

Example: The table below shows 695 children under 15 years of age are cross-classified according to ethnic group and hemoglobin level. Is hemoglobin level associated (related) to ethnicity?

Ethnic Group	Hemoglobin Level (g/100 ml)			Total	Proportion
	≥ 10	9.0 - 9.9	< 9.0		
A	80	100	20	200	200/695
B	99	190	96	385	385/695
C	70	30	10	110	110/695
Total	249	320	126	695	
Proportion	249/695	320/695	126/695		

- If He level is not associated to ethnicity, then the proportion of subjects in population that fall a He group does not depend on ethnicity, i.e., it is the same for each ethnicity group, and vice versa.

Chi-Square test of Homogeneity

Often we are interested in comparing different populations with respect to a variable of interest, e.g., are the populations of carriers and non-carriers of a certain antigen *homogeneous* with respect to blood type?

Example: A sample of 150 carriers of a certain antigen and a sample of 500 non-carriers showed the following blood group distributions:

Blood Group	Carriers	Non-Carriers	Total
O	72	230	302
A	54	192	246
B	16	63	79
AB	8	15	23
Total	150	500	650

Solve this
for X²
in X
not H
not H₀
dependence

Are carriers and non-carriers similar with respect to blood group distributions?

Test of Homogeneity vs. Test of Independence

Comparing the layout of this table with the table for the test of independence, we see that the two layouts are ~~just~~ ^{just} ~~a~~ ^{not} ~~of~~ ^{but} ~~independence.~~ ^{independence.} Thus, mathematically the tests of homogeneity and independence are exactly the same. So, the same formulas apply. However, there are some key conceptual differences.

Sampling procedure:

- *Test of independence:* one overall sample is collected first and then each observation is classified by levels of the two variables. So, neither row nor column totals are fixed in advance.
 - *Test of homogeneity:* several samples are collected from several populations with each sample size fixed in advance. After collecting these pre-determined # of observations, each is classified by various levels of one variable. So, in the above example, column totals are fixed.

Number of variables:

- *Test of independence:* two variables.
- *Test of homogeneity:* one variable. The column / row representing “population” is fixed due to the sampling process.

Hypotheses:

- *Test of independence:* H_0 : two variables are *Indep.*
- *Test of homogeneity:* H_0 : The populations are *identical*.
With the one variable of interest.

R code:

```
x <- c(72, 230, 54, 192, 16, 63, 8, 15)
xmat <- matrix(x, byrow=T, ncol=2)
# > xmat
# [1,] 72 230
# [2,] 54 192
# [3,] 16 63
# [4,] 8 1
# >
# > chisq.test(xmat)
# Pearson's Chi-squared test
```

$$\begin{aligned} & \text{# } \text{value of } \chi^2 \\ & \text{# } \text{df} = 3 \\ & \text{# } \text{p-value} = 0.4927 \\ & \text{# } \text{X-squared} = 2.4052, \text{ df } = 3, \text{ p-value } = 0.4927 \\ & \text{# } > \sum_{i=1}^4 \sum_{j=1}^2 \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \end{aligned}$$

Handwritten notes:
1. χ^2 value
2. degrees of freedom
3. p-value
4. X-squared value
5. df = 3
6. p-value = 0.4927
7. $\sum_{i=1}^4 \sum_{j=1}^2 \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$

Handwritten notes:
1. χ^2 value
2. degrees of freedom
3. p-value
4. X-squared value
5. df = 3
6. p-value = 0.4927
7. $\sum_{i=1}^4 \sum_{j=1}^2 \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$

$$\text{p-value} = \underline{\text{P}}[\chi_3^2 \geq 2.4052] = \underline{\underline{\text{P}}}[\chi_3^2 \geq 2.4052]$$

$$= 1 - \underline{\text{P}}[\chi_3^2 \leq 2.4052] = 1 - \underline{\underline{\text{P}}}[\chi_3^2 \leq 2.4052]$$

↑
calculate this and
verify.

Also a nonparametric procedure because we make no assumptions about the shape of the data dist.

Nonparametric Tests

Not: If N is large, \bar{x} is $N(\mu, \sigma^2/n)$ ~~not~~ μ is not ~~no diff.~~

Issue: We would like test hypothesis on center of a distribution (one-sample problem) or compare centers of two distributions (two-sample problem). But the distributions are not normal — e.g., they are skewed or data has outliers.

