# COMP7180: Quantitative Methods for DAAI

(Credits from Prof. Andrew Ng)

(Credits from HKBU)

## Course Instructors: Dr. Yang Liu and Dr. Bo Han

## Teaching Assistant: Mr. Minhao Li

# Course Contents

- Continuous and Discrete Random Variables (Week 7)

- Conditional Probability and Independence (Week 8)

- Maximum Likelihood Estimation (Week 9) ← Our Focus

- Mathematical Optimization (Week 10)

- Convex and Non-Convex Optimization (Week 11)

- Course Review (Week 12)

# Review: Exercise

- Given three random variables X, Y and Z whose ranges are {0,1},
 if

$$P(X=1|Y=1,Z=1)>P(X=1|Y=0,Z=1)$$

$$P(X=1|Y=1,Z=0)>P(X=1|Y=0,Z=0),$$

is it correct that

$$P(X=1|Y=1)>P(X=1|Y=0)?$$

To answer this issue, we introduce an example

# Review: Exercise



Consider following table of counts that are obtained from an observed sample of individuals both males and females, who had taken Covid-19. X, Y, Z are random variables. Z represents the gender, X represents whether individuals recover from Covid-19, and Y represents whether individuals have been treated.

Z = 1  Male;          Z = 0     Female;
X = 1  Recovery;   X = 0     No Recovery;
Y = 1  Treatment;  Y = 0      No Treatment;

|       | Z=1, X=1 | Z=1, X=0 | Z=0, X=1 | Z=0, X=0 |
|-------|----------|----------|----------|----------|
| Y=1   | 245      | 105      | 315      | 735      |
| Y=0   | 630      | 420      | 70       | 280      |

# Review: Exercise

Please check that

P(X=1|Y=1,Z=1) = 0.7          P(X=1|Y=0,Z=1) = 0.6
P(X=1|Y=1,Z=0)= 0.3           P(X=1|Y=0,Z=0) = 0.2
P(X=1|Y=1) = 0.4               P(X=1|Y=0) = 0.5

P(X=1|Y=1,Z=1)>P(X=1|Y=0,Z=1),
P(X=1|Y=1,Z=0)>P(X=1|Y=0,Z=0).

But P(X=1|Y=1)<P(X=1|Y=0).

# Review: Exercise

It seems that

P(X=1|Y=1,Z=1)>P(X=1|Y=0,Z=1),
P(X=1|Y=1,Z=0)>P(X=1|Y=0,Z=0) and

P(X=1|Y=1)<P(X=1|Y=0) are conflict.

Why does this happen?

It is called Simpson's paradox.

# Review: Exercise

Note that

P (Z=1|Y=1)P(X=1|Y=1,Z=1)+P(Z=0|Y=1)P(X=1|Y=1,Z=0)=P (X=1|Y=1)

P (Z=1|Y=0)P(X=1|Y=0,Z=1)+P(Z=0|Y=0)P(X=1|Y= 0,Z=0)=P(X=1|Y=0)

So let u = P(Z=1|Y=1) and v = P(Z = 1|Y = 0), then we cannot ensure that

$$0.7u + 0.3(1 − u) > 0.6v + 0.2(1 − v).$$

That depends on the values of u and v. If v−u>0.25, we can get that

0.7u+0.3(1−u)<0.6v+0.2(1−v).

But, u=P(Z=1|Y=1)=0.25 and v =P(Z=1|Y=0)=0.75. So v−u=0.5>0.25. That is the reason why this happens.

# Maximum Likelihood (ML) Estimation

- In many artificial intelligence and machine learning applications, the objective is **to estimate the model parameters from the given data**.

- For example, given a distribution class $P(X;\alpha)$, where $\alpha$ is a parameter from a parameter space. Now given data $(x1,x2,...,xn)$ which are drawn from an unknown distribution $P(X;\alpha_0)$, we want to ask that how to select a suitable parameter $\alpha_0$ by given data $(x1,x2,...,xn)$?

- The **Maximum Likelihood Estimation (MLE)** is one of the most widely used methods of estimating the parameters of a model.

# Maximum Likelihood (ML) Estimation

- The method of Maximum Likelihood Estimation (MLE) selects the set of values of the model parameters that maximizes the likelihood function.

- In other words, the basic principle of MLE is to choose values that "explain" the data best by maximizing the probability of the data we've seen as a function of the parameters.

- The **Maximum Likelihood Estimation (MLE)** is one of the most widely used methods of estimating the parameters of a model.

- It answers the question: What values of parameters would make the observations <span style="color:red">most probable</span> ?

# Maximum Likelihood (ML) Estimation

- A distribution class P(X;α), where α is from a parameter space Δ.
- For each αα from the space Δ, P(X;α) corresponds to a distribution.
- We have data  S=(x1,x2,...,xn) , which are drawn from an <span style="color:red">unknown</span> distribution P(X).

<span style="color:red">Problem:</span> what is the <span style="color:red">optimal parameter $α^*$</span> selected from the parameter space <span style="color:red">Δ</span>, such that the seleced distribution $P(X; α^*)$ is the <span style="color:red">most possible distribution sampling data S</span>?

