges

- ① **Represent** the world as a collection of random variables with joint distribution $p(X_1, \dots, X_n, Z_1, \dots, Z_m)$
 - How does one *compactly describe* this joint distribution?
 - Directed graphical models (Bayesian networks)
 - Undirected graphical models (Markov random fields, factor graphs)
not covered in this class – see 6.438
- ② **Learn** the distribution from data
 - Maximum likelihood estimation. Other estimation methods?
 - How much data do we need?
 - How much computation does it take?
- ③ **Perform inference** to compute posterior distribution
 $p(Z | X_1 = x_1, \dots, X_n = x_n)$
 - Exactly doing this is typically computationally intractable
 - What are computationally efficient *approximate* inference algorithms?

Example: Medical diagnosis

- Variable for each **symptom** (e.g. “fever”, “cough”, “fast breathing”, “shaking”, “nausea”, “vomiting”)
- Variable for each **disease** (e.g. “pneumonia”, “flu”, “common cold”, “bronchitis”, “tuberculosis”)
- Diagnosis is performed by **inference** in the model:

$$p(\text{pneumonia} = 1 \mid \text{cough} = 1, \text{fever} = 1, \text{vomiting} = 0)$$

- One famous model, Quick Medical Reference (QMR-DT), has 600 diseases and 4000 findings

Representing the distribution

- Naively, could represent multivariate distributions with table of probabilities for each outcome (assignment)
- How many outcomes are there in QMR-DT? 2^{4600}
- **Estimation** of joint distribution would require a huge amount of data
- **Inference** of conditional probabilities, e.g.

$$p(\text{pneumonia} = 1 \mid \text{cough} = 1, \text{fever} = 1, \text{vomiting} = 0)$$

would require summing over exponentially many variables' values

- Moreover, defeats the purpose of probabilistic modeling, which is to make predictions with *previously unseen observations*

Structure through independence

- If X_1, \dots, X_n are independent, then

$$p(x_1, \dots, x_n) = p(x_1)p(x_2) \cdots p(x_n)$$

- 2^n entries can be described by just n numbers (if $|\text{Val}(X_i)| = 2$)!
- However, this is not a very *useful* model – observing a variable X_i cannot influence our predictions of X_j
- If X_1, \dots, X_n are *conditionally independent* given Y , denoted as $X_i \perp \mathbf{X}_{-i} \mid Y$, then

$$\begin{aligned} p(y, x_1, \dots, x_n) &= p(y)p(x_1 \mid y) \prod_{i=2}^n p(x_i \mid x_1, \dots, x_{i-1}, y) \\ &= p(y)p(x_1 \mid y) \prod_{i=2}^n p(x_i \mid y). \end{aligned}$$

- This is a simple, yet *powerful*, model

Example: naive Bayes for classification

- Classify e-mails as spam ($Y = 1$) or not spam ($Y = 0$)
 - Let $1 : n$ index the words in our vocabulary (e.g., English)
 - $X_i = 1$ if word i appears in an e-mail, and 0 otherwise
 - E-mails are drawn according to some distribution $p(Y, X_1, \dots, X_n)$
- Suppose that the words are conditionally independent given Y . Then,

$$p(y, x_1, \dots, x_n) = p(y) \prod_{i=1}^n p(x_i | y)$$

Estimate the model with maximum likelihood. **Predict** with:

$$p(Y = 1 | x_1, \dots, x_n) = \frac{p(Y = 1) \prod_{i=1}^n p(x_i | Y = 1)}{\sum_{y=\{0,1\}} p(Y = y) \prod_{i=1}^n p(x_i | Y = y)}$$

- Are the independence assumptions made here reasonable?
- Philosophy: Nearly all probabilistic models are “wrong”, but many are nonetheless useful

