# **Testing Goodness of Fit**

Dr. Wolfgang Rolke University of Puerto Rico – Mayaguez Terascale Statistics School 2017 Hamburg – Germany

## **Table of Content**

- The Archetypical Statistics Problem
- Example: Is the die fair?
- Most Famous Answer: Pearson X<sup>2</sup>
- Pearson's Reasoning
- Hypothesis Testing Basics
- Another Derivation of X<sup>2</sup>
- Mendel–Fisher Controversy
- Monte Carlo Simulation
- Fisherian Significance Testing vs Neyman-Pearson

- Overfitting
- Continuous Data
- EDF Methods
- Kolmogorov–Smirnov
- ► X<sup>2</sup> vs K-S
- Probability Plots
- Smooth Tests
- Multidimensional Data
- Special Cases
- <u>2 Sample Problem</u>
- GOFer Online Goodness–of–Fit Testing

#### The Archetypical Statistics Problem:

> There is a theory

> There is data from an experiment

> Does the data agree with the theory?

### Example: Is the die fair?

Theory: die is fair ( $p_i = 1/6$ ) Experiment: roll die 1000 times If die is fair one would expect 1000\*1/6 = 1671's, 2's and so on Data:



Good fit?

#### Most Famous Answer: Pearson X<sup>2</sup>

Sir Karl Pearson 1900, "On the criterion that a given system of deviations from the probable in the case of correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling", Phil. Mag (5) 50, 157–175



Use as measure of deviations

$$X^2 = \sum \frac{(O-E)^2}{E}$$

O: observed counts

E: expected counts

Agreement is bad if  $X^2$  is large

But why 
$$\sum \frac{(O-E)^2}{E}$$
, why not (say)  $\sum \frac{(O-E)^2}{O}$  or  $\sum |O-E|$  or max{ $|O-E|$  ?

	1	2	3	4	5	6
0	187	168	161	147	176	161
E	167	167	167	167	167	167

$$X^{2} = \frac{(187 - 167)^{2}}{167} + ... + \frac{(161 - 167)^{2}}{167} = 5.72$$
  
Is 5.72 "large"?  
If die is fair and rolled 1000 times, how large  
would  $X^{2}$  typically be?

#### Pearson's Reasoning

 $N_i$  = frequency of outcome i, i = 1, ..., k $(N_1,\ldots,N_k) \sim Multinomial(n,p_1,\ldots,p_k)$  $E[N_i] = np_i, Var[N_i] = np_i(1-p_i)$  $\frac{N_i - np_i}{\sqrt{np_i(1-p_i)}} \sim_{app} N(0,1)$  by CLT  $\left(\frac{N_i - np_i}{\sqrt{np_i(1-p_i)}}\right)^2 = \frac{(N_i - np_i)^2}{np_i(1-p_i)} \sim_{app} \chi^2(1)$ so maybe  $\sum_{i=1}^{k} \frac{(N_i - np_i)^2}{np_i(1 - p_i)} \sim_{app} \chi^2$ ? but  $N_1 + ... + N_k = n$  fixed (not independent)

$$k = 2 : (N_1, N_2) = (N, n - N)$$

$$X^2 = \sum_{i=1}^k \frac{(N_i - E_i)^2}{E_i} =$$

$$\frac{(N - np)^2}{np} + \frac{(n - N - n(1 - p))^2}{n(1 - p)} =$$

$$\frac{(N - np)^2}{np} + \frac{(n - N - n + np))^2}{n(1 - p)} =$$

$$\left(\frac{1}{np} + \frac{1}{n(1 - p)}\right)(N - np)^2 =$$

$$\left(\frac{1 - p + p}{np(1 - p)}\right)(N - np)^2 =$$

$$\frac{(N - np)^2}{np(1 - p)} = \left(\frac{N - np}{\sqrt{np(1 - p)}}\right)^2 \sim \chi^2(1)$$

Pearson:  $X^2$  has a chi square distribution with k–1 degrees of freedom (k=number of categories)

Here: mean of  $\chi^2(5) = 5$ 

So our  $X^2 = 5.72$  is not unusually large, die is fair.