Q: Why not simply use large-sample z test?

If N is large, this test is valid, but it does not make sense to look for mean as the measure of center: ~~mean as the measure of center~~.

→ Median

Nonparametric procedures:

- Typically they don't assume a specific distributional form (e.g., normal); only that the distribution is continuous. Some procedures assume that the distribution is symmetric.
- More broadly applicable than parametric procedures that assume specific distributional form.
- Use these when the distributional assumption behind a parametric procedure is clearly violated.

M = Median of X

Sign test

one-sample problem → note: no assumptions about the distribution of the sample if X is given

Data: X_1, \dots, X_n — i.i.d. sample from X .

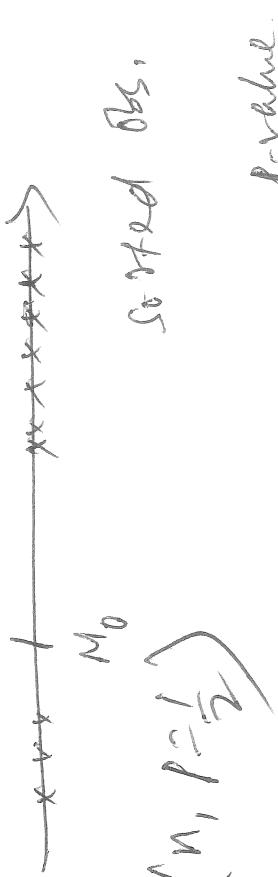
Hypotheses: $H_0: M = M_0$ vs. $H_1: M > M_0$ or $M < M_0$

SigNS: Remove the X 's that are equal to M_0 and reduce the sample size accordingly. (Let $n^* = \# \text{ obs. that are not equal to } M_0$)

Test statistic:

$$S = \# \text{ positive signs}$$

Null distribution: $S \sim \text{Binomial}(n, P = \frac{1}{2})$



So S has n^* possible values.

p-value: $P(S \geq s_{\text{obs}}) \stackrel{\text{Binomial}}{\approx} \sum_{k=s_{\text{obs}}}^{n^*} \binom{n^*}{k} \left(\frac{1}{2}\right)^{n^*}$

When to reject H_0 ?

- $H_1: M > M_0$: Reject H_0 if S is too large, which is compared to what we expect, $\frac{n^*}{2}$
 - $H_1: M < M_0$: Reject H_0 if S is too small
 - $H_1: M \neq M_0$: Reject H_0 if S is either too large or too small
- Reject H_0 if S is either too large or too small \checkmark
- [Drawing of the p-value distribution of the signed rank statistic]

Associate a sign to X :

$$X^* = \begin{cases} 1, & X > M_0 \\ 0, & X < M_0 \end{cases}$$

$$X^* \sim \text{Bernoulli}(P = \underline{P(X > M_0)})$$

$\frac{1}{2}$ if H_0 is true.

Now: Associate a sign to X_i :

$$X_i^* = \begin{cases} 1, & X_i > M_0 \\ 0, & X_i < M_0 \end{cases}$$

Random sample Bernoulli ($P = \underline{P(X > M_0)}$)

$$\Rightarrow X_{11}^* \rightarrow X_n^*$$

\Rightarrow

$$\text{Defn: } S =$$

$$\sum_{i=1}^n X_i^* = \# \text{ obs. in the sample}$$

$= \# \text{ +ve signs in sample} \sim \text{Binomial}(n, P = \underline{P(X > M_0)})$

$= \# \text{ +ve signs in sample} \sim \text{Binomial}(n, P = \frac{1}{2})$

H_0 is true.

