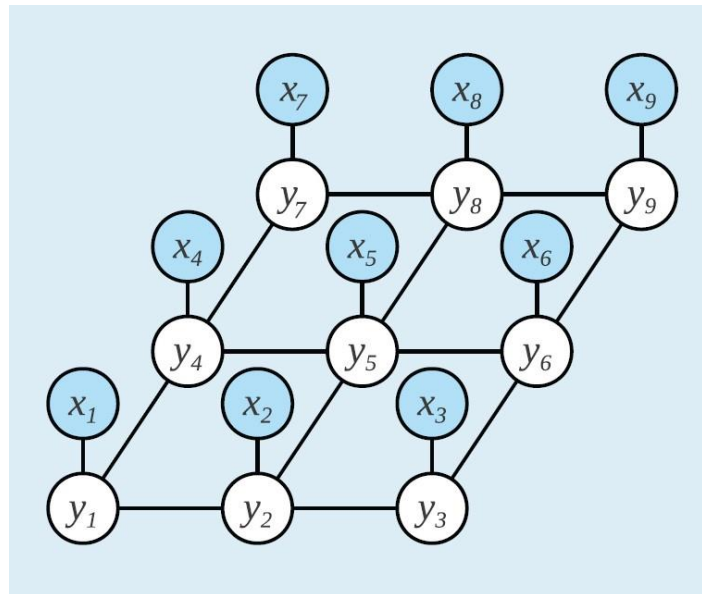


Probabilistic Graphical Models in Bioinformatics

Lecture 1: Introduction



Overview

- **Course website:**
<https://github.com/hesmon/pgm>
- **Teaching assistants:**
 - Naser Elmi and Samaneh Maleknia
- **prerequisites:** Statistics, basic probability theory, programming
- A note to students with background in biology.
- **My Office hours**
 - Sunday 9:00-12:00

- *“Probabilistic graphical models are a marriage between graph theory and probability.”*
- Graphical representations of probability distributions.

Data

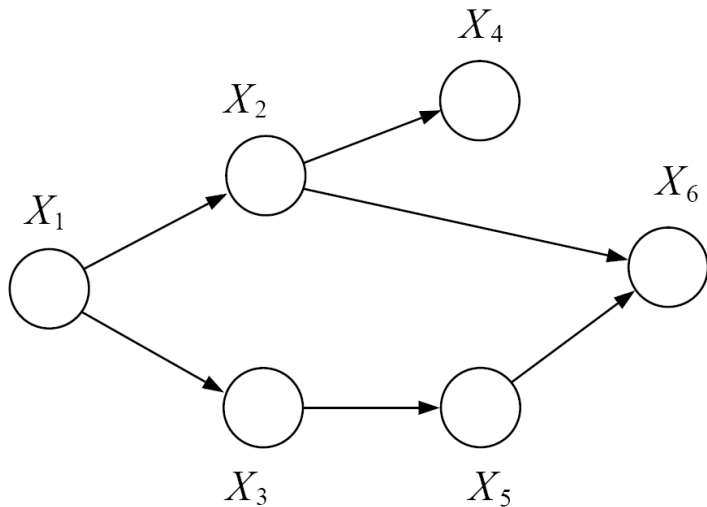


1596	103.3	82.7	80.9	91.3	55.1	123.9	13
4638	119.8	63.1	75.1	160.3	97.5	117.2	78
42.1	116.9	72.9	257.6	243.3	28.9	34.9	164
732.4	788.6	411.3	829.6	680.4	543.4	388.9	488
412.1	473.5	446.1	349.2	430.1	558.2	617.3	351
8706.2	1710.3	1034.2	651.4	954.2	136.4	861.9	5231
5402.7	5913.1	6058.4	4470.5	6159.1	7072.8	5219.3	3293
19502.4	16393.3	14855.6	11297.7	11622.2	19070.2	12580.9	12952
14081.4	13574.0	11698.8	8215	11276.9	16186.2	8365.4	10030
7437.5	6755.5	6002.8	4556.3	6483.2	535.2	528.5	525
147249.2	140345	124028.6	91635.8	121466.1	104034.6	117098.4	122688
179460.8	186706.5	161774.2	125753.2	178065.3	187172.1	152429.9	159575
623	346.2	621.5	600.9	866.4	654.8	488	105
4671	120.6	64.3	67.7	87.5	67.7	1020	438
100.2	100.2	190.7	108.9	315.2	163	478	174
263.3	152.2	73	115.5	104.8	286.3	152.1	110
38.8	67.9	58.4	20.5	101.2	45.2	56	53
17006	1859.7	2019.1	2016.9	2313.1	1640.6	2419.2	1068
763.8	121.1	177.5	1625.6	2043.8	136.8	2491.6	161
10081	970.5	793.2	676.2	1304.8	1348.2	1437.6	1026
2770.3	2605.9	2320.9	1912.8	3441.6	3029.3	2570	2470
223.9	313.1	61.5	178	243.6	323.7	54.9	50
310.5	381.3	153.6	271.8	381.3	341.4	308.3	32
259.6	269.7	33	159.4	242.5	52.3	21.1	138
165.6	264.1	39.8	26.5	166.3	52.5	28.3	135
112.1	28.6	280.7	57.4	339	361.6	70.2	39
46.4	55.4	12.1	31.6	33.5	45.1	161.5	282
106.6	56.6	26.7	44.4	44.4	44.4	68	48

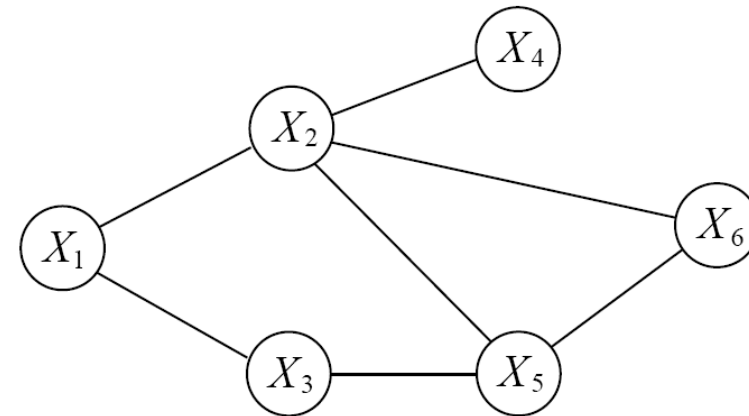
A useful tool for handling uncertainty and complexity