In the derivation of the distribution of  $X^2$  we used the CLT approximation, so this needs a sufficiently large sample size. But how large does it have to be?

```
Famous answer: E \ge 5
```

William G. Cochran The [chi-squared] test of goodness of fit. *Annals of Mathematical Statistics* 1952; 25:315-345.

Seems to have picked 5 for no particular reason. Later research showed this is quite conservative.

## Hypothesis Testing Basics

Type I error: reject true null hypothesis

Type II error: fail to reject false null hypothesis

1: A HT has to have a true type I error probability no higher than the nominal one

2: The probability of committing the type II error should be as low as possible (subject to 1)

Historically 1 was achieved either by finding an exact test or having a large enough sample.

#### Another Derivation of $X^2$

Neyman, Jerzy; Pearson, Egon S. (1933). "On the Problem of the Most Efficient Tests of Statistical Hypotheses". Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences. 231 (694–706):

In a test of a simple vs simple hypotheses likelihood ratio test is most powerful





$$\begin{split} L(p_{1},...,p_{k}) &\sim p_{1}^{N_{1}}...p_{k}^{N_{k}} \\ \Lambda &= \frac{L(p_{1},...,p_{k})}{\max\{L(p_{1},...,p_{k}):n_{1}+..+n_{k}=n\}} = \\ \frac{L(p_{1},...,p_{k})}{L(N_{1}/n,...,N_{k}/n\}} &= \frac{p_{1}^{N_{1}}...p_{k}^{N_{k}}}{\left(\frac{N_{1}}{n}\right)^{N_{1}}...\left(\frac{N_{k}}{n}\right)^{N_{k}}} = \\ \left(\frac{np_{1}}{N_{1}}\right)^{N_{1}}...\left(\frac{np_{k}}{N_{k}}\right)^{N_{k}} \end{split}$$

Samuel S. Wilks: "*The Large–Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses"*, The Annals of Mathematical Statistics, Vol. 9, No. 1 (Mar., 1938), pp. 60–62

$$-2log\Lambda \sim \chi^2(k-1)$$

 $-2\log\Lambda = -2\log\left[\left(\frac{np_1}{N_1}\right)^{N_1}...\left(\frac{np_k}{N_k}\right)^{N_k}\right] =$  $2\sum n_i \log \frac{N_i}{nn_i} =$  $2\sum n_i \log\left(\frac{N_i}{np_i} - 1 + 1\right) =$  $2\sum n_i \log\left(\frac{N_i - np_i}{np_i} + 1\right)$  $log(x + 1) \approx x + x^2/2$  Taylor expansion  $-2\log\Lambda \approx 2\sum n_i \left(\frac{N_i - np_i}{np_i} + \left(\frac{N_i - np_i}{np_i}\right)^2/2\right) =$  $2\sum_{i}n_i\frac{N_i-np_i}{nn_i} + \sum_{i}\left(\frac{N_i-np_i}{nn_i}\right)^2 \approx X^2$ because  $N_i \approx np_i$ , so  $N_i - np_i \approx 0$ so  $2\sum n_i \frac{N_i - np_i}{np_i} \approx 0$ 

#### The Degree of Freedom Controversy

Not

 $H_0: F = Normal(0,1)$  (simple hypothesis) but

 $H_0: F = Normal$ (composite hypothesis)Idea: find estimates of parameters, use those.Any change in test? Pearson said no.In 1915 Greenwood and Yule publish an analysis ofa 2x2 table and note that there is a problem.In 1922, 1924 and 1926 Sir Karl Fisher publishedseveral papers showing that Pearson was wrong.

If m parameters are estimated  $X^2 \sim \chi^2 (k - 1 - m)$ The 1922 paper is the first ever to use the term "degrees of freedom".

In some ways this is an astonishing result: it does not seem to matter how well one can estimate the parameter (aka what sample size is used)

Does it matter what method of estimation is used? Yes, and it has to be minimum chisquare!

Except these days everyone is using maximum likelihood, and then this result can be wrong

Pearson didn't acknowledge Fisher was right until 1935!



### Mendel-Fisher Controversy

Mendel, J.G. (1866). "Versuche über Pflanzenhybriden", *Verhandlungen des naturforschenden Vereines in Brünn*, Bd. IV für das Jahr, 1865, *Abhandlungen*: 3-47

Discovery of Mendelian inheritance

Immediate impact on Science: ZERO!