# Maximum Likelihood (ML) Estimation

The seleced distribution P(X;α*) is the most possible distribution sampling data S=(x1,x2,…,xn), i.i.d.

Understanding above sentence, we can formulate it as follows:

$$\underset{\alpha \in \varDelta}{\mathrm{argmax}} \prod_{i=1}^{n} P(X = \mathrm{xi}; \alpha)$$

here we assume $P(X; \alpha)$ is discrete distributions.

- $\prod_{i=1}^{n} P(X = \mathrm{xi}; \alpha)$ means the largest probability for $P(X; \alpha)$ that S is observed.

# Maximum Likelihood (ML) Estimation

How to address the equation?

Step 1.

$$\underset{\alpha \in \Delta}{\operatorname{argmax}} \prod_{i=1}^{n} P(X = xi; \alpha)$$

$$\underset{\alpha \in \Delta}{\operatorname{argmax}} \prod_{i=1}^{n} P(X = xi; \alpha) = \underset{\alpha \in \Delta}{\operatorname{argmax}} \; \log \prod_{i=1}^{n} P(X = xi; \alpha)$$

# Maximum Likelihood (ML) Estimation

$$\underset{\alpha \in \Delta}{\mathrm{argmax}} \prod_{i=1}^{n} P(X = x_i; \alpha) = \underset{\alpha \in \Delta}{\mathrm{argmax}} \ \log \prod_{i=1}^{n} P(X = x_i; \alpha)$$

Step 2. Using the property of log function:

$$\log \prod_{i=1}^{n} P(X = x_i; \alpha) = \sum_{i=1}^{n} \log P(X = x_i; \alpha)$$

13

# Maximum Likelihood (ML) Estimation

Therefore,

$$\underset{\alpha \in \Delta}{\mathrm{argmax}} \prod_{i=1}^{n} P(X = xi; \alpha) = \underset{\alpha \in \Delta}{\mathrm{argmax}} \sum_{i=1}^{n} \log P(X = xi; \alpha)$$

Step 3. We need to optimize

$$\underset{\alpha \in \Delta}{\mathrm{argmax}} \sum_{i=1}^{n} \log P(X = xi; \alpha) \quad (1)$$

and obtain the optimal solution.

The solution of Eq. 1 is called Maximum Likelihood Estimation.

# Maximum Likelihood (ML) Estimation

If the distribution class consists of continuous distributions, that is P(X;α) is a <span style="color:red">continuous distribution</span> with respect to all α∈ $\Delta$.

Then the <span style="color:red">Maximum Likelihood Estimation</span> is

$$\underset{\alpha\in\Delta}{\mathrm{argmax}} \ \sum_{i=1}^{n} \log p_X(\mathrm{x}i; \alpha) \quad (2)$$

where $p_X(\mathrm{x}; \alpha)$ is the density function of P(X;α).

# Maximum Likelihood (ML) Estimation

How to obtain the solution of

$$\underset{\alpha \in \Delta}{\mathrm{argmax}} \; \sum_{i=1}^{n} \log P(X = xi; \alpha) \; ?$$

- This is related to optimization problem.

- Generally, there are no unviersal approaches to give soultions to all Maximum Likelihood (ML) Estimation.

- The approaches are case by case.

# Maximum Likelihood (ML) Estimation

In this class, we introduce a common used approach.

This approach is based on a <span style="color:red">simple theorem</span>:

If 1) a function f(x1,x2,…,xd) is <span style="color:red">differentiable,</span>
2) $x^* = (x1^*, x2^*, …., xd^*)$ is the <span style="color:red">maximum point</span> of f, then

$$\frac{\partial f}{\partial xi}(x1^*, x2^*, …., xd^*) = 0.$$

# Maximum Likelihood (ML) Estimation

Using this theorem, if $\sum_{i=1}^{n} \log P(X = x_i; \alpha)$ is differentiable, then

Let $\alpha = (\alpha_1, \alpha_2, \ldots, \alpha_d)$,

$$\frac{\partial \sum_{i=1}^{n} \log P(X = x_i; \alpha)}{\partial \alpha_i} = 0, \text{ for } i = 1, \ldots, d$$

Then, addressing above equations.

Check that you've found a <span style="color:red">maximum</span> rather than a <span style="color:red">minimum</span> or <span style="color:red">saddle-point</span>, and be careful if $\alpha$ belongs to $\Delta$.

# Exercises: MLE for Binomial Distribution

Suppose that $x_1, x_2, . . ., x_n$ (i.i.d)
represent the outcomes
of n independent Bernoulli trials
(for example, coin flipping), each
with success probability $\mu$.