PGM

- Nodes represent random variables/states
- Edges represent probabilistic interaction between variables.
 - The missing arcs represent conditional independence assumptions
 - Makes it simpler to write down the **joint distribution** $P(X_1, X_2, X_3, X_4, X_5, X_6)$



Directed graphical models



undirected graphical
models

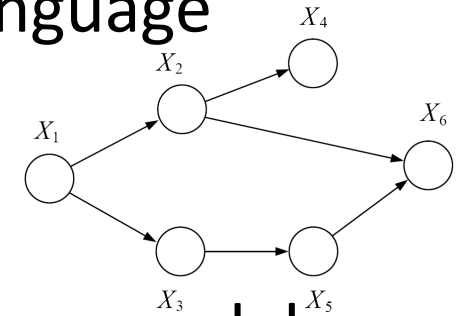
Why joint distribution is important?

- I: intelligence, D: difficulty of the course, G: student's grade
- We can answer different queries using the joint distribution:
 - $P(G = g^1)$
 - $P(D = d_0 \mid I = i_1, G = g_1)$
 - $P(D = d_0 \mid I = i_1)$

I	D	G	Prob.
i^0	d^0	g^1	0.126
i^0	d^0	g^2	0.168
i^0	d^0	g^3	0.126
i^0	d^1	g^1	0.009
i^0	d^1	g^2	0.045
i^0	d^1	g^3	0.126
i^1	d^0	g^1	0.252
i^1	d^0	g^2	0.0224
i^1	d^0	g^3	0.0056
i^1	d^1	g^1	0.06
i^1	d^1	g^2	0.036
i^1	d^1	g^3	0.024

Syllabus overview

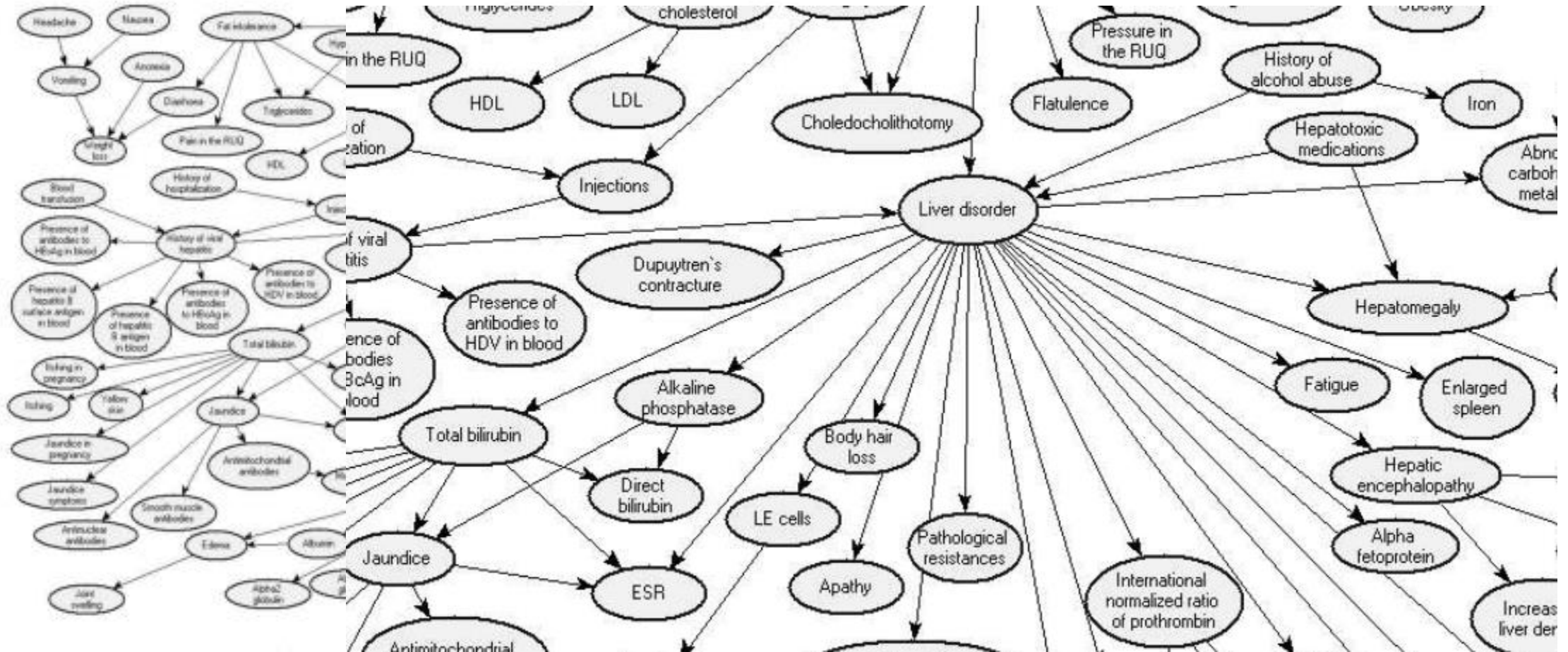
- **Representation:** to represent a distribution in a graphical language
- **Inference:** to answer queries using the joint distribution as our model of world
- **Learning:** learning graphical models from data
- Applications in bioinformatics