Bayesian networks

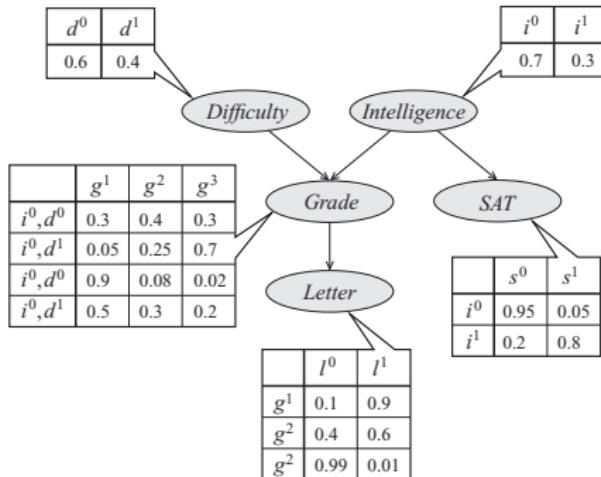
- A **Bayesian network** is specified by a directed *acyclic* graph $G = (V, E)$ with:
 - ① One node $i \in V$ for each random variable X_i ;
 - ② One conditional probability distribution (CPD) per node, $p(x_i | \mathbf{x}_{\text{Pa}(i)})$, specifying the variable's probability conditioned on its parents' values
- Corresponds 1-1 with a particular factorization of the joint distribution:

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i | \mathbf{x}_{\text{Pa}(i)})$$

- Powerful framework for designing *algorithms* to perform probability computations
- Enables use of *prior knowledge* to specify (part of) model structure

Example

- Consider the following Bayesian network:



- What is its joint distribution?

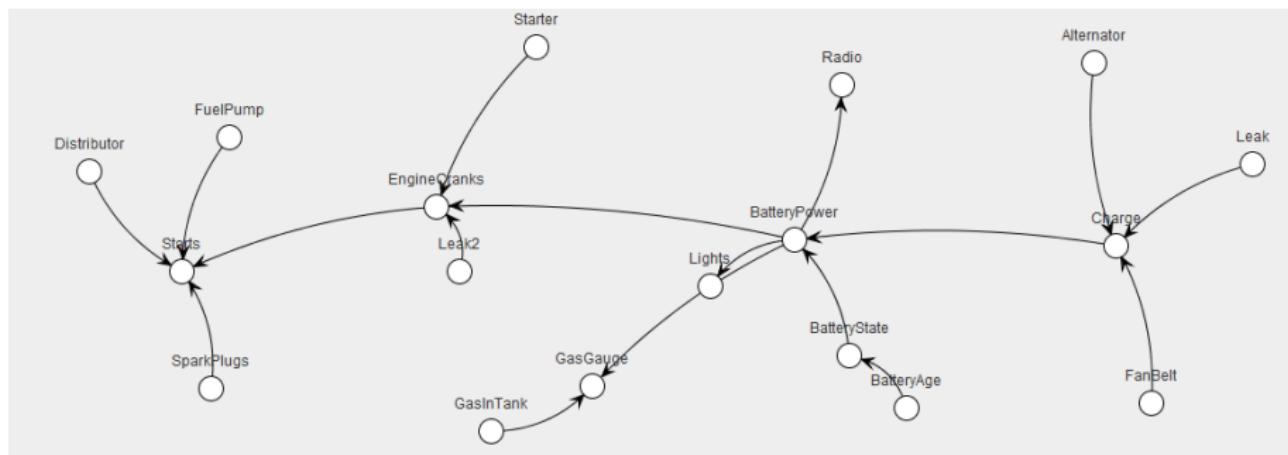
$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i | \mathbf{x}_{\text{Pa}(i)})$$

$$p(d, i, g, s, l) = p(d)p(i)p(g | i, d)p(s | i)p(l | g)$$

More examples

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i | \mathbf{x}_{\text{Pa}(i)})$$

Will my car start this morning?



Heckerman et al., Decision-Theoretic Troubleshooting, 1995

More examples

$$p(x_1, \dots x_n) = \prod_{i \in V} p(x_i | \mathbf{x}_{\text{Pa}(i)})$$

What is the differential diagnosis?

