# Nonparametric Statistics

**References**

Some good references for the topics in this course are

1. Higgins, James (2004), *Introduction to Nonparametric Statistics*

2. Hollander and Wolfe, (1999), *Nonparametric Statistical Methods*

3. Arnold, Steven (1990), *Mathematical Statistics* ( Chapter 17)

4. Hettmansperger, T and McKean, J (1998), *Robust nonparametric Statistical Methodology*

5. Johnson, Morrell, and Schick (1992), Two-Sample Nonparametric Estimation and Confidence Intervals Under Truncation, *Biometrics*, 48, 1043-1056.

6. Beers, Flynn, and Gebhardt (1990), Measures of Location and Scale for Velocities in Cluster of Galaxies-A Robust Approach. *Astron Jr*, 100, 32-46.

7. Lecture Notes by Steven Arnold (2005) Summer School in Statistics for Astronomers. <u>Website</u>: http://astrostatistics.psu.edu/

8. Lecture Notes by Tom Hettmansperger (2007), Nonparametrics.zip. <u>Website</u>: http://astrostatistics.psu.edu/

9. <u>Website</u>: http://www.stat.wmich.edu/slab/RGLM/

[**\*** These notes borrow heavily from the material in 7 and 8]

## 1. Parametric and Nonparametric models

A *parametric statistical model* is a model where the joint distribution of the observations involves several unknown constants called *parameters*. The functional form of the joint distribution is assumed to be known and the only unknowns in the model are the parameters. Two parametric models commonly encountered in astronomical experiments are

1. The Poisson model in which we assume that the observations are independent Poisson random variables with unknown common mean $\theta$.

2. The normal model in which the observations are independently distributed with unknown mean $\mu$ and unknown variance $\sigma^2$.

In the first model $\theta$ is the parameter and in the second $\mu$ and $\sigma^2$ are the parameters. Anything we can compute from the observations is called a *statistic*. In parametric statistics the goal is to use observations to draw inference about the unobserved parameters and hence about the underlined model.

A *nonparametric model* is the one in which no assumption is made about the functional form of the joint distribution. The only assumption made about the observations is that they are independent identically distributed (i.i.d.) from an arbitrary continuous distribution. As a result, the nonparametric statistics is also called *distribution free* statistics. There are no parameters in a nonparametric model.

A *semiparametric model* is the one which has parameters but very weak assumptions are made about the actual form of the distribution of the observations.

Both nonparametric and semiparametric models used to be (and often still are) lumped together and called nonparametric models.

## 2. Why Nonparametric? Robustness

While in many situations parametric assumptions are reasonable (e.g. assumption of Normal distribution for the background noise, Poisson distribution for a photon counting signal of a nonvariable source), we often have no prior knowledge of the underlying distributions. In such situations, the use of parametric statistics can give misleading or even wrong results.

We need statistical procedures which are insensitive to the model assumptions in the sense that the procedures retain their properties in the neighborhood of the model assumptions.

Insensitivity to model assumptions : **Robustness**

In particular, for

- Hypothesis Testing

    We need test procedures where

    – the level of significance is not sensitive to model assumptions (Level Robustness).
    – the statistical power of a test to detect important alternative hypotheses is not sensitive to model assumptions (Power Robustness).

- Estimation

    The estimators such that

    – the variance (precision) of an estimator is not sensitive to model assumptions (Variance Robustness).

Apart from this, we also need procedures which are robust against the presence of outliers in the data.

**Eg:**

1. The sample mean is not robust against the presence of even one outlier in the data and is not variance robust as well. The sample median is robust against outliers and is variance robust.

2. The t-test does not have t-distribution if the underlined distribution is not normal and the sample size is small. For large sample size, it is asymptotically level robust but is not power robust. Also, it is not robust against the presence of outliers.

Procedures derived for nonparametric and semiparametric models are often called *robust* procedures since they are dependent only on very weak assumptions.

## 3. Nonparametric Density Estimation

Let $X_1, X_2, \cdots, X_n$ be a random sample from an unknown probability density function $f$. The interest is to estimate the density function $f$ itself.

Suppose the random sample is drawn from a distribution with known probability density function, say Normal with mean $\mu$ and variance $\sigma^2$. The density $f$ can then be estimated by estimating the values of the unknown parameters $\mu$ and $\sigma^2$ from the data and substituting these estimates in the expression for normal density. Thus the *parametric* density estimator is

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

where $\hat{\mu} = \dfrac{\sum_i x_i}{n}$ and $\hat{\sigma}^2 = \dfrac{\sum_i (x_i - \hat{\mu})^2}{n - 1}$.

