

- We study non-parametric methods to try to avoid specific functional forms. By doing this we expect that the resulting approach becomes useful for a wider class of datasets.
- Discrete data: $PMF = f(x)$, $CDF = F(x)$; continuous data: $pdf = f(x)$, $CDF = F(x)$. For continuous $f(x) = 0$ on intervals $\mathbb{R} \setminus \text{dom } f$, $F(x) = \text{all area to the left of } x$.

- Empirical PMF**: $\hat{F}_n(x)$ used to get **Empirical CDF**: $\hat{F}_n(x)$. $\hat{F}_n(x)$ is an unbiased estimate of $F(x)$. It's saying that as we collect more samples $\hat{F}_n(x) \rightarrow F(x)$.
- Kernel density estimation**: Solves the problem of values from our continuous dataset only showing probability of x_i . To solve this we add bandwidth to our set of numbers so they can give us the probability of getting numbers exactly as x and the ~~closest~~ neighborhood of x . Neighborhood = $(x_i - h, x_i + h)$

- **Uniform Kernel**: $\frac{(x_1 + x_2 + \dots + x_n)}{n}$ where if x is in x_i , neighborhood contains $f(x)$. Then $x_i = \frac{1}{n}$. Give us the ~~mean~~ **neighborhood sample average** at $\hat{F}_n(x)$. Uniform will give a blocky piecewise constant function, as there is no decay from the x to $x_i - h$.
- **Triangular Kernel**: $K_h(x, x_i) = \frac{1}{h} - \frac{|x - x_i|}{h}$ if x is within neighborhood, else 0. Add up the product and divide by n to get $\hat{F}_n(x)$. This function will give a triangle shape graph.
- **Quadratic Kernel Function**: $K_h(x, x_i) = \frac{3}{4h} - \frac{3}{3} \left(\frac{|x - x_i|}{h} \right)^2$ if x is within the neighborhood. More smooth than previous, due to better decay.
- **Gaussian Kernel**: h now plays a role in determining how quickly the information provided by the observation decays as we move away from said observation in either direction. $K_h(x, x_i) = \frac{1}{\sqrt{2\pi}h} e^{-\frac{1}{2} \left(\frac{x - x_i}{h} \right)^2}$. The graph will now look like a normal distribution, where a large h will increase std and small h will decrease std. When $h \rightarrow 0$ the graph will show local peaks and valleys. We want to smooth enough to not show non-existent peaks.
- To find optimized bandwidth use cross-validation methods where we remove parts of the data and estimate the rest of the data to examine how well the ~~approximate~~ density fits the left out data.

- Histogram density estimate**: Create a set of bins that cover the data range. In the bin that holds x_i calculate $\# \text{ of observed data within bin} / (\# \text{ of observations} \times h)$

- Goodness of Fit**: Used to test if we can use parametric model with CDF $F_0(x)$.
- χ^2 gof (chi squared)**: $H_0: F = F_0$. $H_1: F \neq F_0$. Table constructed below with formulae.
- | Values/interval of values Ex: I_1, I_2, \dots | Number of data points within said intervals Ex: n_1, n_2, \dots
- | Probability of interval I_k divided by total χ^2 difference between data and parametric model: $dK = n_k - np_k$ |
- Required criteria: I_1, I_2, \dots, I_K don't overlap, sum of entries in second column is equal to sample size n , sum of entries in 3rd column are equal to 1, sum of entries in 4th column equals 0. Recommended criteria: np_1, np_2, \dots, np_K are all ≥ 5 . Test stat = $T = \frac{d_1^2}{np_1} + \frac{d_2^2}{np_2} + \frac{d_3^2}{np_3} + \dots + \frac{d_K^2}{np_K}$ Low values of T mean there is a low discrepancy between the data and CDF. P-value = right tail prob of value of T under χ^2_{K-1} . Limitation of χ^2 gof is that the answers can vary depending on how many cells used, too many and it will magnify discrepancy.
- Uniform, $DV(a, b)$, no parameter estimates. Normal, $N(\mu, \sigma^2)$, 2 parameter estimates $\hat{\mu} = \bar{x}$ and $\hat{\sigma}^2 = S_x^2$. Poisson, $Poi(\lambda)$, 2 parameter estimate $\hat{\lambda} = \bar{x}$. Exponential, $EXP(\lambda)$, 1 parameter est $\hat{\lambda} = 1/\bar{x}$.

