

Data Analysis and Uncertainty

Part 1: Random Variables

Instructor: Sargur N. Srihari

University at Buffalo
The State University of New York
srihari@cedar.buffalo.edu

Topics

1. Why uncertainty exists?
2. Dealing with Uncertainty
3. Random Variables and Their Relationships
4. Samples and Statistical Inference

Reasons for Uncertainty

1. Data may only be a sample of population to be studied

Uncertain about extent to which samples differ from each other

2. Interest is in making a prediction about tomorrow based on data we have today

3. Cannot observe some values and need to make a guess

Dealing with Uncertainty

- Several Conceptual bases
 1. Probability
 2. Fuzzy Sets
 3. Rough Sets

Lack theoretical backbone and the wide acceptance of probability
- Probability Theory vs Probability Calculus
 - Probability Calculus is well-developed
 - Generally accepted axioms and derivations
 - Probability Theory has scope for perspectives
 - Mapping real world to what probability is

Frequentist vs Bayesian

- Frequentist
 - Probability is objective
 - It is the limiting proportion of times event occurs in identical situations
 - An idealization since all customers are not identical
- Bayesian
 - Subjective probability
 - Explicit characterization of all uncertainty including any parameters estimated from the data
- Frequently yield same results

Random Variable

- Mapping from property of objects to a variable that can take a set of possible values via a process that appears to the observer to have an element of unpredictability
 - Possible values of random variable is its domain
 - Examples
 - Coin toss (domain is the set [heads,tails])
 - No of times a coin has to be tossed to get a head
 - Domain is integers
 - Flying time of a paper aeroplane in seconds
 - Domain is set of positive real numbers

Properties of Univariate (Single) random variable

- X is random variable and x is its value
- Domain is finite:
 - probability mass function $p(x)$
- Domain is real line:
 - probability density function $p(x)$
- Expectation of X
 - $E[X] = \int x p(x) dx$

Multivariate Random Variable

- Set of several random variables
- d -dimensional vector $\mathbf{x} = \{x_1, \dots, x_d\}$
- Density function of \mathbf{x} is the joint density function
 $p(x_1, \dots, x_d)$
- Density function of single variable (or subset) is called a marginal density
 - Derived by summing over variables not included

$$p(x_1) = \iint p(x_1, x_2, x_3) dx_2 dx_3$$

Conditional Probability

- Density of a single variable (or a subset of complete set of variables) given (or ‘conditioned on’) particular values of other variables
 - Conditional density of variable X_1 given $X_2=6$
 - Conditional density of X_1 given some value of X_2 is denoted $f(x_1|x_2)$ and defined as

$$p(x_1 | x_2) = \frac{p(x_1, x_2)}{p(x_2)}$$

Supermarket Data

	Product A	Product B
Customer 1	0	1
Customer 2	1	1
Customer n=100,000		
Total	$n_A=10,000$	$n_B=5000$

Probability that randomly selected customer bought A is $n_A/n=0.1$

Probability that randomly selected customer bought B is $n_B/n=0.05$

n_{AB} = those who bought both A and B=10

$P(B=1|A=1)=10/10,000=0.001$

Probability of customer buying B reduces from 0.05 to 0.001 if we know customer bought product A

Conditional Independence

- Generic problem in data mining is finding relationships between variables
 - Is purchasing item A likely to be related to purchasing item B ?
- Variables are independent if there is no relationship; otherwise they are dependent
- Independent if $p(x,y)=p(x)p(y)$
- Equivalently $p(x|y)=p(x)$ or $p(y|x)=p(y)$ for all values of X and Y
 - (since $p(x,y)=p(x|y)p(y)$)

Conditional Independence: More than 2 variables

- X is conditionally independent of Y
 - given Z if for all values of X, Y, Z we have