In case of the *nonparametric* estimation of the density function, the functional form of the density function is not assumed to be known. We, however, assume that the

underlined distribution has a probability density $f$ and determine its form based on the data at hand.

The oldest and widely used *nonparametric density estimator* is the histogram. Given an origin $x_0$ and a *bandwidth $h$*, we consider the intervals of length $h$, also called *bins*, given by $B_i = [x_0 + mh, x_0 + (m_1)h)$ where $m = 0, \pm 1, \pm 2, \cdots$ and define the histogram by

$$\hat{f}_n(x) = \frac{1}{nh}[\text{ number of observations in the same bin as x}]$$
$$= \frac{1}{nh}\sum_{i=1}^{n} n_j I[x \in B_j]$$

where $n_j$ = number of observations lying in bin $B_j$.

Though it is a very simple estimate, the histogram has many drawbacks, the main one is that we are estimating a continuous function by a discrete function. Also, it is not robust against the choice of origin $x_0$ and bandwidth $h$.

**Kernel Density Estimation**

We consider a specified *kernel function $K(.)$* satisfying the conditions

- K(.) is symmetric around 0

- $\int_{-\infty}^{\infty} K(x)dx = 1$

and define the *kernel density estimator* by

$$\hat{f}_n(x) = \frac{1}{nh} \sum_{i=1}^{n} K\left(\frac{x - X_i}{h}\right).$$

The estimate of $f$ at point $x$ is obtained using a weighted function of observations in the $h$-neighborhood of $x$. The weight given to each of the observation in the neighborhood depends on the choice of kernel function. Some kernel functions are

- Uniform kernel: $K(u) = \frac{1}{2}I[|u| \le 1]$

- Triangle kernel: $K(u) = (1 - |u|)I[|u| \le 1]$

- Epanechnikov kernel: $K(u) = \frac{3}{4}(1 - u^2)I[|u| \le 1]$

- Gaussian kernel: $K(u) = \frac{1}{\sqrt{2\pi}}\exp(-\frac{1}{2}u^2)$

The kernel density estimator satisfies the property

$$\int_{-\infty}^{\infty} \hat{f}_n(x)dx = 1$$

and on the whole gives a better estimate of the underlined density. Some of the properties are

- The kernel estimates do not depend on the choice of the origin, unlike histogram.
- The kernel density estimators are 'smoother' than the histogram estimators since they inherit the property of the smoothness of the kernel chosen.
- The kernel density estimator has a faster rate of convergence.
- Increasing the bandwidth is equivalent to increasing the amount of smoothing in the estimate. Very large $h(\to \infty)$ would give a flat estimate and $h \to 0$ will lead to a needlepoint estimate giving a noisy representation of the data.
- The choice of the kernel function is not very crucial. The choice of the bandwidth, however, is crucial and the optimal bandwidth choice is extensively discussed and derived in the literature. For instance, with Gaussian kernel, the optimal (MISE) bandwidth is

$$h_{\text{opt}} = 1.06\sigma n^{-\frac{1}{5}}$$

where $\sigma$ is the population standard deviation, which is estimated from the data.

## 4. Some Nonparametric Goodness-of-fit Tests

At times, though the samples are drawn from unknown populations, the investigators wish to confirm whether the data fit some proposed model. The Goodness-of-fit tests are the useful procedures to confirm whether the proposed model satisfactorily approximates the observed situation. Apart from the usual Chi-Square goodness of fit test, we have Kolmogorov-Smirnov tests which are discussed here.

### One-sample Kolomogorov-Smirnov Test

This is a test of hypothesis that the sampled population follows some specified distribution.

Suppose we observe $X_1, ..., X_n$ i.i.d. from a continuous distribution function $F(x)$. We want to test the null hypothesis that $F(x) = F_0(x)$ for all $x$, against the alternative that $F(x) \neq F_0(x)$ for some $x$, where $F_0$ is a distribution which is completely specified before we collect the data. Let $\widehat{F}_n(x)$ be the empirical distribution function (e.d.f.) defined by

$$\widehat{F}_n(x) = \frac{1}{n} \sum_{i=1}^{n} I[X_i \leq x]$$

The one sample *Kolmogorov-Smirnov* (KS) statistic is

$$M = \max_x \left| \widehat{F}_n(x) - F_0(x) \right|$$

We want to reject if $M$ is too large.