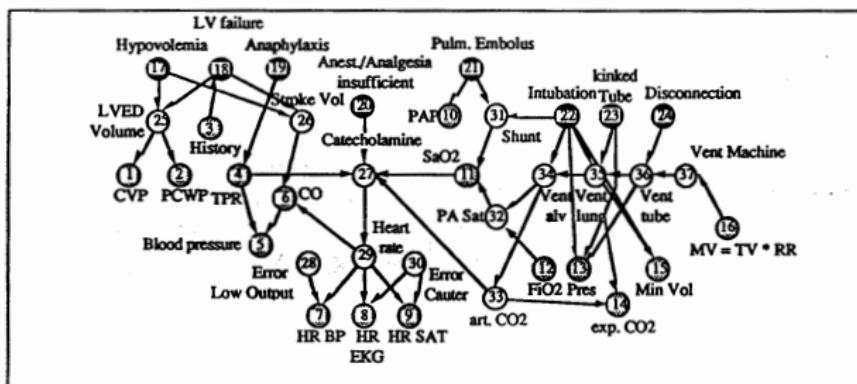
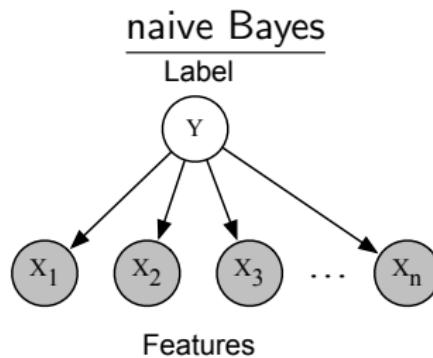


Fig. 1 The ALARM network representing causal relationships is shown with diagnostic (●), intermediate (○) and measurement (◎) nodes. CO: cardiac output, CVP: central venous pressure, LVED volume: left ventricular end-diastolic volume, LV failure: left ventricular failure, MV: minute ventilation, PA Sat: pulmonary artery oxygen saturation, PAP: pulmonary artery pressure, PCWP: pulmonary capillary wedge pressure, Pres: breathing pressure, RR: respiratory rate, TPR: total peripheral resistance, TV: tidal volume

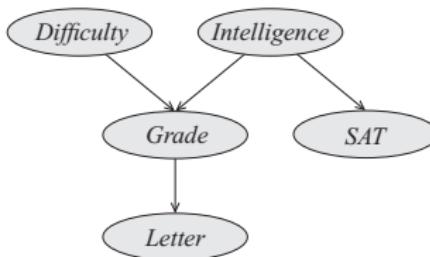
Beinlich et al., The ALARM Monitoring System, 1989

Bayesian networks are *generative models*



- Evidence is denoted by shading in a node
- Can interpret Bayesian network as a **generative process**. For example, to *generate* an e-mail, we
 - ➊ Decide whether it is spam or not spam, by sampling $y \sim p(Y)$
 - ➋ For each word $i = 1$ to n , sample $x_i \sim p(X_i | Y = y)$

Bayesian network structure implies conditional independencies!



- The joint distribution corresponding to the above BN factors as

$$p(d, i, g, s, l) = p(d)p(i)p(g | i, d)p(s | i)p(l | g)$$

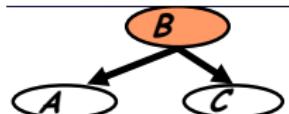
- However, by the chain rule, *any* distribution can be written as

$$p(d, i, g, s, l) = p(d)p(i | d)p(g | i, d)p(s | i, d, g)p(l | g, d, i, g, s)$$

- Thus, we are assuming the following additional independencies:
 $D \perp I$, $S \perp \{D, G\} | I$, $L \perp \{I, D, S\} | G$. What else?

Bayesian network structure implies conditional independencies!

- Generalizing the above arguments, we obtain that a variable is independent from its non-descendants given its parents

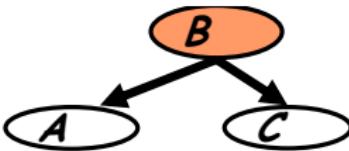


- Common parent** – fixing B *decouples* A and C
- Cascade** – knowing B *decouples* A and C



- V-structure** – Knowing C *couples* A and B
 - This important phenomena is called **explaining away** and is what makes Bayesian networks so powerful

A simple justification (for common parent)



We'll show that $p(A, C | B) = p(A | B)p(C | B)$ for any distribution $p(A, B, C)$ that factors according to this graph structure, i.e.