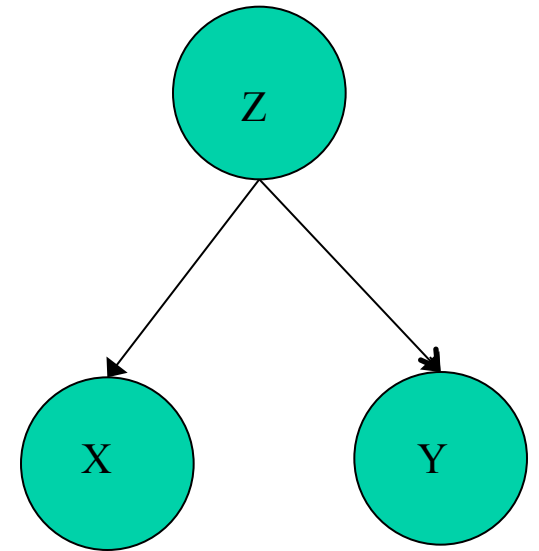
$$p(x, y | z) = p(x | z) p(y | z)$$

- Equivalently

$$p(x | y, z) = p(x | z)$$

Conditional Independence: Example

- Assume bread goes with either butter or cheese
- Person purchases bread ($Z=1$)
- Subsequent purchase of butter ($X=1$) and cheese ($Y=1$) are modeled as conditionally independent
 - Probability of purchasing cheese is unaffected by whether or not butter was purchased once we know bread was purchased



Conditional and Marginal Independence

- Conditional Independence need not imply marginal independence
 - If $p(x,y|z)=p(x|z)p(y|z)$
 - Then it need not imply
$$p(x,y)=p(x)p(y)$$
- We can expect butter & cheese to be dependent since both depend on bread
- Reverse also applies
 - X and Y may be unconditionally independent but conditionally dependent given Z
- Relationship of 2 variables masked by third

Interpreting Conditional Independence

- A and B are two different treatments

- Fraction who recover shown in table

	A	B
Old	2/10	30/90
Young	48/90	10/10

- Treatment B appears better

- Aggregate two rows

	A	B
Total	50/100	40/100

- Known as Simpson's Paradox

- First set conditioned on strata while second is unconditional
 - When two are combined sample size differences larger samples (old B, young A) dominate

Conditional Independence: Sequential Data

- Widely used when next value is dependent on past values in sequence
- Assumption of independence and conditional independence allow factoring joint density into tractable products of simpler densities
- First-Order Markov Model
 - Next value in a sequence is independent of all the past values given the current value in the sequence

$$p(x_1, \dots, x_n) = p(x_1) \prod_{j=2}^n p(x_j | x_{j-1})$$

On Assuming Independence

- Independence is a strong assumption frequently violated in practice
- But provides modeling gains
 - Understandable models
 - Fewer parameters
- Models are approximations of real world
- Benefits of appropriate independence assumptions outweigh more complex but stable models

Dependence and Correlation

- Covariance measures how X and Y vary together:
 - Large positive if large X is associated with large Y and small X with small Y
 - Negative if If large X is associated with small Y
 - Dividing by variance gives correlation
 - Referred to as linear dependency
- Two variables may be dependent but not linearly correlated

Correlation and Causation

- Two variables may be highly correlated without a causal relationship between the two
 - Yellow stained finger and lung cancer may be correlated but causally linked only by a third variable: smoking
 - Human reaction time and earned income are negatively correlated
 - Does not mean one causes the other
 - A third variable “age” is causally related to both