It is not hard to show that the exact null distribution of $M$ is the same for all $F_0$, but different for different $n$. Table of critical values are given in many books. A large sample result is for large $n$

$$P(nM > q) \overset{\bullet}{=} 2 \sum_{i=1}^{\infty} (-1)^{i-1} \exp\left(-2i^2 q^2\right) \overset{\bullet}{=} 2 \exp\left(-2q^2\right)$$

Use of the last formula is quite accurate and conservative. There for a size $\alpha$ test we reject if

$$nM > \left(-\frac{1}{2}\log\left(\frac{\alpha}{2}\right)\right)^{1/2} = M^\alpha$$

We can also construct a confidence band for the distribution as we now show. First note that the distribution of

$$M(F) = \max_x \left|\widehat{F}_n(x) - F(x)\right|$$

is the same as null distribution for the K-S test statistic. Therefore

$$1 - \alpha = P(M(F) \le M^\alpha) = P\left(\left|\widehat{F}_n(x) - F(x)\right| \le \frac{M^\alpha}{n} \text{ for all x}\right)$$

$$= P\left(F(x) \in \widehat{F}_n(x) \pm \frac{M^\alpha}{n} \text{ for all x}\right).$$

On situation in which K-S is misused is in testing for normality. The problem is that for K-S to be applied, the distribution $F_0$ must be completely specified before we collect the data. In testing for normality, we have to choose the mean and the variance based on the data. This means that we have chosen a normal distribution which is a closer to the data than the true $F$ so that $M$ is too small. We must adjust the critical value to adjust for this as we do in $\chi^2$ goodness of fit tests. Lilliefors has investigated the adjustment of p-values necessary to have a correct test for this situation and shown that the test is more powerful than the $\chi^2$ gladness of fit test for normality. The Anderson-Darling and Shapiro-Wilk tests are specifically designed to test for normality.

Another test of this kind for testing $F = F_0$ is the Cramer-von Mises test based on

$$\int_{-\infty}^{\infty} \left(\widehat{F}_n(x) - F_0(x)\right)^2 dF_0$$

**Two-sample Kolmogorov-Smirnov Test**

Alternatively, one may be interested in verifying whether two independent samples come form identically distributed populations.

For this problem, we have two samples $X_1, ..., X_m$ and $Y_1, ..., Y_n$ from continuous distribution functions $F(x)$ and $G(y)$. We want to test the null hypothesis that $F(x) = G(x)$ for all $x$ against the alternative that $F(x) \neq G(x)$ for some $x$. Let $\widehat{F}_n(x)$ and $\widehat{G}_n(y)$ be the empirical distribution functions (edf's) for the $x's$ and $y's$. The two sample *Kolmogorov-Smirnov* (K-S) test is based on

$$M = \max_x \left| \widehat{F}_n(x) - \widehat{G}_n(x) \right|$$

We reject if $M$ is too large. As in the one sample case if $n$ and $m$ are large,

$$P(dM > q) \overset{\bullet}{=} 2 \sum_{i=1}^{\infty} (-1)^{i-1} \exp\left(-2i^2 q^2\right) \overset{\bullet}{=} 2 \exp\left(-2q^2\right)$$

(where $d = 1/\left(\frac{1}{m} + \frac{1}{n}\right)$) so that critical values may be determined easily.

**5. Nonparametric Tests and Confidence Intervals**

The nonparametric tests described here are often called distribution free procedures because their significance levels do not depend on the underlying model assumption i.e., they are level robust. They are also power robust and robust against outliers.

We will mainly discuss the so-called rank procedures. In these procedures, the observations are jointly ranked in some fashion. In using these procedures, it is occasionally important that the small ranks go with small observations, though often it does not matter which order we rank in. The models for these procedures are typically semi-parametric models.

One advantage of using ranks instead of the original observations is that the ranks are not affected by monotone transformations. Hence there is no need of transforming the observations before doing a rank procedure. Another advantage of replacing the observations with the ranks is that the more extreme observations are pulled in closer to the other observations.

As a consequence, a disadvantage is that nearby observations are spread out.

For example

$$
\begin{array}{ccccccccc}
Obs & 1 & 1.05 & 1.10 & 2 & 3 & 100 & 1,000,00 \\
Rank & 1 & 2 & 3 & 4 & 5 & 6 & 7
\end{array}
$$

The main reason we continue to study these rank procedures is the power of the procedures. Suppose the sample size is moderately large. If the observations are really normally distributed, then the rank procedures are nearly as powerful as the parametric ones (which are the best for normal data). In fact it can be shown that Pitman asymptotic relative efficiency (ARE) of the rank procedure to the parametric procedure is

$$3/\pi = .95$$

and in fact the ARE is always greater than $3/\pi$. However the ARE is $\infty$ for some non-normal distributions. What this means is the rank procedure is never much worse that parametric procedure, but can be much better.