- problems with the χ^2 gof: the user decides the bins which can change results of the test. You can get weird results by messing with bins. Secondly, we are ~~actually~~ not using the actual observations, we are instead ~~actually~~ categorizing them within their relative bins and computing said bins frequency.
- Lilliefors Test for Normality**: doesn't depend on the user to select bins & uses the exact values of the observations. This test states that if a dataset is from a normal distribution then the CDF $F(x)$ must be similar to the corresponding normal distributions CDF. $H_0: F = N(\mu, \sigma^2)$ CDF, $H_1: F \neq N(\mu, \sigma^2)$ we test this by looking at the maximum absolute difference between the two CDF's. Table is constructed below.

unique values smallest to largest	$F_N(\hat{\mu}, \hat{\sigma}^2)$ value	\hat{F}_1 value	Differences of $\hat{\mu}$ kind	F_{Previous}	Differences of $\hat{\sigma}^2$ kind
I_1	Given by R	I_1/n	$F_1(\hat{\mu}, \hat{\sigma}^2) - \hat{F}_1$	$\hat{\sigma}$	$F_1(\hat{\mu}, \hat{\sigma}^2) - F_1$ Previous
I_2	:	I_2/n	$F_2(\hat{\mu}, \hat{\sigma}^2) - \hat{F}_2$	I_1/n	$F_2(\hat{\mu}, \hat{\sigma}^2) - F_2$ Previous
I_K	:	I_K/n	$F_K(\hat{\mu}, \hat{\sigma}^2) - \hat{F}_K$	I_{K-1}/n	$F_K(\hat{\mu}, \hat{\sigma}^2) - F_K$ Previous

Lilliefors Test

Test stat

$T_{\text{stat}} = \text{maximum magnitude of 4th or 5th columns}$. Ignore sign of value and choose one.
A large T_{stat} indicates a strong discrepancy between the data we have and the expected Normal CDF. A small T_{stat} indicates low discrepancy. Low discrepancy is what we want to confirm.

K52

K52 Test: goal of this test is to check if two different two different populations on the same continuous random variable have identical probability distributions. To test if

$F_1 = F_2$ we first construct an Empirical PMF (~~CDF~~) we then create the table.

Empirical \Rightarrow
(CDF Derived as
 $F_{n1} \& F_{n2}$)

Pooled observations | which | Empirical CDF of | Empirical CDF of | $F_{n1} - F_{n2}$ at pooled
(smallest to largest) | dataset | dataset 1 | dataset 2 | values

T_{stat} - largest magnitude ignoring sign in the $F_{n1} - F_{n2}$ column. P-value is calculated on table

Sign Test for Median :- To check if median is smaller than a certain threshold, $H_0: M \geq M_0$, $H_1: M < M_0$, $T = \# \text{ of observations that are at or above } M_0$, if $n < 20$ then $\text{Bin}(n, 0.5)$ using R, if $n \geq 20$ then $T + 0.5 - \frac{n}{2}$.

- To check if median is larger than a certain threshold then $H_0: M \leq M_0$, $H_1: M > M_0$, $T = \# \text{ of observations that are at or above } M_0$.

If $n < 20$ then Right tail prob of $\text{Bin}(n, 0.5), n \geq 20$ $T - 0.5 - \frac{n}{2}$

- Type of prob? How to compute using Z-table? $\sqrt{\frac{n}{4}}$

- Right tail prob with a positive value, Compute directly from table. - Left tail prob with a positive value, Find its right tail prob then subtract $1 - \text{said prob}$. - Left tail prob for a negative value, Find the corresponding positive value and Compute directly from table.