Outline

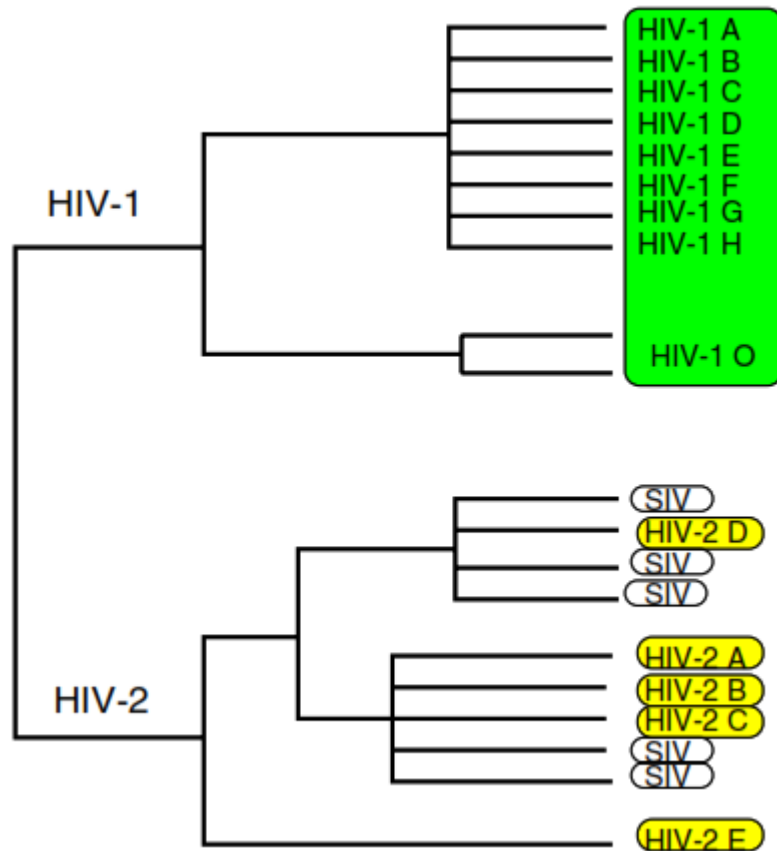
- Motivation: real-world applications of PGMs
- Recap: probability and statistics

Applications-Medical diagnosis system

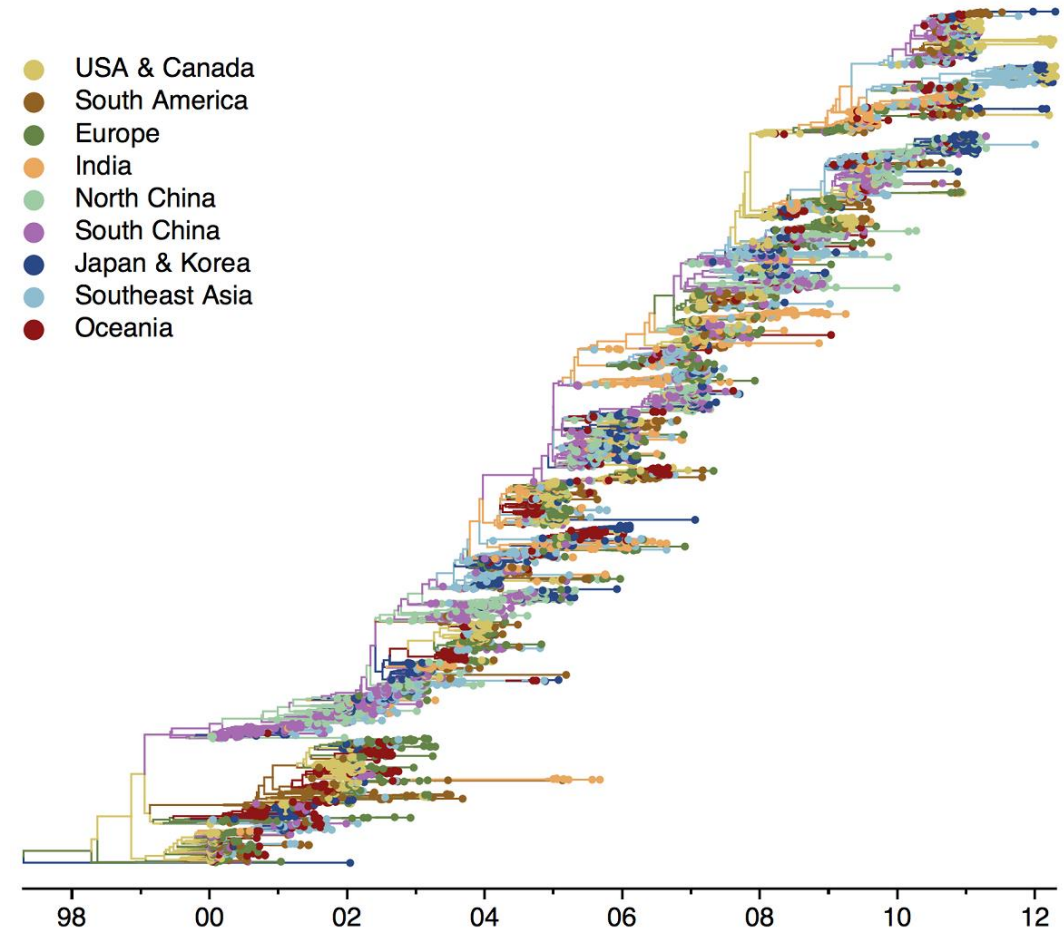


A Bayesian Network Model for Diagnosis of Liver Disorders [Onisko ,1999]

Applications- phylogenetic analysis



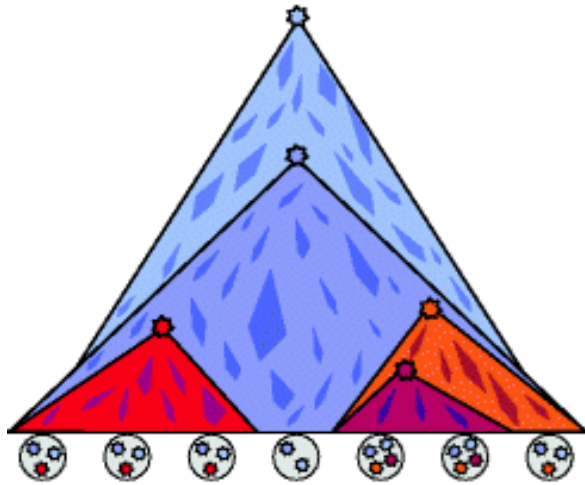
Cross-infection between humans
and monkeys