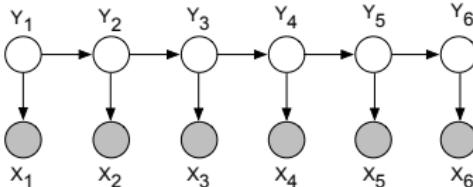
$$p(A, B, C) = p(B)p(A | B)p(C | B)$$

Proof.

$$p(A, C | B) = \frac{p(A, B, C)}{p(B)} = p(A | B)p(C | B)$$



Hidden Markov models

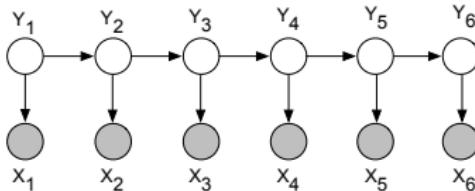


- Widely used in control theory and signal processing. Later applications included speech recognition and computational biology
- Joint distribution factors as:

$$p(\mathbf{y}, \mathbf{x}) = p(y_1)p(x_1 | y_1) \prod_{t=2}^T p(y_t | y_{t-1})p(x_t | y_t)$$

- $p(y_1)$ is the distribution for the starting state
- $p(y_t | y_{t-1})$ is the *transition* probability between any two states
- $p(x_t | y_t)$ is the *emission* probability
- What are the conditional independencies here? For example,
 $Y_1 \perp \{Y_3, \dots, Y_6\} | Y_2$

Hidden Markov models



- Joint distribution factors as:

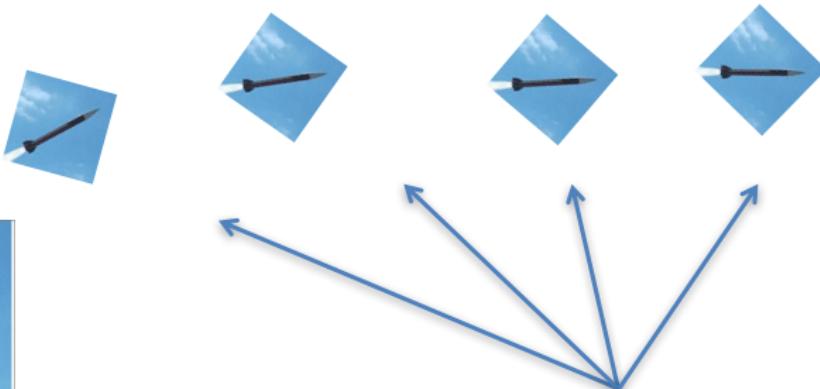
$$p(\mathbf{y}, \mathbf{x}) = p(y_1)p(x_1 | y_1) \prod_{t=2}^T p(y_t | y_{t-1})p(x_t | y_t)$$

- A **homogeneous** HMM uses the same parameters (β and α below) for each transition and emission distribution (**parameter sharing**):

$$p(\mathbf{y}, \mathbf{x}) = p(y_1)\alpha_{x_1, y_1} \prod_{t=2}^T \beta_{y_t, y_{t-1}}\alpha_{x_t, y_t}$$

How many parameters need to be learned?

Example application of HMMs: Tracking



Observe noisy measurements of
missile location: X_1, X_2, \dots



Radar

Where is the missile **now**? Where will it be in 10 seconds?

Inference

- To find out where the missile is *now*, we do **marginal inference**:

$$p(y_T \mid x_1, \dots, x_T)$$

- To find the most likely *trajectory*, we do **MAP (maximum a posteriori) inference**:

$$\arg \max_y p(y_1, \dots, y_T \mid x_1, \dots, x_T)$$

- Key challenge: efficiently computing these quantities
- Suppose Y_t takes k states. Naively, would seem to require enumerating or summing over k^{T-1} sequences:

$$\begin{aligned} p(y_T \mid x_1, \dots, x_T) &\propto p(y_T, x_1, \dots, x_T) \\ &= \sum_{y_1, \dots, y_{T-1}} p(y_1, \dots, y_T, x_1, \dots, x_T) \end{aligned}$$

Marginal inference in HMMs

- Use **dynamic programming**

$$\begin{aligned} p(y_T, x_1, \dots, x_T) &= \sum_{y_{T-1}} p(y_{T-1}, y_T, x_1, \dots, x_T) \\ &= \sum_{y_{T-1}} p(y_{T-1}, x_1, \dots, x_{T-1}) p(y_T, x_T \mid y_{T-1}, x_1, \dots, x_{T-1}) \\ &= \sum_{y_{T-1}} p(y_{T-1}, x_1, \dots, x_{T-1}) p(y_T, x_T \mid y_{T-1}) \\ &= \sum_{y_{T-1}} p(y_{T-1}, x_1, \dots, x_{T-1}) p(y_T \mid y_{T-1}) p(x_T \mid y_{T-1}, y_T) \\ &= \sum_{y_{T-1}} p(y_{T-1}, x_1, \dots, x_{T-1}) p(y_T \mid y_{T-1}) p(x_T \mid y_T) \end{aligned}$$