Darwin could have used this when he wrote On The Origin of Species. His cousin Francis Galton (inventor of regression!) could have told him.







Around 1900, <u>Hugo de Vries</u> and <u>Carl Correns</u> first independently repeat some of Mendel's experiments and then rediscover Mendel's writings and laws.

Finally Mendel becomes the "Father of Genetics"

Fisher, R.A. (1936). <u>"Has Mendel's work been</u> <u>rediscovered?</u>". Annals of Science. 1 (2): 115–137.

Fisher re-analyzed Mendel's data and applied the  $X^2$  test to all of them together. He finds an (almost) perfect agreement. But inheritance is intrinsically random, the agreement should not be that good.

Fisher's words: "to good to be true"

X<sup>2</sup> large (blue area)
→ difference between O and E to large
→ theory doesn't agree with data

X<sup>2</sup> small (red area)
→ difference between O and E to small
→ Cheating!



More than 50 papers published since 1936 have tried to figure out what happened.

For a long time: it was the Gardener!

Another explanation, which seems to have gained momentum in recent years: It was early in the history of experimentation, modern ideas of how to avoid (even unconscious) biases were not yet developed.

Allan Franklin, A. W. F. Edwards, Daniel J. Fairbanks, Daniel L. Hartl and Teddy Seidenfeld. *"Ending the Mendel–Fisher Controversy",* University of Pittsburgh Press, 2008.

#### Variations on $X^2$

Cressie-Read $\frac{1}{n\lambda(\lambda-1)}\sum O\left\{\left(\frac{O}{E}\right)^{\lambda}-1\right\}$ Pearson  $(\lambda = 1)$  $\sum \frac{(O-E)^2}{E}$ log likelihood ratio  $(\lambda = 0)$  $2\sum O\log(\frac{O}{E})$ Freeman-Tukey  $(\lambda = -1/2)$  $4\sum \left[\sqrt{O} - \sqrt{E}\right]^2$ Neyman modified  $X^2$   $(\lambda = -2)$  $\sum \frac{(O-E)^2}{O}$ modified likelihood ratio  $(\lambda = -1)$  $2\sum E\log(\frac{E}{O})$ 

Question used to be: which converges fastest to  $\chi^2$ ? But these day null distribution can be found most easily using Monte Carlo simulation!

#### **Monte Carlo Simulation**

function(B=1e4) { O<-c(187,168,161,147,176,161) E < -rep(1,6)/6\*1000TS.Data<-rep(0,5)TS.Data[1]<-sum((O-E) $^2/E$ ) TS.Data[2] < -2\*sum(O\*log(O/E))TS.Data[3] < -4\*sum( (sqrt(O)-sqrt(E))^2) TS.Data[4] < -sum( (O-E)^2/O) TS.Data[5] < -2\*sum(E\*log(E/O))TS.Sim < -matrix(0,B,5)for(i in 1:B) { O<-table(sample(1:6,size=1000,replace=T TS.Sim[i,1] < sum( (O-E)^2/E) TS.Sim[i,2] < -2\*sum(O\*log(O/E)) TS.Sim[i,3] < -4\*sum((sqrt(O)-sqrt(E))^2) TS.Sim[i,4] < sum( (O-E)^2/O) TS.Sim[i,5] < -2\*sum(E\*log(E/O))

list(TS.Data,apply(TS.Sim,2,quantile,0.95))

Method	Data	95 <sup>th</sup>
Pearson	5.72	10.95
log likelihood ratio	5.76	10.97
Freeman-Tukey	5.75	10.95
Neyman modified	5.73	11.08
modified likelihood ratio	5.73	11.00

#### Question today: Which Method has highest Power?

function(B=1e4) { crit95<-c(10.95, 10.97, 10.95, 11.08, 11.00)E < -rep(1,6)/6\*1000TS.Sim < -matrix(0,B,5)for(i in 1:B) { 0<table mple(1:6,size=1000,replace=T, TS.Sim[i,1] < sum( (O-E) $^2/E$ ) TS.Sim[i,2] < -2\*sum(O\*log(O/E))TS.Sim[i,3] < -4\*sum( (sqrt(O)-sqrt(E))^2 TS.Sim[i,4]<-sum((O-E) $^2$ /O) TS.Sim[i,5] < -2\*sum(E\*log(E/O)) power < -rep(0,5)for(i in 1:5) power[i]<sum(TS.Sim[,i]>crit95[i])/B S power

Method	Power
Pearson	55.47%
log likelihood ratio	53.95%
Freeman-Tukey	53.33%
Neyman modified	50.50%
modified likelihood ratio	52.26%

#### Fisherian Significance Testing vs Neyman-Pearson

Fisher's question: does data agree with theory?

