

Data Analysis and Uncertainty

Part 3: Hypothesis Testing/Sampling

Instructor: Sargur N. Srihari

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

Topics

1. Hypothesis Testing

1. t-test for means
2. Chi-Square test of independence
3. Komogorov-Smirnov Test to compare distributions

2. Sampling Methods

1. Random Sampling
2. Stratified Sampling
3. Cluster Sampling

Motivation

- If a data mining algorithm generates a potentially interesting hypothesis we want to explore it further
- Commonly, value of a parameter
 - Is new treatment better than standard one
 - Two variables related in a population
- Conclusions based on a sample of population

Classical Hypothesis Testing

1. Define two complementary hypotheses

- Null hypothesis and alternative hypothesis
 - Often null hypothesis is a point value, e.g., draw conclusions about a parameter θ
 - Null Hypothesis is $H_0: \theta = \theta_0$ Alternative Hypothesis is $H_1: \theta \neq \theta_0$

2. Using data calculate a statistic

- Which depends on nature of hypotheses
- Determine expected distribution of chosen statistic
 - Observed value would be one point in this distribution

3. If in tail then either unlikely event or H_0 false

- More extreme observed value less confidence in null hypothesis

Example Problem

- Hypotheses are mutually exclusive
 - if one is true, other is false
- Determine whether a coin is fair
 - $H_0: P = 0.5$
 $H_a: P \neq 0.5$
- Flipped coin 50 times: 40 Heads and 10 Tails
 - Inclined to reject H_0 and accept H_a
- Accept/reject decision focuses on single test statistic

Test for Comparing two Means

- Whether population mean differs from hypothesized value
 - Called one-sample t test
- Common problem
 - Does Your Group Come from a Different Population than the One Specified?
- One-sample t-test (given sample of one population)
- Two-sample t-test (two populations)

One-sample t test

- Fix significance level in $[0,1]$, e.g., 0.01, 0.05, 1.0
- Degrees of freedom, $DF = n - 1$
 - n is no of observations in sample
- Compute test statistic (t-score) $t = \frac{\bar{x} - x}{s/\sqrt{n}}$
 - \bar{x} : sample mean, x : hypothesized mean (H_0), s std dev of sample
- Compute p-value from student's t-distribution
 - reject null hypothesis if p-value < significance level
- Used when population variances are equal / unequal, and with large or small samples.

Rejection Region

- Test statistic
 - mean score, proportion, t-score, z-score
- One and Two tail tests

Hyp Set	Null hyp	Alternative hyp	No of tails
1	$\mu = M$	$\mu \neq M$	2
2	$\mu \geq M$	$\mu < M$	1
3	$\mu \leq M$	$\mu > M$	1

- Values outside region of acceptance is **region of rejection**
- Equivalent Approaches: p-value and region of acceptance
- Size of region is significance level

Power of a Test

- Compare different test procedures
- Power of Test
 - Probability it will correctly reject a false null hypothesis ($1-\beta$)
 - False Negative Rate
- Significance of Test
 - Test's probability of *incorrectly* rejecting the null hypothesis (α)
 - True Negative Rate

Type 1 and Type 2 errors are denoted α and β

	Null is True	Null is False
Accept Null	$1-\alpha$ True Positive	β False Positive
Reject Null	α True Negative	$1-\beta$ False Negative

Likelihood Ratio Statistic

- Good strategy to find statistic is to use the Likelihood Ratio

- Likelihood Ratio Statistic to test hypothesis

$H_0: \theta = \theta_0$ $H_1: \theta \neq \theta_0$ is defined as

$$\lambda = \frac{L(\theta_0 | D)}{\sup_{\varphi} L(\varphi | D)} \text{ where } D = \{\mathbf{x}(1), \dots, \mathbf{x}(n)\}$$

i.e., Ratio of likelihood when $\theta = \theta_0$ to the largest value of the likelihood when θ is unconstrained

- Null hypothesis rejected when λ is small
- Generalizable to when null is not single point

Testing for Mean of Normal

- Given a sample of n points drawn from Normal with unknown mean and unit variance

- Likelihood under null hypothesis

$$L(0 | x(1), \dots, x(n)) = \prod_{i=1}^n p(x(i) | 0) = \prod_i \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x(i) - 0)^2\right)$$

- Maximum likelihood estimator is sample mean

$$L(\mu | x(1), \dots, x(n)) = \prod_{i=1}^n p(x(i) | \mu) = \prod_i \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x(i) - \bar{x})^2\right)$$