- $P(X=1; \mu) = \mu$
- $P(X=0; \mu) = 1-\mu$

So $P(X=x_i; \mu) = \mu^{x_i}(1-\mu)^{1-x_i}$

**MLE: For which $\mu$ is $x_1, x_2, \ldots x_n$ most likely?**

# Exercises: MLE for Binomial Distribution

Maximum Likelihood (ML) Estimation：

$$\underset{0 \le \mu \le 1}{\text{argmax}} \ \sum_{i=1}^{n} \log P(X = x_i; \mu)$$

$$\sum_{i=1}^{n} \log P(X = x_i; \mu) \ = \ \sum_{i=1}^{n} \log \mu^{x_i}(1 - \mu)^{1 - x_i}$$

$$= \sum_{i=1}^{n} x_i \log \mu \ + \sum_{i=1}^{n}(1 - x_i)\log(1 - \mu)$$

# Exercises: MLE for Binomial Distribution

Derivation of MLE

$$\frac{\partial \sum_{i=1}^{n} \log P(X=x_i;\mu)}{\partial \mu},$$

- $\frac{\partial \sum_{i=1}^{n} \log P(X=x_i;\mu)}{\partial \mu} = \sum_{i=1}^{n} \frac{x_i}{\mu} - \sum_{i=1}^{n} \frac{(1-x_i)}{1-\mu} = 0.$

- We have that $\mu = \sum_{i=1}^{n} x_i / n$

# Exercises: MLE for Gaussian Distribution

Suppose you have x1,x2,…,xn (i.i.d)  $N(\mu, \sigma^2)$

$$\sqrt{\frac{1}{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x-\mu)^2)$$



- But you don't know $\mu$ (you do know $\sigma^2$)

**MLE: For which μ is x1, x2, …, xn most likely?**

# Exercises: MLE for Gaussian Distribution

Compute the MLE $\quad \underset{\mu \in R}{\operatorname{argmax}} \; \sum_{i=1}^{n} \log p_X(\mathsf{x}i; \mu)$

$$\underset{\mu \in R}{\arg\max} \frac{1}{\sqrt{2\pi}\,\sigma} \sum_{i=1}^{n} -\frac{(\mathsf{x}i - \mu)^2}{2\sigma^2}$$

$$= \underset{\mu \in R}{\arg\min} \sum_{i=1}^{n} (\mathsf{x}i - \mu)^2$$

# Exercises: MLE for Gaussian Distribution

Derivation the equation $\arg\min_{\mu\in R} \sum_{i=1}^{n} (x_i - \mu)^2$

$$\frac{d \sum_{i=1}^{n}(x_i - \mu)^2}{d\mu} = 2\sum_{i=1}^{n}(x_i - \mu) = 0$$



So the solution is

$$\mu = \frac{\sum_{i=1}^{n} x_i}{n}$$

# Exercises: MLE for Gaussian Distribution

•In conclusion, the best estimate of the mean of a gaussian distribution is the mean of the sample!

$$\mu = \frac{\sum_{i=1}^{n} x_i}{n}$$

# Exercises: MLE for high-dimensional Gaussian Distribution

- Given a 2 $\times$ 2 positive semi-definite matrix $\Sigma$ and a 2 $\times$ 1 vector $\mu$, a three dimensional normal distribution $N(\mu, \Sigma)$ can be represented as follows: the density function of this distribution is



Multivariate Normal Distribution

$$p_{XY}(x, y; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^2 |\Sigma|}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^{\mathrm{T}} \Sigma^{-1} (\mathbf{x}-\mu)}$$

where $|\Sigma|$ is the determinant of $\Sigma$ and $\mathbf{x}=(x, y)^{\mathrm{T}}$.

# Exercises: MLE for high-dimensional Gaussian Distribution

- If μ=(a1,a2) and Σ is a diagonal matrix with eigenvalues λ1, λ2 (λ1 > 0, λ2 > 0),

$$\Sigma = \begin{bmatrix} \lambda 1 & 0 \\ 0 & \lambda 2 \end{bmatrix}$$

then $f(x, y)$ can be writeen as:

Multivariate Normal Distribution

$$p_{XY}(x, y; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^2 \lambda 1 * \lambda 2}} e^{-\frac{1}{2\lambda 1}(x-a1)^2 - \frac{1}{2\lambda 2}(y-a2)^2}$$

# Exercises: MLE for high-dimensional Gaussian Distribution

- If

$$\Sigma = \begin{bmatrix} \lambda1 & 0 \\ 0 & \lambda2 \end{bmatrix}$$

and we have n data $(x1,y1),...,(xn,yn)$ sampled from a two-dimensional Gaussian Distribution $N(\mu, \Sigma)$, i.i.d., calculate $\mu$ by the maximum likelihood estimation method.



Multivariate Normal Distribution

# Exercises: MLE for high-dimensional Gaussian Distribution

- Maximum Likelihood (ML) Estimation：

$$\underset{\mu}{\operatorname{argmax}} \ \sum_{i=1}^{n} \log p_{XY}(\mathrm{x}i, \mathrm{y}i; \mu, \Sigma)$$

It is equal to

$$\underset{a1,a2}{\operatorname{argmin}} \sum_{i=1}^{n} \left( \frac{1}{2\lambda 1}(\mathrm{x}i - a1)^2 + \frac{1}{2\lambda 2}(\mathrm{y}i - a2)^2 \right)$$

# Exercises: MLE for high-dimensional Gaussian Distribution

Derivation the equation $\quad G(a1, a2) = \sum\limits_{i=1}^{n}(\dfrac{1}{2\lambda 1}(xi - a1)^2 + \dfrac{1}{2\lambda 2}(yi - a2)^2)$

$$\dfrac{\partial G}{\partial a1} = \sum\limits_{i=1}^{n}\dfrac{a1 - xi}{\lambda 1} = 0, \qquad\qquad \dfrac{\partial G}{\partial a2} = \sum\limits_{i=1}^{n}\dfrac{a2 - yi}{\lambda 2} = 0$$

So $a1 = \dfrac{1}{n}\sum_{i=1}^{n} xi$ , $\quad a2 = \dfrac{1}{n}\sum_{i=1}^{n} yi$

# Exercises: MLE for high-dimensional Gaussian Distribution

- In conclusion, the best estimate of the mean of a two-dimensional gaussian distribution is the mean of the sample!