Right tail probability of a negative value, printed its not negative and calculate the prob with said prob do $1 - \text{prob}$ to get correct answer.

K52
for
Medians

MWW Test :- This test compares two datasets medians, similar to K52. Firstly, to use this test we must assume $f_1 \& f_2$ are in a "location family". This basically states that the shape of f_1 and f_2 are identical, only their locations are different. ALSO the population with a smaller n is Population One OR f_1 and Population with a bigger n is f_2 .

- To test if first median is smaller than the second median, $H_0: M_1 \geq M_2$, $H_1: M_1 < M_2$

$T = \text{sum of ranks for observations for Population 2}$, If at least one sample $n \leq 10$ then Left tail probability of the value of T using table for MWW test. If both $n > 10$ then

use a "+" sign on the T . - To check if the first

median is larger than the second median, $H_0: M_1 \leq M_2$,

$H_1: M_1 > M_2$, T same as above, if at least one of the

samples $n \leq 10$ then use the MWW table, if

both samples $n > 10$ then use a "-" for the T formula.

$$T + 0.5 - \frac{n_1(n_1+n_2+1)}{2}$$

$$\sqrt{\frac{n_1 n_2 (n_1+n_2+1)}{12}}$$

- To check if the two medians are the same / coincide, $H_0: M_1 = M_2$, $H_1: M_1 \neq M_2$, $T = \text{same as above}$, If one of the samples $n \leq 10$ then find the minimum between the left & right tail probabilities then multiply by 2 and look up T in the MWW table. If both samples $n > 10$ then do the same except with the above formula. To find T_{stat} for about 5 sorted pooled data | Is it from 1st pop? | order of pooled data from 1 to n_1+n_2 | Add order # of 1st pop data. The last column sum will give the T_{stat} .

Population 2

One Sample T-Test
Unpaired T-Test
Paired T-Test

- **One Way ANOVA:** Used to check if ~~multiple~~ ^{MULTIPLE} ~~variables~~ (on the same random variable) populations have the same median for the probability distribution of said random variable.
- $H_0: M_1 = M_2 = \dots = M_K$ where M_K is median ~~of~~ of the distribution of populations.
- $H_1:$ At least one of the medians is different from the rest. First we must make an assumption being, that f_1, f_2, \dots, f_K are all shifted versions of each other, AKA all populations are members of a location family.
- **KW One Way ANOVA:** Table constructed below:

Pooled data sorted small to large | Which population it comes from | Rank; 1, 2, 3, ..., n | \Rightarrow

Calculation of dataset-specific sum of ranks, $R_1, R_2, R_3, \dots, R_K$ where $R = \text{sum of Ranks}$

Rank column for each population. $| \Rightarrow$ ~~TEST STATISTICS~~

TEST stat = $\frac{12}{n(n+1)} \cdot \left(\frac{R_1^2}{n_1} + \frac{R_2^2}{n_2} + \dots + \frac{R_K^2}{n_K} \right) - 3 \cdot (n+1)$ where large T supports H_0 , small T lacks support for H_1 .

P-value: if each of K datasets has more than 5 observations then P-value can be approximated

• **Multivariate datasets:** Stats with datasets containing more than one variable. Univariate dataset: stats pertaining to datasets on only one variable. Bi-variate dataset: two variables.

• **Nonparametric measures of association:** Dependence/Independence between two variables. AKA if one variable increases should we expect an increase/decrease/no change in the other variable.

• **Nonparametric regression:** Prediction, If we have the ~~value~~ value for one of the two variables can we predict the value of the other variable for the same observation with reasonable accuracy.

• **Measuring association:** We cannot use the correlation coefficient r , because it can only calculate linear associations which cannot tell us if $X \& Y$ are independent. Thus we have to assume a smooth parametric model.