**Ties:**

We assume that the underlined probability distribution is continuous for the rank procedures and hence, theoretically, there are no ties in the sample. However, the samples often have ties in practice and procedures have been developed for dealing with these ties. They are rather complicated and not uniquely defined so we do not discuss them here. (refer Higgins for details).

## 5.1 Single Sample Procedures

We introduce the concept of *location parameter* first.

A population is said to be located at $\mu_0$ if the population median is $\mu_0$.

Suppose $X_1, \cdots, X_n$ is a sample from the population. We say that $X_1, \cdots, X_n$ is located at $\mu$ if $X_1 - \mu, \cdots, X_n - \mu$ is located at 0.

Thus any statistic

$$S(\mu) = S(X_1 - \mu, \cdots, X_n - \mu)$$

is useful for the location analysis if $E[S(\mu_0)] = 0$ when the population is located at $\mu_0$. This simple fact leads to some test procedures to test the hypothesis of population locations.

## Sign Test

This is one of the oldest nonparametric procedure where the data are converted to a series of plus and minus signs. Let $S(\mu)$ be the sign statistic defined by

$$
\begin{aligned}
S(\mu) &= \sum_{i=1}^{n} sign(X_i - \mu) \\
&= [\#X_i > \mu] - [\#X_i < \mu] \\
&= S^+(\mu) - S^-(\mu) \\
&= 2S^+(\mu) - n
\end{aligned}
$$

To find a $\hat{\mu}$ such that $S(\hat{\mu}) = 0$, we get $\hat{\mu} = \text{median}(X_i)$. Thus if $\mu_0$ is the median of the population, we expect $E[S(\mu_0)] = 0$.

Suppose we wish to test the hypothesis that the population median is $\mu_0$. Thus we have

$$H_0 : \mu = \mu_0 \quad \text{against} \quad H_1 : \mu \neq \mu_0$$

and based on $S(\mu_0)$, the proposed decision rule is:

$$\text{Reject } H_0 \text{ if } |S(\mu_0)| = |2S^+(\mu_0) - n| \geq c$$

where c is chosen such that

$$P_{\mu_0}[|2S^+(\mu_0) - n| \geq c] = \alpha.$$

It is easy to see that under $H_0 : \mu = \mu_0$, the distribution of $S^+(\mu_0)$ is Binomial$\left(n, \frac{1}{2}\right)$ irrespective of the underlined distribution of $X_i$'s and hence $c$ can be chosen appropriately. Equivalently, we reject $H_0$ if

$$S^+(\mu_0) \leq k \quad \text{or} \quad S^+(\mu_0) \geq n - k$$

where

$$P_{\mu_0}[S^+(\mu_0) \leq k] = \frac{\alpha}{2}.$$

This fact can be used to construct a confidence interval for the population median $\mu$. Consider

$$P_d[k < S^+(d) < n - k] = 1 - \alpha$$

and find the smallest $d$ such that [number of $X_i > d] < n - k$. Suppose we get

$$d = X_{(k)} \quad : \quad [\#X_i > X(k)] = n - k$$
$$d_{min} = X_{(k+1)} \quad : \quad [\#X_i > X(k+1)] = n - k - 1.$$

On the same lines, we find $d_{max} = X_{(n-k)}$. Then a $(1 - \alpha)100\%$ distribution-free confidence interval for $\mu$ is given by $[X_{(k+1)}, \ X_{(n-k)}]$

Note that the median is a robust measure of location and does not get affected by the outliers. The sign test is also robust and insensitive to the outliers and hence the confidence interval is robust too.

**Wilcoxon Signed Rank test**

The sign test above utilizes only the signs of the differences between the observed values and the hypothesized median. We can use the signs as well as the ranks of the differences, which leads to an alternative procedure.

Suppose $X_1, \cdots, X_n$ is a random sample from an unknown population with median $\mu$. Further, suppose we wish to test the hypothesis that $\mu = \mu_0$ against the alternative that $\mu \neq \mu_0$.