- Ratio simplifies to $\lambda = \exp(-n(\bar{x} - 0)^2 / 2)$

- Rejection region: $\{\lambda | \lambda \leq c\}$ for a suitably chosen c

- Expression written as $\bar{x} \geq \sqrt{-\frac{2}{n} \ln c}$

– Compare sample mean with a constant

Types of Tests used Frequently

- Differences between means
- Compare variances
- Compare observed distribution with a hypothesized distribution
 - Called goodness-of-fit test
- t-test for difference between means of two independent groups

Two sample t-test

- Whether two means have the same value

$x(1), \dots, x(n)$ drawn from $N(\mu_x, \sigma^2)$, $y(1), \dots, y(n)$ drawn from $N(\mu_y, \sigma^2)$

- $H_0: \mu_x = \mu_y$
- Likelihood Ratio statistic

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{s^2(1/n + 1/m)}}$$

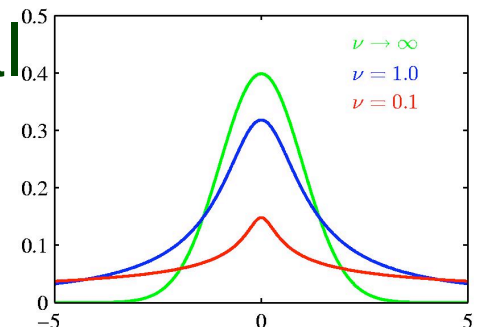
Difference between sample means adjusted by standard deviation of that difference

– with $s = s_x^2 \frac{n-1}{n+m-2} + s_y^2 \frac{m-1}{n+m-2}$

weighted sum of sample variances

– where $s_x^2 = \sum (x - \bar{x})^2 / (n-1)$

- t has t-distribution with $n+m-2$ degrees of freedom
- Test robust to departures from normal
- Test is widely used

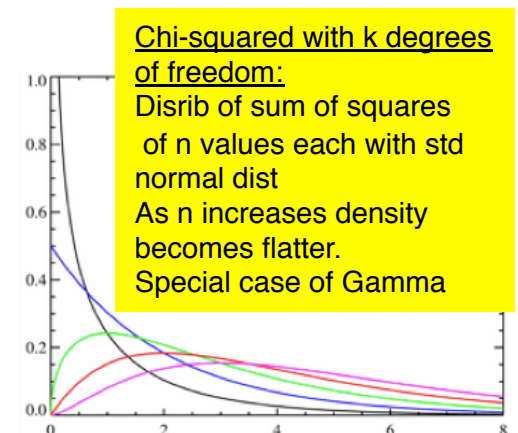


Test for Relationship between Variables

- Whether distribution of value taken by one variable is independent of value taken by another
- Chi-squared test
 - Goodness-of-fit test with null hypothesis of independence
- Two categorical variables
 - x takes values $x_i, i=1,\dots,r$ with probabilities $p(x_i)$
 - y takes values $y_j, j=1,\dots,s$ with probabilities $p(y_j)$

Chi-squared Test for Independence

- If x and y are independent $p(x_i, y_j) = p(x_i)p(y_j)$
- $n(x_i)/n$ and $n(y_j)/n$ are estimates of probabilities of x taking value x_i and y taking value y_j
- If independent estimate of $p(x_i, y_j)$ is $n(x_i)n(y_j)/n^2$
 - We expect to find $n(x_i)n(y_j)/n$ samples in (x_i, y_j) cell
 - Number the cells 1 to t where $t = r \cdot s$
 - E_k is expected number in k th cell, O_k is observed number
- Aggregation given by
$$X^2 = \sum_{k=1, t} \frac{(E_k - O_k)^2}{E_k}$$
- If null hyp holds X^2 has χ^2 distrib
 - With $(r-1)(s-1)$ degrees of freedom
 - Found in tables or directly computed



Testing for Independence Example

Outcome\Hospital	Referral	Non-referral
No improvement	43	47
Partial Improvement	29	120
Complete Improvement	10	118

- Medical data
- Whether outcome of surgery is independent of hospital type

- Total for referral=82
- Total with no improvement=90
- Overall total=367
- Under independence, top left cell has expected no= $82 \times 90 / 367 = 20.11$
- Observed number is 43. Contributes $(20.11 - 43)^2 / 20.11$ to χ^2
- Total value is $\chi^2 = 49.8$
- Comparing with χ^2 distribution with $(3-1)(2-1) = 2$ degrees of freedom reveals very high degree of significance
- Suggests that outcome is dependent on hospital type