In fact, it also holds for high-dimensional gaussian distribution.



Multivariate Normal Distribution

# Maximum A Posteriori (MAP) Estimation

- In Bayesian statistics, a maximum a posteriori probability (MAP) estimate is an estimate of an unknown quantity, that equals the mode of the posterior distribution.

- The MAP can be used to obtain a point estimate of an unobserved quantity on the basis of empirical data. It is closely related to the method of maximum likelihood (ML) estimation, but employs an augmented optimization objective which incorporates a prior distribution (that quantifies the additional information available through prior knowledge of a related event) over the quantity one wants to estimate.

# Maximum A Posteriori (MAP) Estimation

- Consider a distribution class P(X;α).

- MAP regards the <span style="color:red">parameter</span> α as a <span style="color:red">random variable</span>.

- Therefore, we can rewrite P(X;α) as the conditional distribution <span style="color:red">P(X|α)</span>.

- In MAP, we aim to estimate the parameter <span style="color:red">α</span>, given the data x1,...,xn i.i.d. from an unknown distribution $P(X;α^*)$:

$$\underset{α}{\mathrm{argmax}} \quad P(α|x1,x2,...,xn)$$

# Maximum A Posteriori (MAP) Estimation

- In MAP, we mainly use <span style="color:red">Bayesian rule</span>:

  $$P(a|x1,x2,...,xn) = (P(x1,x2,...,xn|a) \, P(a))/P(x1,x2,...,xn)$$

- Note that when the data x1,x2,...,xn are given, $P(x1,x2,...,xn)$ is a <span style="color:red">constant</span>.

- Therefore,

  $$\underset{a}{\text{argmax}} \; P(a|x1,x2,...,xn) = \underset{a}{\text{argmax}} \; P(x1,x2,...,xn|a) \, P(a)$$

# Maximum A Posteriori (MAP) Estimation

- $\operatorname*{argmax}_{\alpha}$ P(x1,x2,...,xn|α) P(α)

is called maximum a poserior estimation.

Because x1,...,xn are Independent and identically distributed, then
$$\operatorname*{argmax}_{\alpha} \text{P(x1,x2,...,xn|α) P(α)} = \operatorname*{argmax}_{\alpha} \prod_{i=1}^{n} P(X = x\text{i}|\text{α})\text{P(α)}$$

- Compared with ML estimation, MAP estimation has a prior distribution P(α).

- MAP estimation: find the most likely parameter settings under the posterior.

# Exercises: Maximum A Posteriori (MAP) Estimation

- Suppose we need to determine if a patient has a rare disease, given a laboratory test of that patient.

- We consider a set of two random variables: α=1 (disease) and α=-1 (no disease).

- X is the random variable related to the laboratory test. X=1 means the positive in the laboratory test, and X=0 means the negative in the laboratory test.

Credicts from Book Foundations of Machine Learning

# Exercises: Maximum A Posteriori (MAP) Estimation

- Suppose that the disease is rare, say P(α=1)=0.005

- The laboratory is relatively accurate: P(X=1|α=1) = 0.98, P(X=0|α=0) = 0.95.

- If the test is positive, what should be the diagnosis?

  In other words, we have a data x which is equal to 1.

# Exercises: Maximum A Posteriori (MAP) Estimation

- Using maximum a poserior estimation

$$\underset{\alpha}{\mathrm{argmax}} \ P(X=1|\alpha) \ P(\alpha)$$

- If α = 1, then P(X=1|α=1) P(α=1) = 0.98*0.005 = 0.0049.

- If α = 0, then P(X=1|α=0) P(α=0) = 0.05*0.995 = 0.4975.

Because P(X=1|α=0) P(α=0)>P(X=1|α=1) P(α=1), we obtain that the solution is

$$0 = \underset{\alpha}{\mathrm{argmax}} \ P(X=1|\alpha) \ P(\alpha)$$

# Exercises: Maximum A Posteriori (MAP) Estimation

Thus, in this case, the MAP prediction is no disease: according to the MAP solution,with the values indicated, a patient with a positive test result is nonetheless more likely not to have the disease!