We define $Y_i = X_i - \mu_0$ and first rank the absolute values of $|Y_i|$. Let $R_i$ be the rank of the absolute value of $Y_i$ corresponding to the $i^{th}$ observation. The signed rank of an observation is the rank of the observation times the sign of the corresponding $Y_i$. Let

$$
S_i = \begin{cases} 1 & \text{if } (X_i - \mu_0) > 0 \\ 0 & \text{otherwise.} \end{cases}
$$

By arguments similar to the one mentioned for earlier test, we can construct a test using the statistic

$$
WS = \sum_{i=1}^{n} S_i R_i.
$$

$WS$ is called the *Wilcoxon signed rank statistic.*

Note that $WS$ is the sum of ranks with positive sign of $Y_i$, i.e. the positive signed ranks. If $H_0$ is true, the probability of observing a positive difference $Y_i = X_i - \mu_0$ of given magnitude is equal to the probability of observing a negative difference of same magnitude. Under the null hypothesis the sum of positive signed ranks is expected to have the same value as that of negative signed ranks. Thus a large or a small value of $WS$ indicates a departure from the null hypothesis. Hence we reject the null hypothesis if $WS$ is too large or too small.

The critical values of the Wilcoxon Signed Rank test statistic are tabulated for various

sample sizes. The tables of exact distribution of $WS$ based on permutations is given in Higgins(2004).

**Normal approximation**

It can be shown that for large sample, the null distribution of $WS$ is approximately normal with mean $\mu$ and variance $\sigma^2$ where

$$\mu = \frac{n(n+1)}{4}, \ \ \sigma^2 = \frac{n(n+1)(2n+1)}{24}$$

and the Normal cut-off points can be used for large values of $n$.

**Hodges-Lehmann confidence Interval for $\mu$**

We can construct a $(1-\alpha)100\%$ confidence interval for population median $\mu$ using Wilcoxon Signed rank statistic, under the assumption that the underlined population is symmetric around $\mu$.

Let

$$W_{ij} = \frac{X_i + X_j}{2}, \ \ n \geq i \geq j \geq 1.$$

be the average of the $i^{th}$ and $j^{th}$ original observations, called a *Walsh average.*

For example, consider a single sample with 5 observations $X_1, \cdots, X_5$ given by $-3, 1, 4, 6, 8$. Then the Walsh averages are

|     | -3  | 1   | 4   | 6   | 8   |
|-----|-----|-----|-----|-----|-----|
| -3  | -3  | -1  | .5  | 1.5 | 2.5 |
| 1   |     | 1   | 2.5 | 3.5 | 4.5 |
| 4   |     |     | 4   | 5   | 6   |
| 6   |     |     |     | 6   | 7   |
| 8   |     |     |     |     | 8   |

We order the $W_{ij}$ according to their magnitude and let $U_{[i]}$ be the $i^{th}$ largest $W_{ij}$.

The median of $W_{ij}$'s provides a point estimation of the population median $\mu$. This median of Walsh averages is known as the *Hodges-Lehmann* estimator of the population median $\mu$.

Using the Walsh averages, it is easy to see that another representation for the Wilcoxon Signed Rank statistic is

$$WS = \# \left( W_{ij} \geq 0 \right)$$

(Note that this definition gives $WS = 13$ for the example.)

Now suppose that we do not know $\mu$. Define

$$WS \left( \mu \right) = \# \left( W_{ij} \geq \mu \right)$$

Then the general distribution of $WS \left( \mu \right)$ is the same as null distribution $WS$ statistic. Suppose that a size $1 - \alpha$ two-sided Wilcoxon Signed Rank test for $\mu = 0$ accepts the null hypothesis if

$$a \leq WS < b,$$

where $a$ and $b$ depend on $\alpha$. Then a $(1 - \alpha)100\%$ confidence interval for $\mu$ is

$$a \leq WS \left( \mu \right) < b \quad \Leftrightarrow \quad U_{[a]} < \mu \leq U_{[b]}$$

This confidence interval is called the *Hodges-Lehmann confidence interval* for $\theta$

For the data above, it can be seen from the table values that the acceptance region for a $\alpha = .125$ test is

$$2 \leq WS < 14$$

so that

$$U_{[2]} < \mu \leq U_{[14]} \quad \Leftrightarrow \quad -1 < \mu \leq 7$$

is a 87.5% confidence interval for $\mu$. Note that the assumed continuity implies that the inequality can be replaced by an equality in the last formula (but not the one before it) or vice versa.

Note that the H-L interval is associated with the Wilcoxon signed rank test in that the two-sided Wilcoxon test rejects $\mu = 0$ iff 0 is not in the confidence interval. Also note that there is no problem with ties in either the H-L confidence interval or H-L estimator.

## 5.2 Two Sample Procedures

Suppose we observe two independent random samples: $X_1, ..., X_n$ from distribution function $F(x)$, and $Y_1, ..., Y_n$ from distribution $G(y)$ where both $F$ and $G$ are continuous distributions.