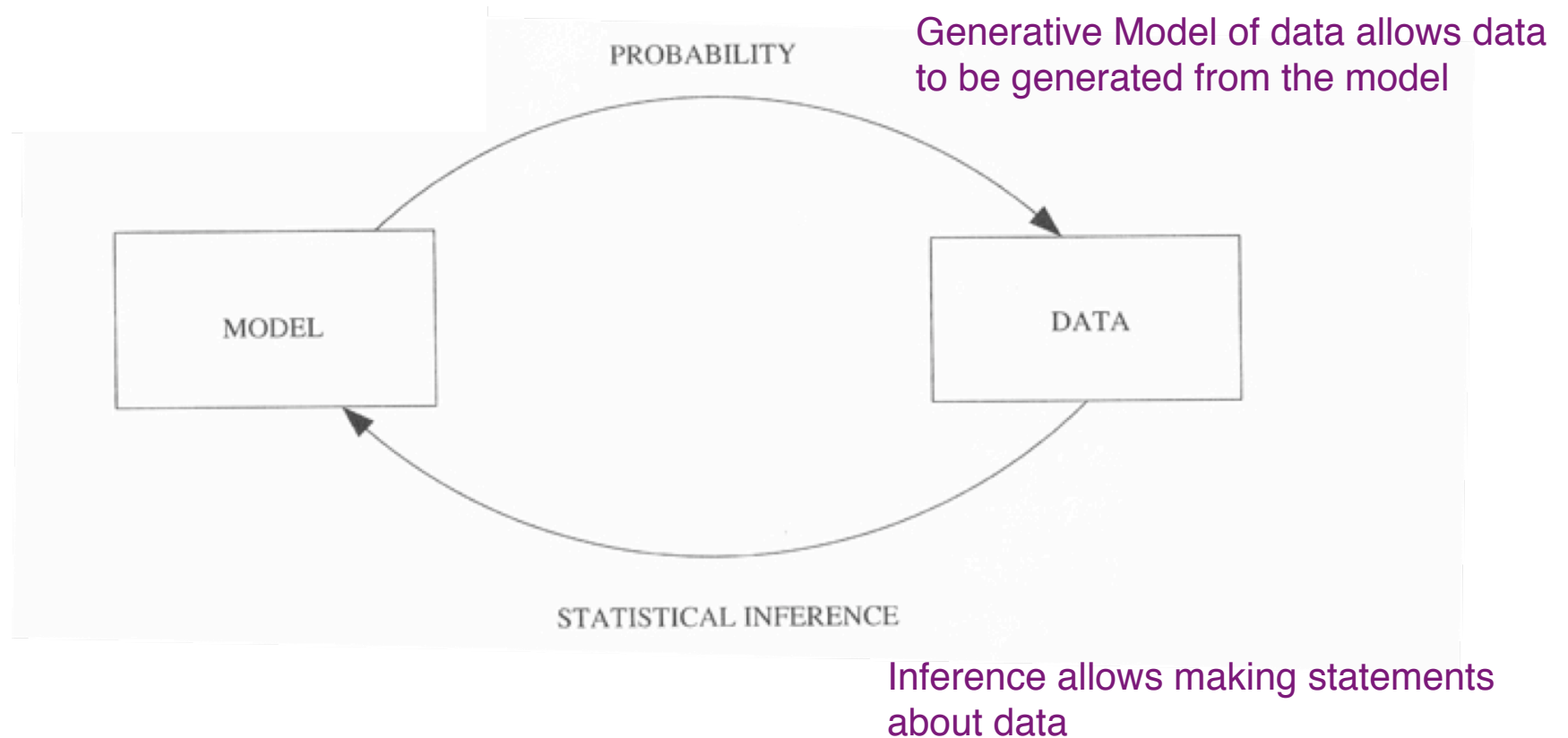
Causality Example: Hospitals

- In-house coronary bypass mortality rates
 - Regression: hospitals with more operations have lower rates
 - Conclusion: close low-surgery units
- Issues
 - Large hospitals might degrade with volume
 - Correlation because superior performance attracts more cases
 - No of cases and outcome are related by some other factor

Samples and Statistical Inference

- Samples Can Be Used To Model the Data
- Less appropriate if the goal is to detect small deviations from the bulk of the data

Dual Role of Probability and Statistics in Data Analysis



Testing goodness of the Model

- If $p(\mathbf{x}(i))$ is the probability of individual i having vector measurement $\mathbf{x}(i)$ (p could be probability density function)
- If the probability of each member of the population being included is independent
 - The overall probability of observing the entire distribution of values in the sample is the likelihood function

$$p(D | \theta, M) = \prod_{i=1}^n p(\mathbf{x}(i) | \theta, M)$$

where M is the model and θ are the parameters of the model

- Based on this probability we can decide how realistic the model is

Using a hypothesis test to accept or reject the model