# Exercises: MAP for Gaussian Distribution

Suppose you have x1,x2,...,xn (i.i.d) $N(\mu, \sigma^2)$ with density

$$p(x|\mu) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x-\mu)^2)$$

$$p(\mu) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(\mu-\mu_0)^2)$$



**MAP: For which $\mu$ is?**

# Exercises: MAP for Gaussian Distribution

Using maximum a poserior estimation $\underset{\mu}{\text{argmax}} \prod_{i=1}^{n} p(\text{xi}|\mu)p(\mu)$

$$= \underset{\mu}{\text{argmax}} \prod_{i=1}^{n} \exp(-\frac{1}{2\sigma^2}(\text{xi}-\mu)^2)\exp(-\frac{1}{2\sigma^2}(\mu-\mu_0)^2)$$

$$= \underset{\mu}{\text{argmax}} \log(\prod_{i=1}^{n} \exp(-\frac{1}{2\sigma^2}(\text{xi}-\mu)^2)\exp(-\frac{1}{2\sigma^2}(\mu-\mu_0)^2))$$

$$= \underset{\mu}{\text{argmax}} -\sum_{i=1}^{n}(\text{xi}-\mu)^2 - (\mu-\mu_0)^2$$

$$= \underset{\mu}{\text{argmin}} \sum_{i=1}^{n}(\text{xi}-\mu)^2 + (\mu-\mu_0)^2$$

# Exercises: MAP for Gaussian Distribution

Addressing this optimization problem:

$$\underset{\mu}{\mathrm{argmin}} \sum_{i=1}^{n} (\mathrm{xi} - \mu)^2 + (\mu - \mu_0)^2$$



Derivation

$$\frac{d(\sum_{i=1}^{n}(\mathrm{xi}-\mu)^2+(\mu-\mu_0)^2)}{d\mu} = 2\sum_{i=1}^{n}(\mu - \mathrm{xi}) + 2(\mu - \mu_0) = 0$$

- The solution is

$$\mu = \frac{\sum_{i=1}^{n} \mathrm{xi} + \mu_0}{n+1}$$

# Exercises: Maximum A   Posteriori (MAP) Estimation

- Imagine you sent a message α to your friend that is either 1 or 0 with probability p and 1−p, respectively.

- Unfortunately that message gets corrupted by Gaussian noise N with zero mean and unit variance. Then what your friend would receive is a message X given by X=α+N.

# Exercises: Maximum A    Posteriori (MAP) Estimation

- Given that what your friend observed was that X takes a particular value x, that is X=x, <span style="color:red">he wants to know which was, probably, the value of α that you sent to him</span>.

By MAP, we sholud compute

$$\operatorname*{argmax}_{\alpha} \; p_X(x|\alpha) \, P(\alpha)$$

# Exercises: Maximum A Posteriori (MAP) Estimation

- $p_X(x|a=1)$

$$= p_N(x-1) = \sqrt{\frac{1}{2\pi}} \exp\left(-\frac{1}{2}(x-1)^2\right)$$

- $p_X(x|a=0)$

$$= p_N(x) = \sqrt{\frac{1}{2\pi}} \exp\left(-\frac{1}{2}(x)^2\right)$$

# Exercises: Maximum A Posteriori (MAP) Estimation

Therefore,

$$p_X(x|a=1) \, P(a=1) = p\sqrt{\frac{1}{2\pi}} \exp(-\frac{1}{2}(x-1)^2)$$

$$p_X(x|a=0) \, P(a=0) = (1-p)\sqrt{\frac{1}{2\pi}} \exp(-\frac{1}{2}(x)^2)$$

Then, $p_X(x|a=1) \, P(a=1) > p_X(x|a=0) \, P(a=0)$

**if and only if**

$$x > 0.5 + \log(1-p) - \log(p)$$

# Exercises: Maximum A Posteriori (MAP) Estimation

Therefore, the solution is

$$\underset{a}{\mathrm{argmax}} \ p_X(x|a) \ P(a) = 1,$$

$$\text{if } x > 0.5 + \log(1-p) - \log(p)$$

$$\underset{a}{\mathrm{argmax}} \ p_X(x|a) \ P(a) = 0,$$

$$\text{if } x < 0.5 + \log(1-p) - \log(p)$$

# Graphical Models

**What are Graphical Models?**

• They are diagrammatic representations of probability distributions – marriage between probability theory and graph theory.

• Also called probabilistic graphical models

# Graphical Models

**Why Graphical Models?**

- A graphical model is a method of modeling a probability distribution for reasoning under uncertainty, which is needed in applications such as speech recognition and computer vision.

# Graphical Models

**Why Graphical Models?**

- Graphical models allow us to address three fundamental questions:

1. How should I represent my data in a way that reflects domain knowledge while acknowledging uncertainty?

2. How do I make inferences from this data?

3. How can I learn the 'right' model for this data?

# Graphical Models

- Each of these questions can be rephrased as a question about probability distributions:

1. What is the joint distribution over my input variables? Which state configurations of the distribution are relevant to the problem?

2. How can we obtain the state probabilities? Do we use maximum-likelihood estimation, or can we use domain knowledge?

3. How can we compute conditional distributions of unobserved (latent) variable without needing to sum over a large number of state configurations?

Credicts from Professor Eric Xing

# Graphical Models

- Each of these questions can be rephrased as a question about probability distributions:

1. What is the joint distribution over my input variables? Which state configurations of the distribution are relevant to the problem?

2. How can we obtain the state probabilities? Do we use maximum-likelihood estimation, or can we use domain knowledge?

3. How can we compute conditional distributions of unobserved (latent) variable without needing to sum over a large number of state configurations?

# Graphical Models

What is a Graph?

- Consists of nodes (also called vertices) and links (also called edges or arcs).