Global influenza phylogeny

[Bedford et al, 2015]

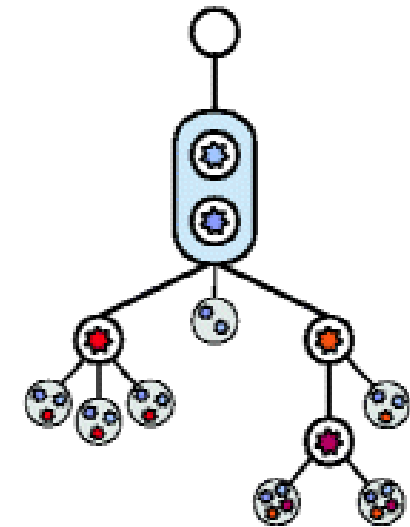
Applications- tumor evolution



tumor evolution with time
progressing downwards

		cells						
		s ₁	s ₂	s ₃	s ₄	s ₅	s ₆	s ₇
								
mutations		1	1	-	1	1	-	1
		1	1	1	1	0	1	1
		1	1	1	0	0	0	0
		1	0	0	0	1	1	1
		0	0	-	0	0	1	0

Mutation data



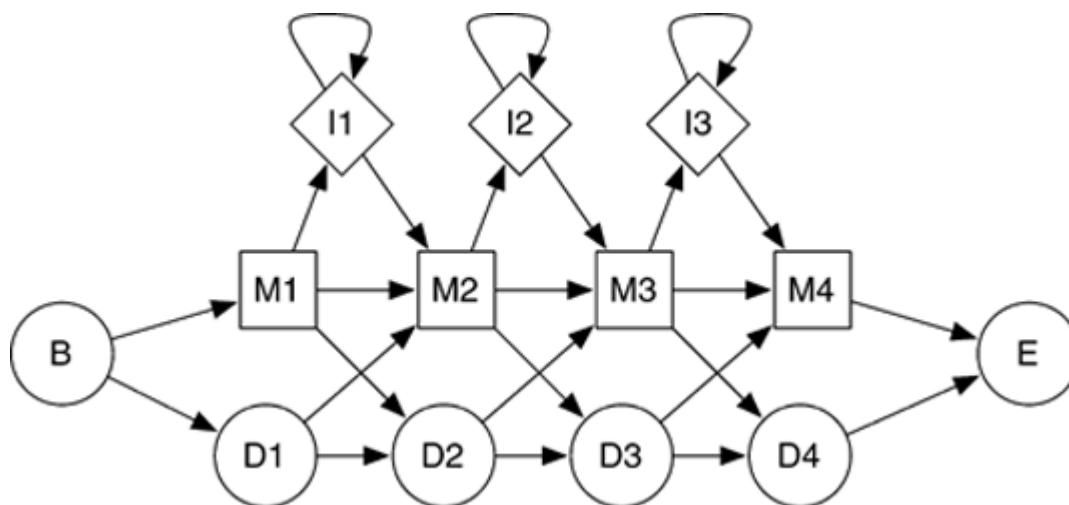
Mutation tree

Applications- multiple sequence alignment

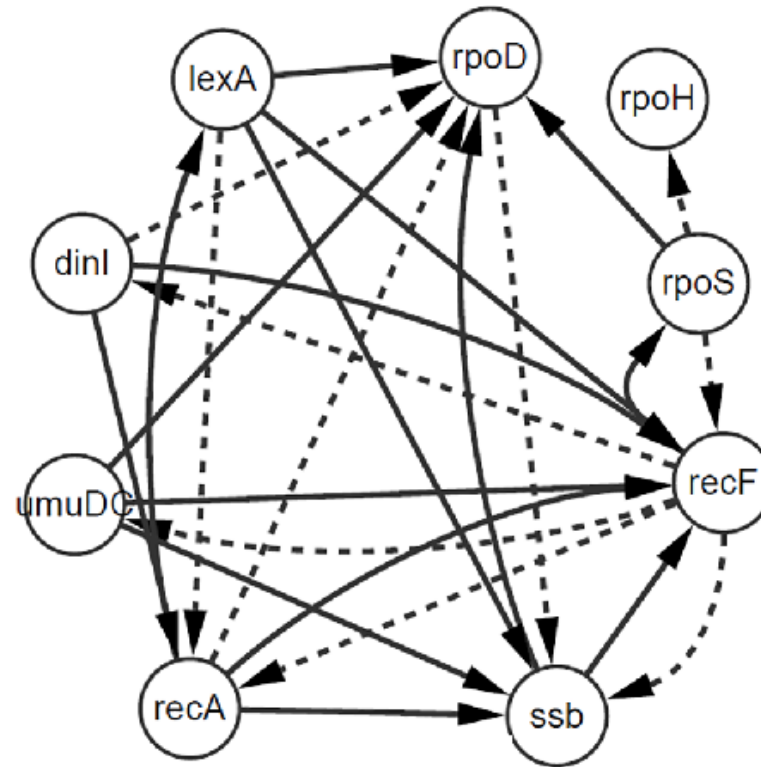
```

Q5E940_BOVIN -----*-----*-----
RLA0_HUMAN -----*-----*-----
RLA0_MOUSE -----*-----*-----
RLA0_RAT -----*-----*-----
RLA0_CHICK -----*-----*-----
RLA0_RANSY -----*-----*-----
Q7ZUG3_BRARE -----*-----*-----
RLA0_ICTPU -----*-----*-----
RLA0_DROME -----*-----*-----
RLA0_DICDI -----*-----*-----
Q54LP0_DICDI -----*-----*-----
RLA0_PLAF8 -----*-----*-----
RLA0_SULAC -----*-----*-----
RLA0_SULTO -----*-----*-----
RLA0_SULSO -----*-----*-----
RLA0_AERPE -----*-----*-----
RLA0_PYRAE -----*-----*-----
RLA0_METAC -----*-----*-----
RLA0_METMA -----*-----*-----
RLA0_ARCFU -----*-----*-----
RLA0_METKA -----*-----*-----