R code:

```
# Time between keystrokes data from Example 10.9

x <- c(0.24, 0.22, 0.26, 0.34, 0.35, 0.32, 0.33, 0.29,
      0.19, 0.36,
      0.30, 0.15, 0.17, 0.28, 0.38, 0.40, 0.37, 0.27)

# Histogram and boxplot

par(mfrow=c(1,2)) # 2 plots in 1 row

hist(x)
qqnorm(x)
qqline(x)

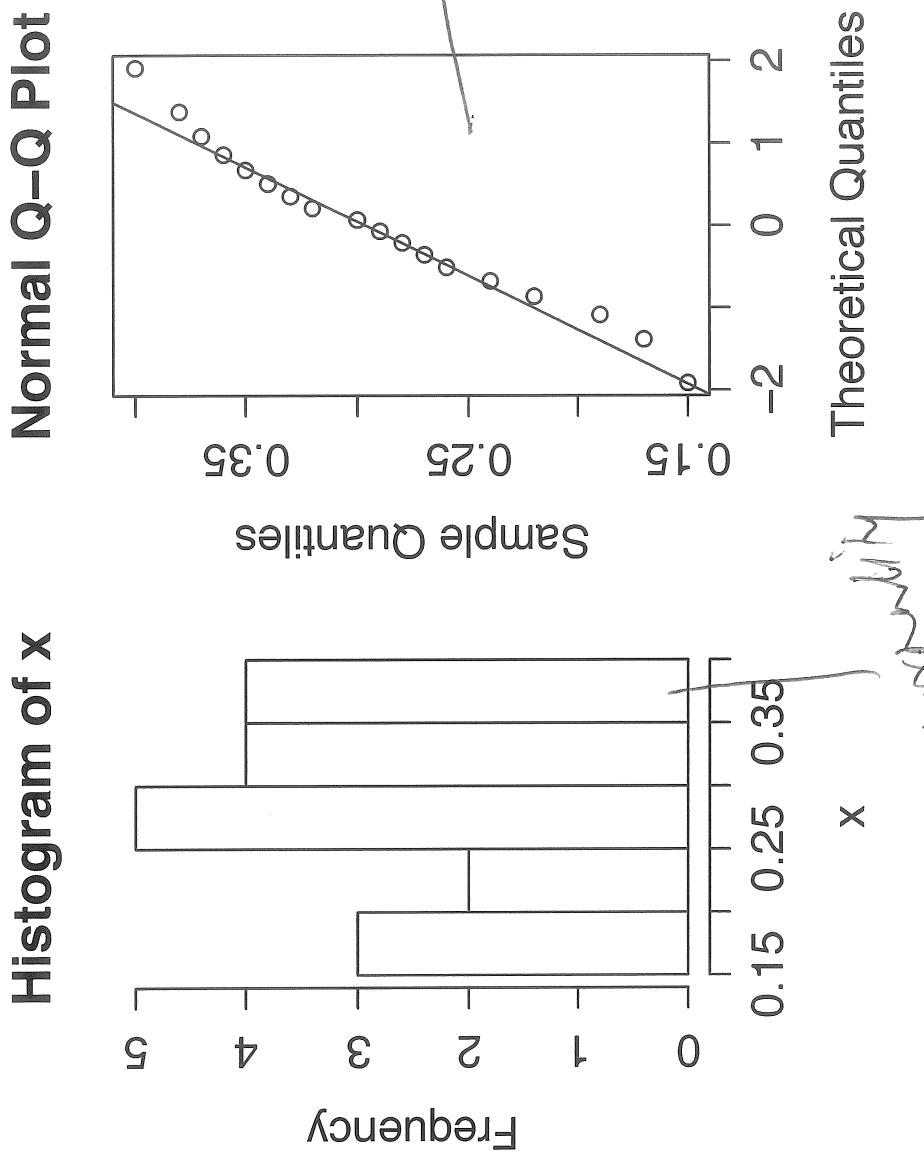
library(nortest)
```

```
> shapiro.test(x)
```

Shapiro-Wilk normality test

```
data: x  
W = 0.9611, p-value = 0.6233  
>
```

Histogram of X



The situation is not so simple.
size of n -
Let's apply the
non-parametric
procedure
see what
we conclude.

Sign test of H0: $M = 0.2$ vs M is not equal to 0.2

sign.stat <- sum(x > 0.2)

> sign.stat
[1] 15 # Obs. > 0.2
> from doing sign test
> binom.test(sign.stat, n=sum(x != 0.2), p = 0.5,
alt="two.sided", conf.level=0.95)

↓ $H_1: M \neq 0.2$

Exact binomial test

Sign test

Exact binomial test

number of successes = 15, number of trials = 18, p-value = 0.007538

alternative hypothesis: true probability of success is not equal to 0.5

Verify that it works with off.

95 percent confidence interval:

$$0.5858225 \quad 0.9642149$$

sample estimates:

probability of success

$$0.8333333$$

>

$$\text{for } p = P[X > m]$$