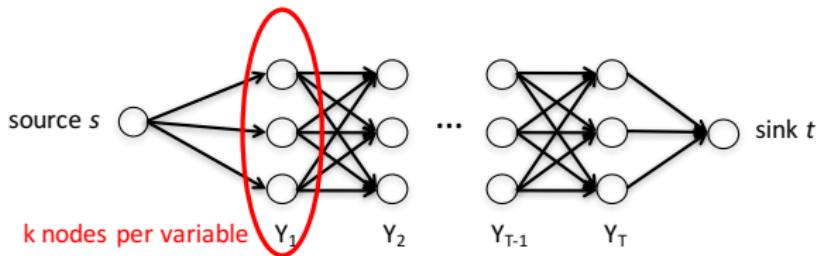
- For $T = 1$, initialize $p(y_1, x_1) = p(y_1)p(x_1 \mid y_1)$
- Total running time is $O(Tk^2)$ – linear time! **Easy to do filtering**

MAP inference in HMMs

- MAP inference in HMMs can *also* be solved in linear time!

$$\begin{aligned}\arg \max_y \Pr(y_1, \dots, y_T \mid x_1, \dots, x_T) &= \arg \max_y p(y_1, \dots, y_T, x_1, \dots, x_T) \\&= \arg \max_y \log p(y_1, \dots, y_T, x_1, \dots, x_T) \\&= \arg \max_y \log [p(y_1)p(x_1 \mid y_1)] + \sum_{t=2}^T \log [p(y_t \mid y_{t-1})p(x_t \mid y_t)]\end{aligned}$$

- Formulate as a shortest paths problem (called the *Viterbi algorithm*)



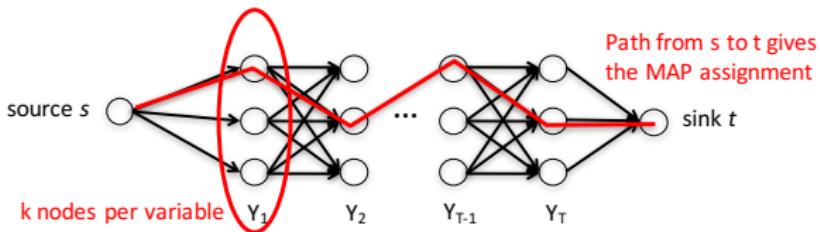
- What should the weight for edges (s, y_1) and (y_T, t) be?
- What should the weight for edges (y_{t-1}, y_t) be?

MAP inference in HMMs

- MAP inference in HMMs can *also* be solved in linear time!

$$\begin{aligned}\arg \max_y \Pr(y_1, \dots, y_T \mid x_1, \dots, x_T) &= \arg \max_y p(y_1, \dots, y_T, x_1, \dots, x_T) \\&= \arg \max_y \log p(y_1, \dots, y_T, x_1, \dots, x_T) \\&= \arg \max_y \log [p(y_1)p(x_1 \mid y_1)] + \sum_{t=2}^T \log [p(y_t \mid y_{t-1})p(x_t \mid y_t)]\end{aligned}$$

- Formulate as a shortest paths problem (called the *Viterbi algorithm*)



- Weight for edge (s, y_1) is $-\log [p(y_1)p(x_1 \mid y_1)]$. Weight for edge $(y_T, t) = 0$
- Weight for edge (y_{t-1}, y_t) is $-\log [p(y_t \mid y_{t-1})p(x_t \mid y_t)]$

Summary

- **Bayesian networks** given by (G, P) where P is specified as a set of local **conditional probability distributions** associated with G 's nodes
- One interpretation of a BN is as a **generative model**, where variables are sampled in topological order
- Factorization implies conditional independencies
- Computing the probability of any assignment is obtained by multiplying CPDs
 - **Bayes' rule** is used to compute conditional probabilities
 - Marginalization or **inference** is often computationally difficult