Chi-squared Goodness of Fit Test

- Chi-squared test is more versatile than t-test.
Used for categorical distributions
- Used for testing normal distribution as well
 - E,g.
 - 10 measurements: $\{x_1, x_2, \dots, x_{10}\}$. They are supposed to be "normally" distributed with mean μ and standard deviation σ . You could calculate

$$\chi^2 = \sum_i \frac{(x_i - \mu)^2}{\sigma^2}$$

- we expect the measurements to deviate from the mean by the standard deviation, so: $|x_i - \mu|$ is about the same thing as σ .
- Thus in calculating chi-square we add-up 10 numbers that would be near 1. We expect it to approximately equal k , the number of data points. If chi-square is "a lot" bigger than expected something is wrong. Thus one purpose of chi-square is to compare observed results with expected results and see if the result is likely.

Randomization/permutation tests

- Earlier tests assume random sample drawn, to:
 - Make probability statement about a parameter
 - Make inference about population from sample
- Consider medical example:
 - Compare treatment and control group
 - H_0 : no effect (distrib of those treated same as those not)
 - Samples may not be drawn independently
- What difference in sample means would there be if difference is consequence of imbalance of populatns
- Randomization tests:
 - Allow us to make statements conditioned on input samples

Distribution-free Tests

- Other Tests assume form of distribution from which samples are drawn
 - Distribution-free tests replace values by ranks
 - Examples
 - If samples from same distribution: ranks well-mixed
 - If mean is larger, ranks of one larger than other
 - Test statistics, called *nonparametric tests*, are:
 - Sign test statistic
 - Kolmogorov- Smirnov Test Statistic
 - Rank sum test statistic
 - Wilcoxon test statistic

Kolmogorov Smirnov Test

- Determining whether two data sets have the same distribution
- To compare two CDFs, $S_{N_1}(x)$ and $S_{N_2}(x)$,
 - Compute the statistic D , the max of absolute difference between them:

$$D = \max_{-\infty < x < \infty} |S_{N_1}(x) - S_{N_2}(x)|$$

- which is mapped to a probability of similarity, P

$$P = Q_{KS} \left(\sqrt{N_e} + 0.12 + \left(0.11 \frac{0.11}{\sqrt{N_e}} \right) D \right)$$

- where the Q_{KS} function is defined as

$$Q_{KS}(\lambda) = 2 \sum_{j=1}^{\infty} (-1)^{j-1} e^{-2j^2\lambda^2}, \quad Q_{KS}(0) = 1, \quad Q_{KS}(\infty) = 0$$

- and N_e is the effective number of data points $N_e = \frac{N_1 N_2}{N_1 + N_2}$
 - where N_1 and N_2 are the no of data points in each

Comparing Hypotheses: Bayesian Perspective

- Done by comparison posterior probabilities

$$p(H_i | x) \propto p(x | H_i) p(H_i)$$

- Ratio leads to factorization in terms of prior odds and likelihood ratio

$$\frac{p(H_0 | x)}{p(H_1 | x)} \propto \frac{p(H_0)}{p(H_1)} \cdot \frac{p(x | H_0)}{p(x | H_1)}$$

- Some complications:
 - Likelihoods are marginal likelihoods obtained by integrating over unspecified parameters
 - Prior probs zero if parameter has specific value
 - Assign discrete non-zero values to given values of θ

Hypothesis Testing in Context

- In Data mining things are more complicated
 1. Data sets are large
 - Expect to obtain statistical significance
 - Slight departures from model will be seen as significant
 2. Sequential model fitting is common
 3. Results have various implications
 - Many models will be examined
 - e.g., discrete values for mean
 - If m true independent hypotheses are examined at 0.05 probability of incorrect rejection, with 100 hypotheses we are certain to reject one hypothesis
 - If we set family error rate as 0.05 , we get $\alpha = 0.0005$

Simultaneous Test Procedures

- Address difficulties of multiple hypotheses
- Bonferroni Inequality

$$p(a_1, a_2, \dots, a_n) \geq p(a_1) + p(a_2) + \dots + p(a_n) - n + 1$$

- By including other terms in the expansion, we can develop more accurate bounds but require knowledge of dependence relationships

Sampling in Data Mining

- Differences between Statistical inference and data mining
 - Data set in data mining may consists of entire population
 - Experimental design in statistics is concerned with optimal ways of collecting data
- More often database is sample of population
 - Contain records for every object in population but the analysis is based on a sample