• **Spearman's Rank Correlation Coefficient r_s :** Assume no ties, AKA repetitions in X and Y values.

Sorted pairs (based on X values from smallest to largest) | Rank from 1, 2, 3, ..., n for Y values |

| Rank from 1, 2, 3, ..., n-y for X values (already sorted) | Difference between ranks $D_{rank} = X_{rank} - Y_{rank}$ |

Accuracy check, sum of last column ~~must~~ must = 0. The final column is known as $D_1, D_2, D_3, \dots, D_n$

$$r_s = 1 - \frac{6 \cdot (D_1^2 + D_2^2 + \dots + D_n^2)}{n \cdot (n^2 - 1)}$$

(r_s can never be higher than 1 and lower than -1.)

Interpretation of r_s : ± 0.5 = independence, $\pm 0.3, 0.7$ = weak, $\pm 0.9, 1.0$ = strong, ± 1 = perfect

-hypothesis testing: To check if $X \& Y$ have = association. $H_0: P \leq 0, H_1: P > 0, T = r_s$, if $n \leq 30$

then LT probability using Spearman's rank table. $\text{Norm if } n > 30$ Then LT probability of

$T = r_s$. If $n \leq 30$ then RT prob of Spearman's rank table. If $n > 30$ the RT prob

of $T = r_s$, using Z-table. To check if $X \& Y$ have a positive association. $H_0: P = 0$

$H_1: P \neq 0, T = r_s$. If $n \leq 30$ Z = the min value between left & right tail prob. If

$n > 30$ then Z = min of RT & LT prob of $T = r_s$ using Z-table.

• **Regression:** The use of one variable to predict the other, X is the covariate, Y is the response, referred to as the regression of Y on X.

• **Simple Linear Regression: LINE**: $y = a + bx + e$ where a is the intercept, b is the slope, and e is the error in the regression equation. The median of the distribution always has $e = 0$, thus median of $Y = a + bx$. This makes the line linear. Thus, we must assume median of $Y = a + bx$ to use this test. To find b AKA the slope:

Pairs Sorted by X Values	(x_1, y_1)	(x_2, y_2)	(x_3, y_3)	...	(x_n, y_n)
(x_1, y_1)		$y_2 - y_1$ $x_2 - x_1$	$y_3 - y_1$ $x_3 - x_1$		$y_n - y_1$ $x_n - x_1$
(x_2, y_2)			$y_3 - y_2$ $x_3 - x_2$		$y_n - y_2$ $x_n - x_2$
(x_3, y_3)					$y_n - y_3$ $x_n - x_3$
\vdots					$y_n - y_n$ $x_n - x_n$
(x_n, y_n)					

* find the middle most value, if b is odd, If b is even find the 2 most middle & average.

$$a_1 = y_1 - \hat{b}x_1 \quad a_2 = y_2 - \hat{b}x_2 \quad a_3 = y_3 - \hat{b}x_3 \quad \dots \quad a_n = y_n - \hat{b}x_n$$

Similarly, we then find the median estimation of the collection, a turns to \hat{a}

→ we now are able to predict the ~~value~~ median of the distribution y_{new} .

ToH for simple linear regression KT test! Testing if or if not the linear regression

$H_0: b = 0 \quad H_1: b \neq 0, \quad T = (\# \text{of positive pairwise slopes}) - (\# \text{of negative pairwise slopes.})$

If sample size $n \leq 30$ then $Z = \min(\text{LT prob of } T, \text{RT prob of } \bar{T})$

using the KT table. If $n > 30$ first compute: $T = \sqrt{n}b$, then
 $Z = \min(\text{LT prob}, \text{RT prob})$ using standard Z score table. $\sqrt{n(n-1)(2n+5)}$

*interpretation of \hat{b} : if x increases by 1 unit then median y is estimated to increase/decrease by \hat{b} .