We discuss the nonparametric procedure for making inference about the difference between the two location parameters of $F$ and $G$ here. In particular, we make the assumption that the distribution functions of the two populations differ only with respect to the location parameter, if they differ at all. This can alternatively be stated by expressing $G(y) = F(y + \delta)$ where $\delta$ is the difference between the medians here. This situation is often called a *shift* family.

There is no symmetry assumption in the two sample model. The continuity of the distributions implies there will be no ties.

### Wilcoxon rank sum statistic

Consider testing that $\delta = 0$ against $\delta \neq 0$. We first combine and jointly rank all the observations. Let $R_i$ and $S_j$ be the ranks associated with $X_i$ and $Y_j$. Then we could compute a two-sample t based on these ranks. However, an equivalent test is based

16

on

$$H = \sum_{i=1}^{n} R_i$$

Note that if $\delta > 0$, then the $X_i's$ should be greater than the $Y's$, hence the $R_i's$ should be large and hence $H$ should be large. A similar motivation works when $\delta < 0$. Thus we reject the null hypothesis $H_0 : \delta = 0$ if $H$ is too large or too small. This test is called the *Wilcoxon rank-sum test*.

Tables of permutation (exact) distribution of $H$ are available in Higgins (p 340).

For example, suppose we have two independent random samples of size 4 and 3. Suppose further that we observe 37, 49, 55, 57 in the first sample and 23, 31, 46 in the second. We get

| *obs* | 37 | 49 | 55 | 57 | 23 | 31 | 46 |
|-------|----|----|----|----|----|----|----|
| *rank* | 3 | 5 | 6 | 7 | 1 | 2 | 4 |

Therefore, for the observed data

$$H = 21$$

Again we reject if the observed $H$ is one of the two largest or two smallest values. Based on the exact permutation distribution, we reject the null hypothesis as the p-value is $2 \times 2/35 = .101$.

**Normal approximation**

It can be shown that for large sample the null distribution of $H$ is approximately normal with mean $\mu$ and variance $\sigma^2$ where

$$\mu = \frac{m(m+n+1)}{2}, \ \sigma^2 = \frac{mn(m+n+1)}{12}$$

Suppose, as above, we compute $H = 21$ based on a samples of size 4 and 3. In this case $\mu = 16$, $\sigma^2 = 8$, so the approximate p-value is (using a continuity correction)

$$2P(H \geq 21) = 2P(H \geq 20.5) =$$

17

$$2P\left(\frac{Q-16}{\sqrt{8}} \geq \frac{20.5-16}{\sqrt{8}}\right) = 2P\left(Z \geq 1.59\right) = .11$$

which is close to the true p-value derived above even for this small sample size.

**Mann-Whitney test**

Let

$$V_{ij} = X_i - Y_j,$$

We define

$$U = \#\left(V_{ij} > 0\right)$$

which is the *Mann-Whitney* statistic. The Mann-Whitney test rejects the null hypothesis $H_0 : \delta = 0$ if $U$ is too large or too small.

For our example we see that

|    | 23 | 31 | 46 |
|----|----|----|-----|
| 37 | 14 | 6  | −9  |
| 49 | 26 | 18 | 3   |
| 55 | 32 | 24 | 9   |
| 57 | 34 | 26 | 11  |

Therefore, for this data set $U = 11$.

It can be shown that there is a relationship between the Wilcoxon rank sum $H$ and the Mann-Whitney $U$ :

$$H = U + \frac{m\left(m+1\right)}{2}.$$

Hence the critical values and p-values for $U$ can be determined from those for $H$.

**The Hodges-Lehmann confidence interval for $\delta$**

Analogous to the single sample procedure, we can construct a $(1-\alpha)100\%$ confidence interval for $\delta$ using the Mann-Whitney procedure.

We order $V_{ij}$ according to their magnitude and let $V_{[i]}$ be the $i^{th}$ largest $V_{ij}$. Then the *Hodges Lehmann estimator* for $\delta$ is the median of the $V_{ij}$.

Let

$$U(\delta) = \# (V_{ij} > \delta).$$

Then the general distribution of $U(\delta)$ is the same as the null distribution of $U$. Suppose that two-sided size $\alpha$ test the $\delta = 0$ against $\delta \neq 0$ accepts the null hypothesis if

$$a \leq U < b$$

Then a $(1-\alpha)100\%$ confidence region for $\delta$ is given by

$$a \leq U(\delta) < b \quad \Leftrightarrow \quad V_{[a]} < \delta \leq V_{[b]}$$

which is the Hodges-Lehmann confidence interval for $\delta$. In our example the estimator is the average of the 6th and 7th largest of the $V_{ij}$, giving

$$\widehat{\delta} = 16$$

The parametric estimator is $\overline{X} - \overline{Y} = 16.2$.