Systematic Sampling

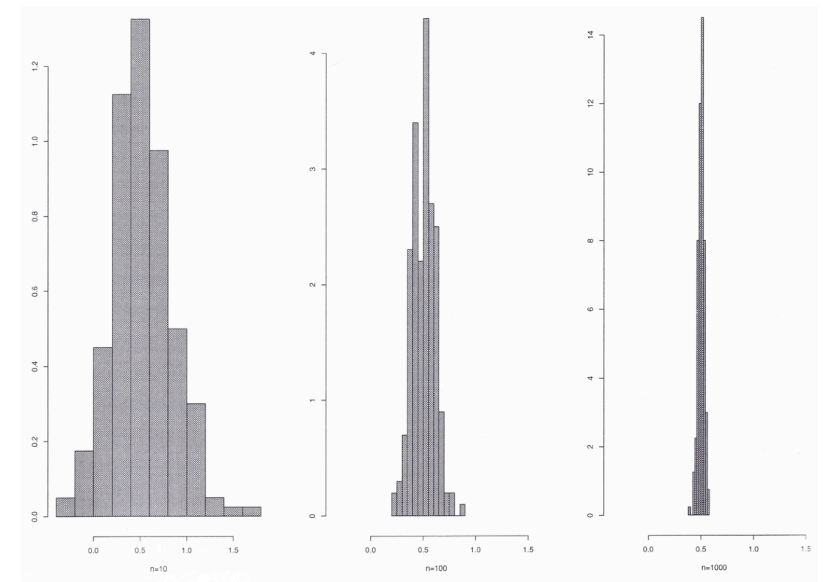
- Strategy for Representativeness
 - Taking one out of every two records
 - Sampling fraction =0.5
- Can lead to unexpected problems when there are regularities in database
 - E.g., data set where records are of married couples, one at a time

Random Sampling

- Avoiding regularities
- Epsem Sampling:
 - Equal probability of selection method
 - Each record has same probability of being chosen
 - Simple random sampling
 - n records are chosen from database of size N such that each set of n records has the same probability of being chosen

Results of Simple Random Sampling

- Population True mean = 0.5
- Samples of size $n = 10, 100$ and 1000
- Repeat procedure 200 times and plot histograms
 - Larger the sample more closely are values of sample mean are distributed about true mean



$n=10$

$n=100$

$n=1000$

Variance of Mean of Random Sample

- If variance of population of size N is σ^2 variance of mean of a simple random sample of size n without replacement is

$$\frac{\sigma^2}{n} \left(1 - \frac{n}{N}\right)$$

- Since second term is small, variance decreases as sample size increases
- When sample size is doubled standard deviation is reduced by factor of $\sqrt{2}$
- Estimate σ^2 from the sample using

$$\sum (x(i) - \bar{x})^2 / (n - 1)$$

Stratified Random Sampling

- Split population into non-overlapping subpopulations or strata
- Advantages:
 - Enables making statements about subpopulations
 - If strata are relatively homogeneous, most of variability is accounted for by differences between strata

Mean of Stratified Sample

- Total size of population is N
- k^{th} stratum has N_k elements in it
- n_k are chosen for the sample from this stratum
- Sample mean within k^{th} stratum is \bar{x}_k
- Estimate of Population Mean:

$$\sum \frac{N_k \bar{x}_k}{N}$$

- Variance of estimator

$$\frac{1}{N^2} \sum N_k^2 \text{var}(\bar{x}_k)$$

Cluster Sampling

- Hierarchical Data
 - Letters occur in words which lie in sentences grouped into paragraphs occurring in chapters forming books sitting in libraries
- Simple random sample is difficult
- To draw a sample of elements
 - Instead draw a sample of units that contain several elements

Unequal Cluster Sizes

- Size of random sample from k^{th} cluster is n_k
- Sample mean r

$$\sum x_k / \sum n_k$$

- Overall sampling fraction is f
 - (which is small and ignorable)
- Variance of r

$$\frac{1-f}{(\sum n_k)^2} \frac{a}{1-a} (\sum s_k^2 + r^2 \sum n_k^2 - 2r \sum s_k n_k)$$

Conclusion

- Nothing is certain
- Data mining objective is to make discoveries from data
- We want to be as confident as we can that our conclusions are correct
- Fundamental tool is probability
 - Universal language for handling uncertainty
 - Allows us to obtain best estimates even with data inadequacies and small samples