```



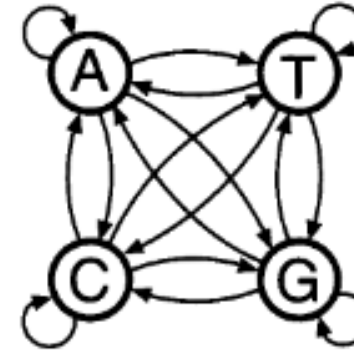
Applications- reconstruction of gene regulatory network



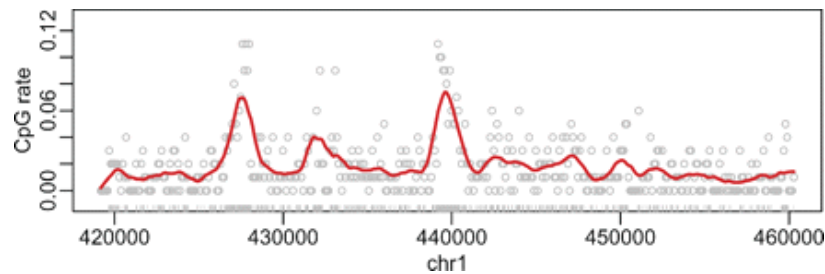
inferred SOS DNA repair network using dynamic
Bayesian network [Liu et al, 2016]

Applications- genome annotation

- Methylation of **CpG islands** plays an important role in regulating gene expression.
- Identifying CpG islands by graphical models



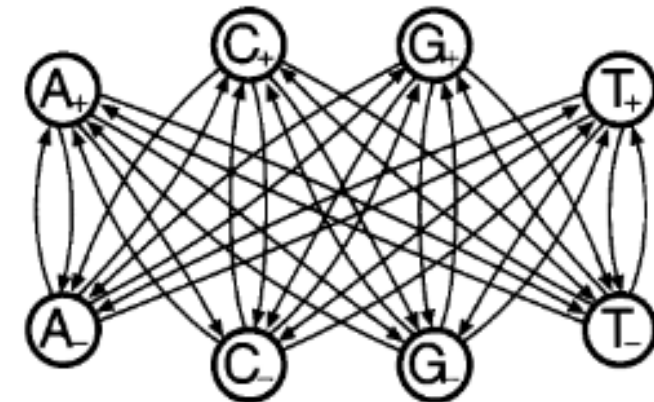
A Markov chain for DNA



[Wo et al, 2010]

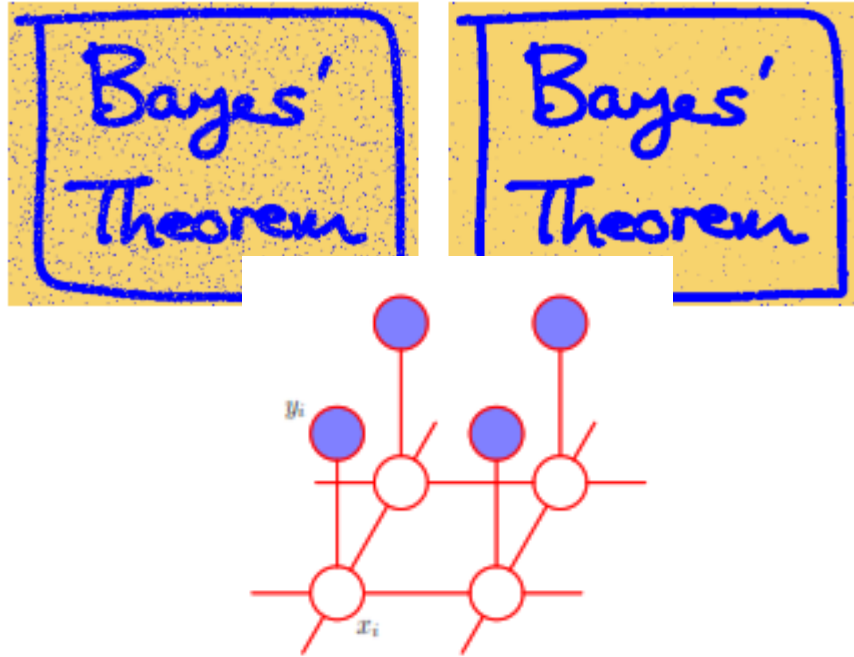
CpG Island

background



[Durbin, 2005]