To find the confidence interval, note that $H = U + 10$

$$.89 = P(12 \leq H < 21) = P(2 \leq U < 11)$$

Therefore the .89 Hodges-Lehmann confidence interval for $\delta$ is

$$V_{[2]} \leq \delta < V_{[11]} \Leftrightarrow 3 \leq \delta < 32$$

The classical (t) confidence interval for the data based on t-statistics is $1.12 < \delta \leq 31.22$.

**Paired data**

Analogous to the paired t-test in parametric inference, we can propose a nonparametric test of hypothesis that the median of the population of differences between pairs of observations is zero.

Suppose we observe a sequence of i.i.d. paired observations $(X_1, Y_1), ..., (X_n, Y_n)$. Let $\mu_D$ be the median of the population of differences between the pairs. The goal is to draw inference about $\mu_D$. Let

$$D_i = X_i - Y_i$$

The distribution of $D_i$ is symmetric about $\mu_D$. Therefore, we may used the procedures discussed earlier for the one-sample model, based on the observations $D_i$.

### 5.3 $k$-Sample Procedure

Suppose we wish to test the hypothesis that the $k$ samples are drawn from the populations which all have equal location parameter. The Mann-Witney-Wilcoxon procedure discussed above can be extended to analyze the data from $k$ independent samples. The test procedure we consider is the *Kruskal-Wallis Test* which is the nonparametric analogue of the parametric one-way analysis of variance procedure.

Suppose we have $k$ independent random samples of sizes $n_i, i = 1, \cdots, k$ each, represented by $X_{ij}, j = 1, \cdots, n_i; \ i = 1, \cdots, k$. Let the underlined location parameters be denoted by $\mu_i, i = 1, \cdots, k$. The null hypothesis to test is that the $\mu_i$ are all equal against the alternative that at least one pair $\mu_i, \ \mu_{i'}$ is different.

For the Kruskal Wallis test procedure, we combine the $k$ samples and rank the obser-

vations.

Let $R_{ij}$ be the rank associated with $X_{ij}$ and let $\overline{R}_{i.}$ be the average of the ranks in the $i^{th}$ sample. If the null hypothesis is true, the distribution of ranks over different samples will be random and no sample will get a concentration of large or small ranks. Thus under the null hypothesis, the average of ranks in each sample will be close to the average of ranks for under the null hypothesis.

The Kruskal-Wallis test statistic is given by

$$KW = \frac{12}{N(N+1)} \sum n_i \left( \overline{R}_{i.} - \frac{N+1}{2} \right)^2$$

If the null hypothesis is not true, the test statistic $KW$ is expected to be large and hence we reject the null hypothesis of equal locations for large values of $KW$.

We generally use a $\chi^2$ distribution with $k-1$ degrees of freedom as an approximate sampling distribution for the statistic.

## 6. Permutation tests

The parametric test statistics can also be used to carry out the nonparametric test procedures. The parametric assumptions determine the distribution of the test statistic and hence the cut-off values under the null hypothesis. Instead, we use permutation tests to determine the cutoff points.

We give an example below.

Consider a two sample problem with 4 observations $X_1, X_2, X_3, X_4$ in the first sample from cdf $F(x)$ and 3 observations $Y_1, Y_2, Y_3$ in the second sample from cdf $G(y)$. We want to test the null hypothesis $F(x) = G(x)$ against the alternative hypothesis $F(x) \neq G(x)$

Suppose we observe 37, 49, 55, 57 in the first sample and 23, 31, 46 in the second

(Section 5.2). Suppose we want a test with size .10.

1. The parametric test for this situation is the two-sample t-test which rejects if

$$|T| = \left| \frac{\overline{X} - \overline{Y}}{S_p \sqrt{\frac{1}{4} + \frac{1}{3}}} \right| > t_5^{05} = 2.015$$

   For this data set, $T = 2.08$ so we reject (barely). The p-value for these data is .092. Note that this analysis depends on the assumptions that the data are normally distributed with equal variances.

2. We now look at rearrangements of the data observed. One possible rearrangement is 31, 37 46, 55 in the first sample and 23, 49, 57 in the second. For each rearrangement, we compute the value of the $T$. Note that there are

$$\binom{7}{4} = 35$$

   such rearrangements. Under the null hypothesis (that all 7 observations come from the same distribution) all 35 rearrangements are equally likely, each with probability 1/35. With the permutation test, we reject if the value of T for the original data is one of the 2 largest or 2 smallest. This test has $\alpha = 4/35 = .11$ The p-value for the permutation test is twice the rank of the original data divided by 35.