Neyman-Pearson's question: should one reject the null hypothesis in favor of some specific alternative?

Main advantage of Neyman-Pearson style test: can decide which method is better (aka has a higher power)

Today's procedure is a hybrid of both

GOF testing much closer to Fisherian significance testing, except when we have a specific alternative in mind George Box: All models are wrong, but some are useful

Probability models are theoretical constructs, one can not expect them to be perfect fits in real life ("there is no perfect circle in nature")

→ how close an agreement between null and data is needed depends on context

→ related to choice of type I error probability  $\alpha$ , 5%? 1%? 5 $\sigma$  (I hope not!)

## Overfitting

Usual question: is our theory a good enough model for the data?

We also should worry about: is our model better than it should be?

> Overfitting!

#### Exponential Model – Good Fit?

 $\chi^2$  (6 bins): p value = 0.111

# KS test: p value = 0.117



Typical procedure (especially for background fits):

Start with low degree polynomial

Add higher degrees until fit looks good

If GOF done by "visual inspection" quite likely leads to overfitting.

Additional problem: Version of "look-elsewhere effect" (aka simultaneous inference)

### **Continuous Data**

Need to bin the data In principle any binning is ok

Two Questions:

What kind of bins?
 How many bins?

#### What kind of bins? Equi-distant - Equi-probable



Most studies suggest equi-probable is better

One advantage: E=1/k >> 5 for all bins, no need to adjust binning

Major advantage: In general leads to tests with higher power

Bins can be found easily as quantiles of F or as quantiles of data

## How many bins?

Textbook answer:  $k = 2n^{2/5}$ D'Agostini and Stephens (1986) "*Goodness-of-Fit Techiques*"

But: really depends on alternative Example:  $H_0: X \sim U[0,1] vs H_a: X \sim Linear$ Optimal k: k=2!

### **EDF** Methods

#### **EDF**: Empirical Distribution Function

$$\widehat{F}(x) = \frac{1}{n} \sum_{i=1}^{n} I_{(-\infty,x]}(X_i) = \frac{\text{\# of events } \le x}{n}$$

 $\hat{F}(x) \rightarrow F(x)$  uniformly (Glivenko–Cantelli lemma)

SO

#### $D\{\widehat{F}(x),F(x)\}$

where D is some distance measure on function space can be used for goodness-of-fit test.

**Theorem**: (Probability Integral Transform) Let X be a continuous random variable with distribution function F, then the random variable

Y = F(X) has a uniform (0,1) distribution.

Consequence: D is distribution free, aka does not depend on F.

One table to rule them all!

Except this does not work if parameters are estimated from data!

#### Kolmogorov–Smirnov

$$KS = max\{|\hat{F}(x) - F(x)|; x\} = max\{\left|\frac{i}{n} - F(X_{(i)})\right|, \left|F(X_{(i)}) - \frac{i-1}{n}\right|\}$$

Kolmogorov A (1933). "Sulla determinazione empirica di una legge di distribuzione". G. Ist. Ital. Attuari. 4: 83-91.

Smirnov N (1948). "Table for estimating the goodness of fit of empirical distributions". Annals of Mathematical Statistics. 19: 279-281



#### Alternatives

#### Anderson-Darling

Anderson, T. W.; Darling, D. A. (1952). "Asymptotic theory of certain "goodnessof-fit" criteria based on stochastic processes". Annals of Mathematical Statistics. 23: 193-212.

$$AD = n \int_{-\infty}^{\infty} \frac{(\widehat{F}(x) - F(x))^2}{F(x)[1 - F(x)]} dF(x)$$

#### Cramer-vonMises

Cramér, H. (1928). "On the Composition of Elementary Errors". Scandinavian Actuarial Journal. 1928 (1): 13–74. doi:10.1080/03461238.1928.10416862. von Mises, R. E. (1928). Wahrscheinlichkeit, Statistik und Wahrheit. Julius Springer.