Other applications



Markov random fields

Image analysis

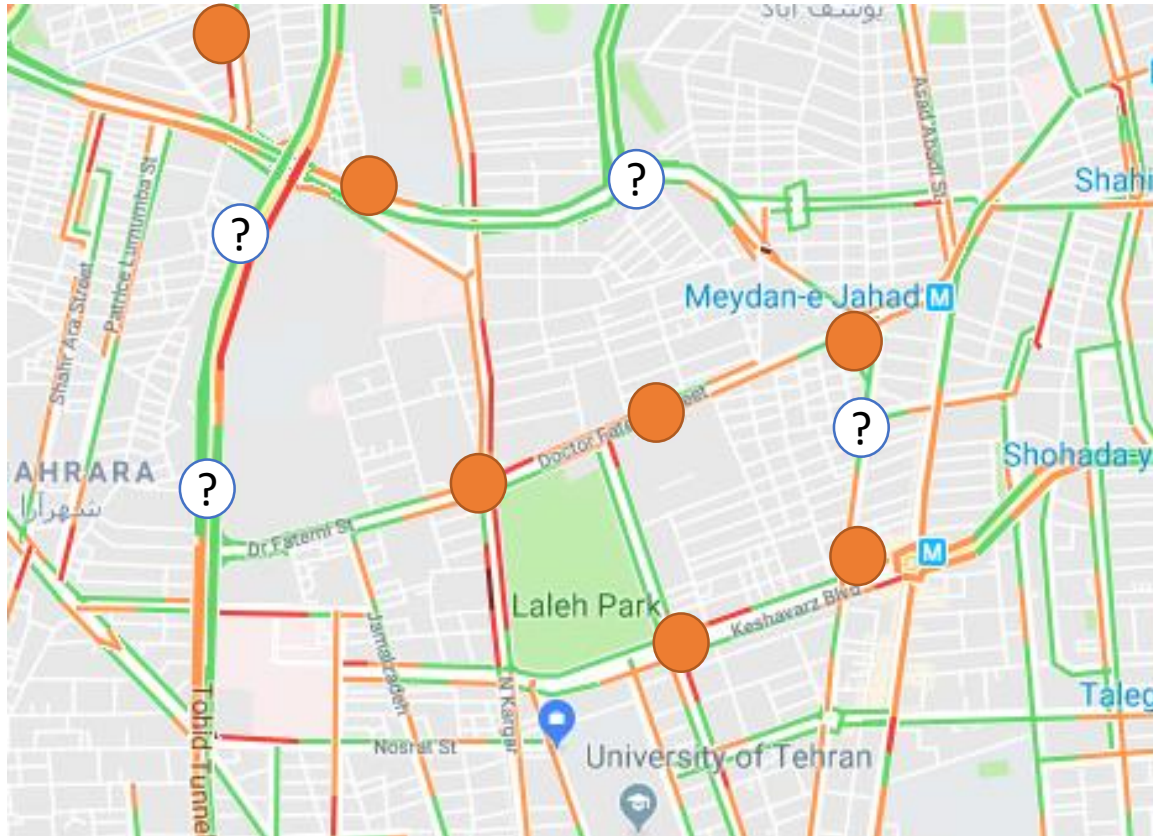


"Do you really like graphical models?"

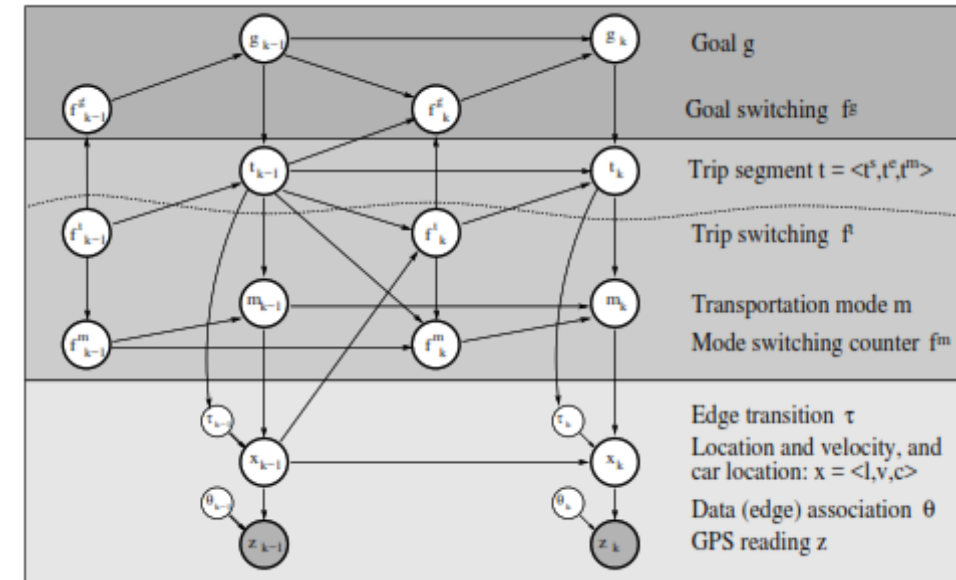
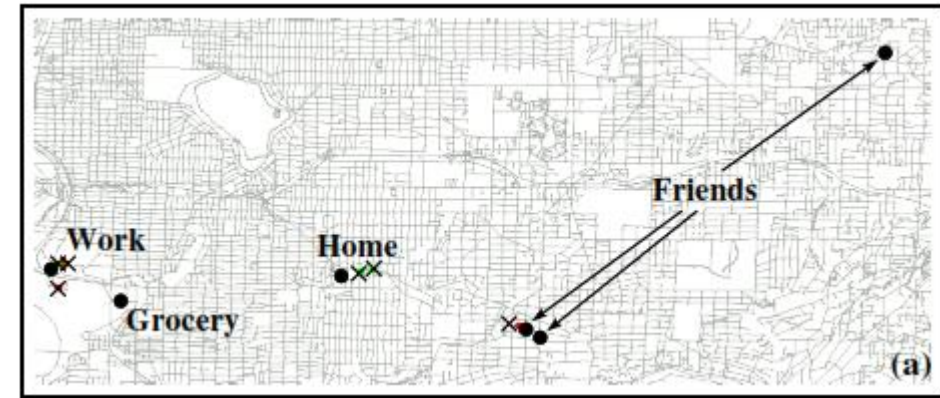
By hidden Markov models

Speech recognition

Other applications-2



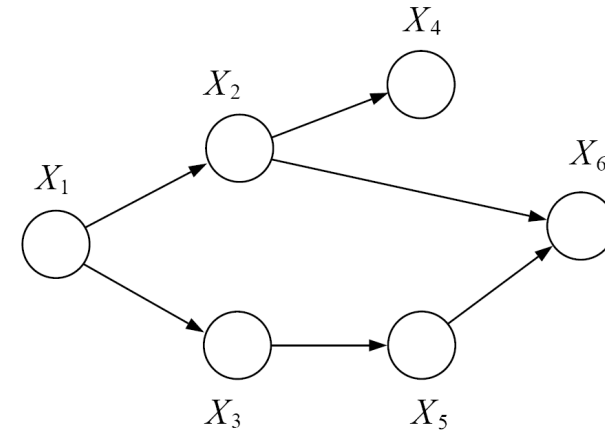
Estimating unmonitored road speeds
from monitored speeds



Hierarchical graphical model representing person's
outdoor movement. Data: raw GPS coordinates

Key ideas

- **Representation:** represent the problem as a collection of random variables X_1, \dots, X_n with joint distribution $P(X_1, \dots, X_n)$ with some conditional independences.