3. If we do this to the data above, we see that the original data gives the second largest value for $T$. (Only the rearrangement 46, 49, 55, 57 and 23, 31, 37 gives a higher $T$.) Therefore we reject the null hypothesis. The p-value is $2 \times 2/35 = .11$. Note that the only assumption necessary for these calculations to be valid is that under the null hypothesis the two distributions be the same

22

(so that each rearrangement is equally likely). That is, the assumptions are much lower for this nonparametric computation.

These permutation computations are only practical for small data sets. For the two sample model with m and n observations in the samples, there are

$$\binom{m+n}{m} = \binom{m+n}{n}$$

possible rearrangements. For example

$$\binom{20}{10} = 184,756$$

so that if we had two samples of size 10, we would need to compute $V$ for a total of 184,756 rearrangements. A recent suggestion is that we don't look at all rearrangements, but rather look a randomly chosen subset of them and estimate critical values and p-values from the sample.

What most people who use these tests would do in practice is use the t-test for large samples, where the t-test is fairly robust and use the permutation calculation in small samples where the test is much more sensitive to assumptions.

## 7. Correlation coefficients

### Pearson's r

The parametric analysis assumes that we have a set of i.i.d. two-dimensional vectors, $(X_1, Y_1), ..., (X_n, Y_n)$ which are normally distributed with correlation coefficient

$$\rho = \frac{cov\,(X_i, Y_i)}{\sqrt{var\,(X_i)\,var\,(Y_i)}}.$$

$\rho$ is estimated by the sample correlation coefficient (Pearson's r)

$$r = \frac{\sum \left( X_i - \overline{X} \right) \left( Y_i - \overline{Y} \right)}{\sqrt{\sum \left( X_i - \overline{X} \right)^2 \sum \left( Y_i - \overline{Y} \right)^2}}$$

The null hypothesis $\rho = 0$ is tested with the test statistic

$$t = \sqrt{\frac{n-2}{1-r^2}} r \sim t_{n-2}$$

under the null hypothesis.

To make this test more robust, we can use a permutation test to get nonparametric critical values and p-values. To do the rearrangements for this test, we fix the $X's$ and permute the $Y's$.

**Some Semiparametric correlation coefficients**

A semiparametric model alternative for the normal correlation model above is to assume that the $(X_1, Y_1), ..., (X_n, Y_n)$ are i.i.d. from a continuous bivariate distribution, implying no ties.

**Spearman's rank correlation**

We rank the X's and Y's separately getting ranks $R_i$ and $S_i$. The sample correlation coefficient between the $R_i$ and $S_i$ is called *Spearman's rank correlation*. Suppose, for example the we observe

| $x$ | 1 | 3 | 6 | 9 | 15 |
|-----|---|---|---|---|-----|
| $r$ | 1 | 2 | 3 | 4 | 5 |
| $y$ | 1 | 9 | 36 | 81 | 225 |
| $s$ | 1 | 2 | 3 | 4 | 5 |

Then the rank correlation $r_S$ is obviously one. Note that this happens because $Y = X^2$. Since $Y$ is not a linear function of $X$, the correlation coefficient is less than 1. In fact the correlation coefficient is .967.

We often want to test that $X$ and $Y$ are independent. We reject if $r_S$ is too large or too small. We determine the critical values and p-values from the permutation test as described above. For reasonably large sample sizes, it can be shown that under the null hypothesis

$$r_S \overset{\bullet}{\sim} N\left(0, \frac{1}{n-1}\right)$$

**Kendall's coefficient of concordance**

We say two of the vectors $(X_i, Y_i)$ and $(X_{i*}, Y_{i*})$ are concordant if

$$(X_i - Y_i)(X_{i*} - Y_{i*}] > 0$$

Kendall's $\tau$ is defined by

$$\tau = 2P\left[(X_i - Y_i)(X_{i*} - Y_{i*}) > 0\right) - 1$$

We estimate Kendall's $\tau$ by

$$r_K = 2\frac{\#\,(concordant\ pairs)}{\binom{n}{2}} - 1$$

To test $\tau = 0$, we would use $r_K$. One and two sided (exact) critical values can be determined from permutation arguments. Approximate critical value and p-values can be determined from the fact that for reasonably large $n$, the null distribution is

$$r_K \overset{\bullet}{\sim} N\left(0, \frac{4n+10}{9(n^2 - n)}\right).$$