- In a graphical model

**Example of Graph**

– Each node represents a random variable (or group of random variables)
– Links express probabilistic relationships between variables

# Graphical Models



**Example of Graph**

**Graphical Models in Engineering**

- Natural tool for handling Uncertainty and Complexity

– which occur throughout applied mathematics and engineering

- Fundamental to the idea of a graphical model is the notion of modularity

– a complex system is built by combining simpler parts

# Graphical Models

Role of Graphical Models in Machine Learning

- Simple way to visualize
– Structure of probabilistic model

- Insights into properties of model
– Conditional independence properties by inspecting graph

- Complex computations
– Required to perform inference and learning expressed as graphical manipulations

# Graphical Models

## Graph Directionality

- Directed graphical models
– Directionality associated with arrows



- Bayesian networks (BNs)
– Express causal relationships between random variables

# Graphical Models

## Graph Directionality

- Undirected graphical models
– Links without arrows



- Markov random fields (MRFs)
– Better suited to express soft constraints between variables

# Graphical Models

Bayesian Networks

- A Bayesian network (BN) is a graphical model for representing knowledge about an uncertain domain
– Where each node corresponds to a random variable and each edge represents the conditional probability for the corresponding random variables

- Directed graphs
– Used to describe probability distributions

Credits from Professor Sargur N. Srihari   58

# Graphical Models: Bayesian Networks

- Example: given random variables X, Y, and Z



This graph is associated with P(X,Y,Z) = P(Z|Y)P(Y|X)P(X)

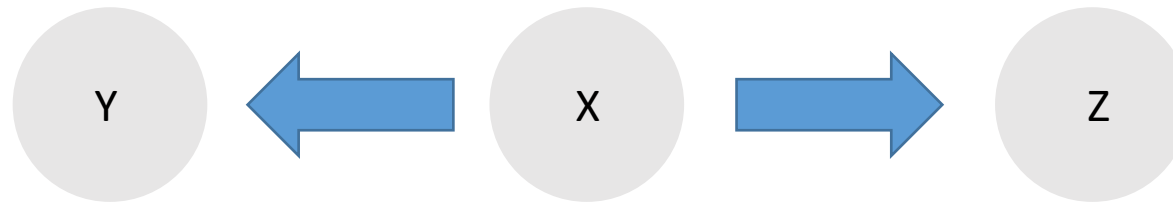It is not Bayesian Network

# Graphical Models: Bayesian Networks
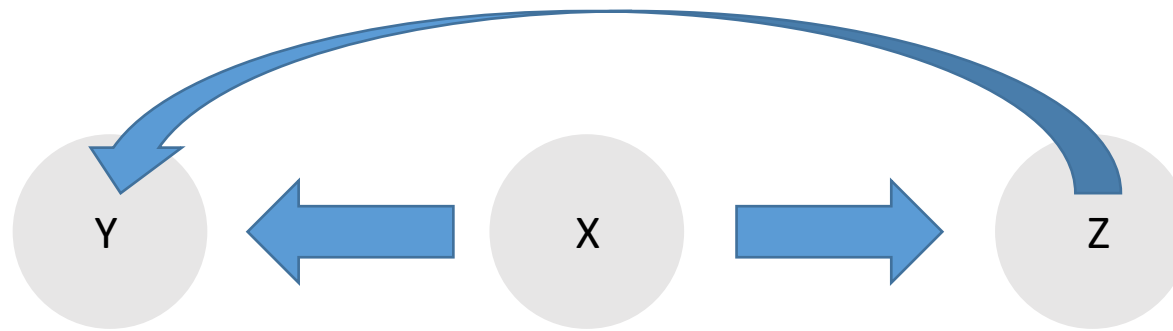
- Example: given random variables X, Y, and Z



This graph is associated with P(X,Y,Z) = P(Z|X,Y)P(Y|X)P(X)

It is Bayesian Network

# Graphical Models: Bayesian Networks

- Example: given random variables X, Y, and Z



This graph is associated with P(X,Y, Z) = P(Y|X)P(Z|X)P(X)

# Graphical Models: Bayesian Networks
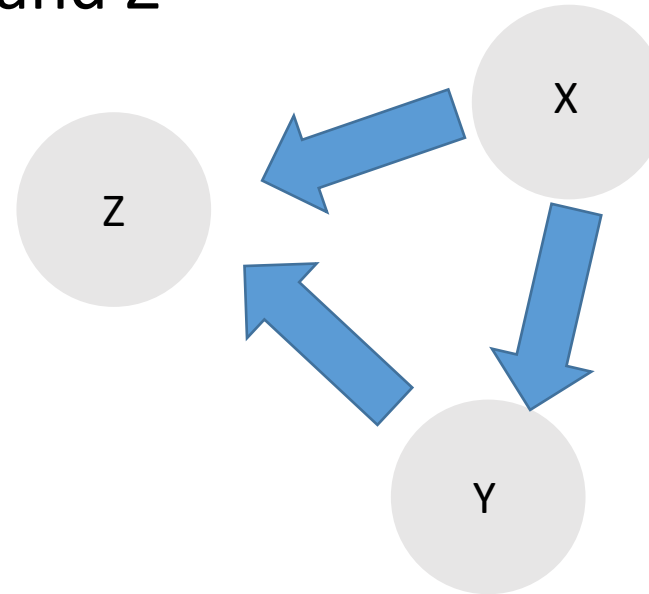
- Example: given random variables X, Y, and Z