- **Inference:** compute conditional probabilities given some evidences $P(X_i \mid E = e)$.
- **Learning:** estimate the parameters and structure of a Bayesian network from data

Recap: probability and statistics

Random variables

- *A random variable is a variable whose possible values (D) are outcomes (Ω) of a random phenomenon (from Wikipedia).*
- A random variable X is a mapping $X: \Omega \rightarrow D$.
- Example: a random variable *intelligence* that takes as values either *high* or *low*.
- Random variables can be either
 - Categorical or discrete
 - Or continuous

Examples

- Binary-valued random variable: “biased coin flips”

$$D = \{H, T\}$$

$$P(X = H) = p. \text{ Hence, } P(X=T) = 1-p$$

The distribution of such a random variable is called a Bernoulli distribution and denoted by $X \sim \text{Bernoulli}(p)$

- Categorical random variable: “ k -sided dice”

$$D = \{1, \dots, k\}$$

$$P(X = i) = p_i \text{ subject to } \sum_{i=1}^k p_i = 1$$

The distribution of such a random variable is called a multinomial distribution and denoted by $X \sim \text{Mult}(p)$

Joint distribution

- An example of joint distribution $P(\text{intelligence, grade})$

		Intelligence		
		low	high	
Grade	A	0.07	0.18	0.25
	B	0.28	0.09	0.37
	C	0.35	0.03	0.38
		0.7	0.3	1

- Marginal distribution:

$$P(G = A) = P(G = A, I = \text{low}) + P(G = A, I = \text{high}) = 0.07 + 0.18 = 0.25$$

- **Question:** suppose X_1, \dots, X_n are binary-valued variables. How many parameters do we need to specify the joint distribution?
 - $2^n - 1$ parameters for 2^n different assignments of values x_1, \dots, x_n .

Marginal distributions

Discrete case:
$$P(X) = \sum_Y P(X, Y), \quad P(Y) = \sum_X P(X, Y)$$

Continuous case:
$$P(X) = \int_Y P(X, Y) dY, \quad P(Y) = \int_X P(X, Y) dX$$

- Given the joint distribution $P(X_1, \dots, X_n)$
- Then

$$P(X_i = x_i) = \sum_{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n} P(X_1 = x_1, \dots, X_n = x_n)$$

- **Question:** if all X_i binary: How many terms?

Conditional probabilities

- The conditional probability of X and Y is

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

- G : overexpression of an oncogene
- C : the presence of a tumor
- $P(G, C)$: Prob. of oncogene overexpression and the presence of a tumor.
- $P(G / C)$: Prob. of oncogene overexpression in cancer patients (can be assessed by counting).
- $P(C / G)$: Prob. of cancer given gene expression measurement (**might be difficult to assess**).

Bayes' rule (important)

- Since

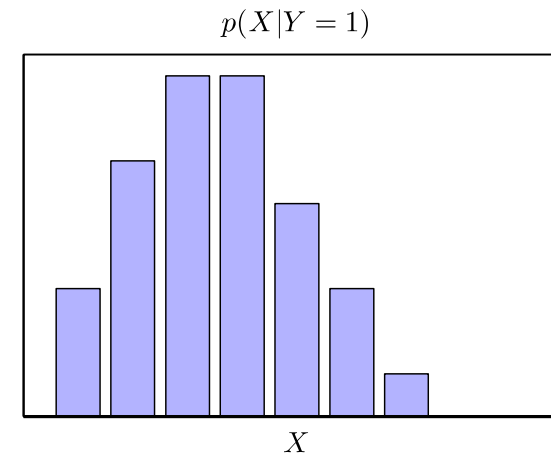
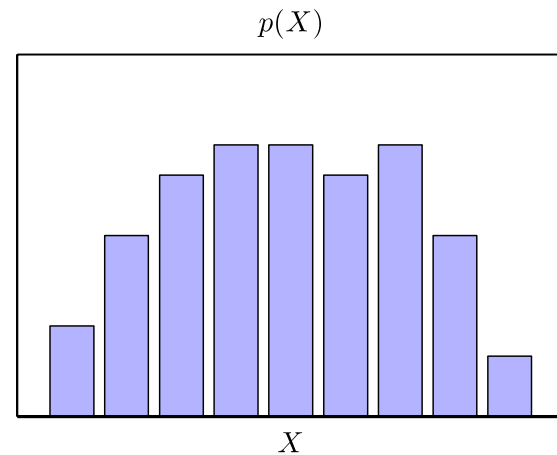
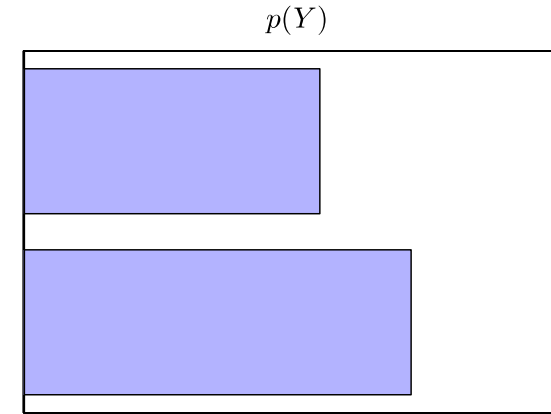
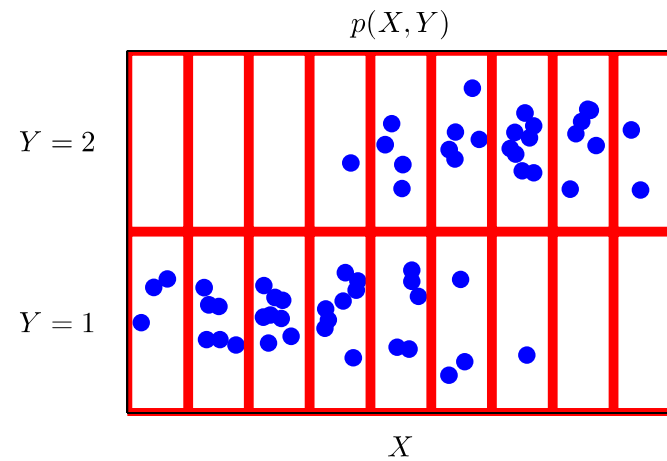
$$P(G, C) = P(G|C)P(C) = P(C|G)P(G)$$

- Hence

$$P(C|G) = \frac{P(G|C)P(C)}{P(G)}$$

- Hence, the diagnostic conditional probability $P(C|G)$ can be computed without having to be determined **explicitly**.