This graph is associated with P(X,Y, Z) = P(Y|X)P(Z|X)P(X)

It is not Bayesian Network

# Graphical Models: Bayesian Networks

- Example: given random variables X, Y, and Z



This graph is associated with P(X,Y, Z) = P(Y|X)P(Z|X)P(X)

# Graphical Models: Bayesian Networks

- Example: given random variables X, Y, and Z



This graph is associated with P(X,Y, Z) = P(Y|X)P(Z|X)P(X)

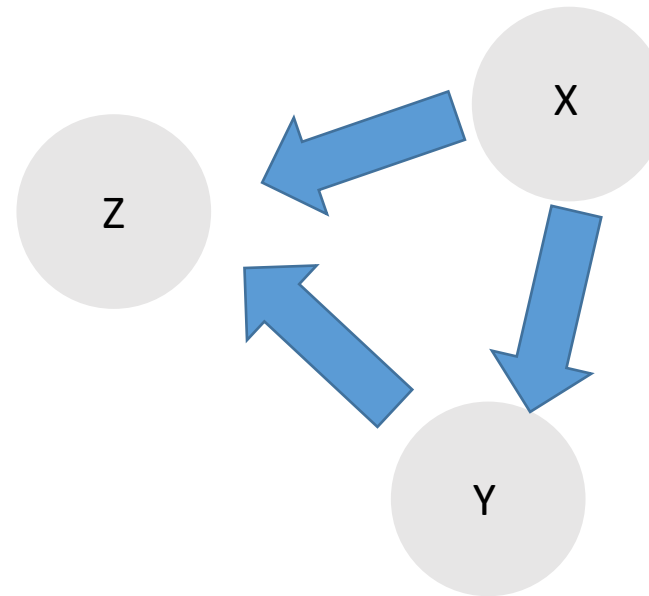It is not Bayesian Network

# Graphical Models: Bayesian Networks

- Example: given random variables X, Y, and Z



This graph is associated with P(X,Y,Z) = P(Y|X,Z)P(Z|X)P(X)

It is a Bayesian Network.

# Graphical Models: Bayesian Networks

• Example: given random variables X, Y, and Z



This graph is associated with P(X,Y,Z) = P(Z|X,Y)P(X)P(Y)

# Graphical Models: Bayesian Networks

- Example: given random variables X, Y, and Z



This graph is associated with P(X,Y,Z) = P(Z|X,Y)P(X)P(Y)

It is <span style="color:red">not</span> a Bayesian Network,

# Exercises: Bayesian Networks

- given random variables X, Y, and Z



Please compute P(X,Y,Z)

# Exercises: Bayesian Networks

– Introduce a node for each
random variable

– Associate each node with
conditional distribution

Answer: P(X,Y,Z) = P(Z|X,Y)P(Y|X)P(X)

# Exercises: Bayesian Networks

Given random variables X1, X2, X3, X4



Please compute P(X1,X2,X3,X4)

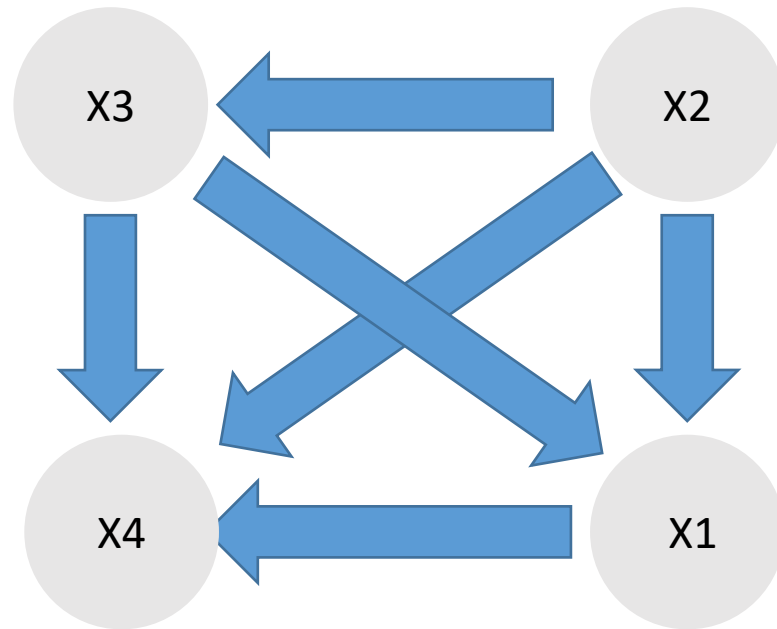# Exercises: Bayesian Networks



Answer: P(X1,X2,X3,X4) = P(X4|X1,X2,X3)P(X3|X1,X2)P(X1|X2)P(X2)