Example: joint, marginal and conditional probabilities



Chain rule (important)

$$P(X_1, \dots, X_n) = P(X_1)P(X_2 \mid X_1) \dots, P(X_n \mid X_1, \dots, X_{n-1})$$

Independent random variables

- For n independent random variables X_1, \dots, X_n we have

$$P(X_1, \dots, X_n) = P(X_1) \dots P(X_n)$$

- **Question:** how many parameters in this case?
- Something in between full dependence and full independence assumptions?

Conditional independence

- It is not often that we encounter two independent random variables. A more common situation is when two random variables are independent given an additional random variable.
- Proposition: $X \perp Y \mid Z \iff P(X, Y \mid Z) = P(X \mid Z)P(Y \mid Z)$
- An alternative way to prove conditional independence of X and Y given Z

$$P(X \mid Y, Z) = P(X \mid Z)$$



Why conditional independence is useful?

- By chain rule

$$P(X_1, \dots, X_n) = P(X_1)P(X_2 | X_1) \dots P(X_n | X_1, \dots, X_{n-1})$$

Question: how many parameters?

- Suppose $X_i \perp X_2, \dots, X_{i-1}, X_{i+1}, \dots, X_n \mid X_1$ for all $i \in \{2, \dots, n\}$. Then

$$P(X_1, \dots, X_n) = ?$$

Why conditional independence is useful?

- By chain rule

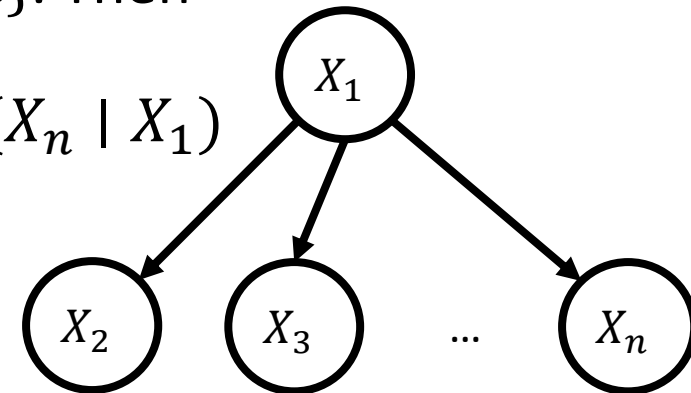
$$P(X_1, \dots, X_n) = P(X_1)P(X_2 | X_1) \dots, P(X_n | X_1, \dots, X_{n-1})$$

Question: how many parameters?

- Suppose $X_i \perp X_2, \dots, X_{i-1}, X_{i+1}, \dots, X_n | X_1$ for all $i \in \{2, \dots, n\}$. Then

$$P(X_1, \dots, X_n) = P(X_1)P(X_2 | X_1)P(X_3 | X_1) \dots, P(X_n | X_1)$$

Question: how many parameters?



Querying a distribution

- Probability queries
 - **Evidence:** a subset E of random variables in the model, and an instantiation e to these variables.
 - **Query variables:** a subset Y of random variables.
 - Query: $P(Y \mid E = e)$
 - **Question:** $P(G \mid i^0)$?
- Map queries:
 - To find the most likely assignment to the variables W given the evidence $E=e$

$$MAP(W|e) = \arg \max_w P(w \mid e)$$
 - **Question:** $MAP(D, G \mid i^1)$?

I	D	G	Prob.
i^0	d^0	g^1	0.126
i^0	d^0	g^2	0.168
i^0	d^0	g^3	0.126
i^0	d^1	g^1	0.009
i^0	d^1	g^2	0.045
i^0	d^1	g^3	0.126
i^1	d^0	g^1	0.252
i^1	d^0	g^2	0.0224
i^1	d^0	g^3	0.0056
i^1	d^1	g^1	0.06
i^1	d^1	g^2	0.036
i^1	d^1	g^3	0.024

Continuous random variables

- Probability density function f for continuous random variable X

$$\int_x f(x)dx = 1$$

- Cumulative distribution P

$$P(X \leq a) = \int_{-\infty}^a f(x)dx$$

- Similarly for multiple random variables

$$P(a_1 \leq X_1 \leq b_1, \dots, a_n \leq X_n \leq b_n) = \int_{a_1}^{b_1} \dots \int_{a_n}^{b_n} f(x_1, \dots, x_n)dx_1 \dots dx_n$$

Project 1

- Application of Bayesian networks on gene expression data
- Deadline: Esfand 19, 1397 (March 9, 2019)

Coursework

- **Grading based on**
 - Data analysis projects: 25% (two students per group)
 - Problem sets: 25%
 - Reading assignments: 10%
 - Final exam: 40%
- **Reading assignments:**
 - Deadline: the day before the lecture by 17:00.
- **Late submission policy**
 - In total, 5 late days for problem sets. 3 late days for projects. 6 late days for reading assignments.
 - No further extension. Zero point for submitting assignments after the deadline.
 - Regardless, in order to pass the course **you have to submit all assignments.**
- Discussing assignments is encouraged, but you must turn in your own solutions.

Bibliography

1. Daphne Koller and Nir Friedman. Probabilistic graphical models: principles and techniques, 2011, MIT press.
2. Dirk Husmeier. Probabilistic Modeling in Bioinformatics and Medical Informatics, 2010, Springer.
3. Christopher M. Bishop. Pattern Recognition and Machine Learning, 2016, Springer.