# Exercises: Bayesian Networks

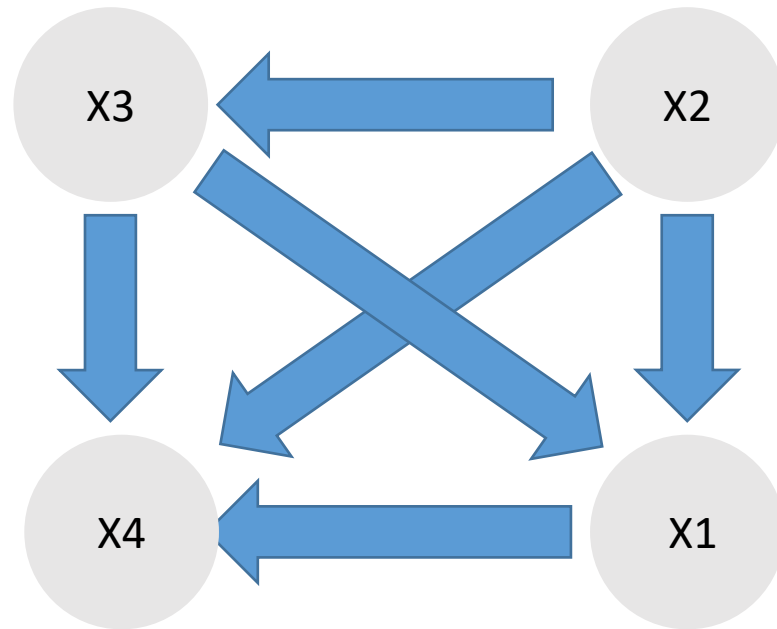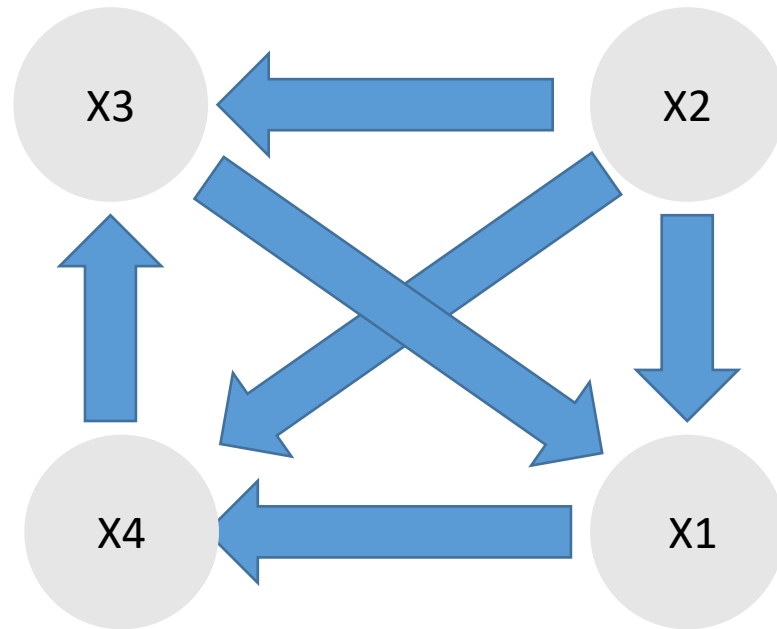Given random variables X1, X2, X3, X4



Please compute P(X1,X2,X3,X4)

# Exercises: Bayesian Networks



Answer: P(X1,X2,X3,X4) = P(X4|X1,X2,X3)P(X3|X2)P(X1|X2,X3)P(X2)
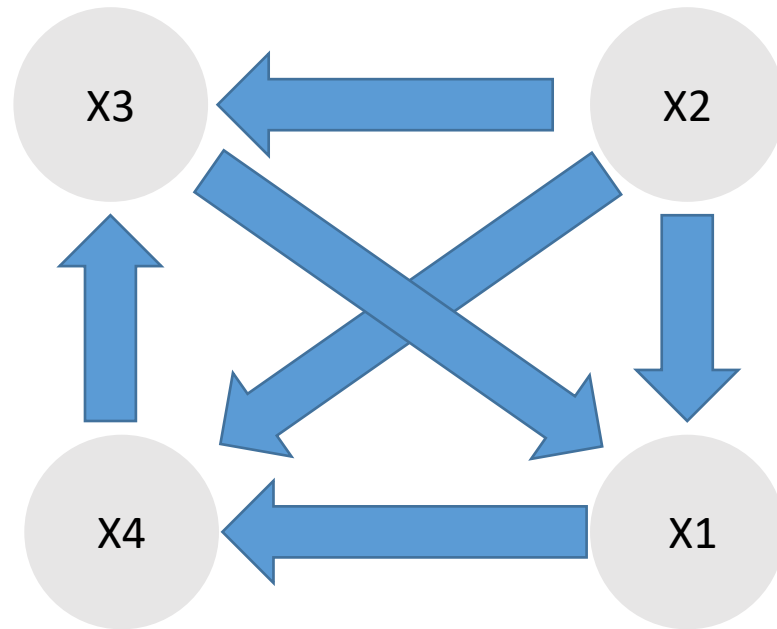
# Exercises: Bayesian Networks

Given random variables X1, X2, X3, X4



Please compute P(X1,X2,X3,X4)

# Exercises: Bayesian Networks



Answer: P(X1,X2,X3,X4) = P(X4|X1,X2)P(X1|X2,X3)P(X3|X2,X4)P(X2)

# Thank You!