And more...

Modern theory based on convergence of  $\hat{F}$  to Gaussian process

$$CM = \int_{-\infty}^{\infty} \left(\widehat{F}(x) - F(x)\right)^2 d\widehat{F}(x)$$

None of these allows estimation of parameters except in some special cases:

*H*<sub>0</sub>: *X*~*Normal* Hubert Lilliefors (1967), "*On the Kolmogorov– Smirnov test for normality with mean and variance unknown*", Journal of the American Statistical Association, Vol. 62. pp. 399–402.

*H*<sub>0</sub>: *X*~*Exponential* Hubert Lilliefors (1969), "On the Kolmogorov– Smirnov test for the exponential distribution with mean unknown", Journal of the American Statistical Association, Vol. 64. pp. 387–389.

But then again, just find null distribution via Monte Carlo!

#### R package KScorrect

Uses maximum likelihood to estimate parameters and Monte Carlo simulation to estimate null distribution

#### Example:

> x<-rexp(1000,1)
> LcKS(x,"pexp")\$p.value
[1] 0.3998

- "pnorm" for normal,
- "pmixnorm" for (univariate) normal mixture,
- "plnorm" for lognormal (log-normal, log normal),
- "punif" for uniform,
- "plunif" for loguniform (log-uniform, log uniform),
- "pexp" for exponential,
- "pgamma" for gamma,
- "pweibull" for Weibull.

# Later I will show you another way to do GOF testing!

 $X^2$  vs K-S

 $H_0: F = Uniform[0,1]$ 

 $H_a$ : F = Linear [0,1]

Sample size: n=1000

X<sup>2</sup>: 32 bins (here Equi-distant = Equi-probable )





### **Probability Plots**

Plot quantiles of F vs sample quantiles

If F is correct model, points form a straight line



## Turn this into a formal test

Again Probability Integral Transform:  $X \sim F \rightarrow F(X) \sim U[0,1]$ 

 $(U_1, ..., U_n)$  *iid* U[0,1]

Order Statistic

$$U_{(1)} < \ldots < U_{(n)}$$

 $U_{(k)} \sim Beta(k, n-k+1)$ 

Find pointwise confidence intervals from quantiles of Beta distribution

Turn into simultaneous confidence band by adjusting nominal confidence level

Sivan Aldor-Noima, Lawrence D. Brown, Andreas Buja, Robert A. Stine and Wolfgang Rolke, "*The Power to See: A New Graphical Test of Normality",* The American Statistician (2013), Vol 67/4

Andreas Buja, Wolfgang Rolke "Calibration for Simultaneity: (Re) Sampling Methods for Simultaneous Inference with Applications to Function Estimation and Functional Data", Technical Report, Wharton School of Business, Univ. of Pennsylvania

R routines: http://academic.uprm.edu/wrol ke/research/publications.htm



#### **Smooth Tests**

Old idea – goes back to Neyman (1937) – but with some recent improvements.

Basic idea: embed density f in family of densities  $\{g_k\}$  indexed by some parameter vector  $\Theta = (\theta_1, \dots, \theta_k)$  which includes true density for some k and such that

 $H_0$ : true density is  $f \leftrightarrow H_0$ :  $\Theta = \mathbf{0}$ 

$$g_k(x;\theta,\beta) = C(\theta,\beta) \exp\left\{\sum_{j=1}^k \theta_j h_j(x;\beta)\right\} f(x;\beta)$$

# ${h_j}$ should be orthonormal family of functions, i.e.

$$\int_{-\infty}^{\infty}h_i(x)h_j(x)dx=\delta_{ij}$$

#### optimal choice of $\{h_j\}$ depends on f!

Typical choices for  $\{h_j\}$ : Legendre Polynomials, Fourier series,  $h_j(\mathbf{x}) = \sqrt{2} \cos(j\pi x)$ , Haar functions, ....

Basics of the test:

$$U_{j} = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} h_{j}(X_{i})$$
$$T_{k} = \sum_{j=1}^{k} U_{j}$$
$$T_{k} \rightarrow_{d} \chi_{k}^{2}$$

Interesting feature: partial tests  $(\theta_1, ..., \theta_m) = 0$  for m < k can give insight into HOW null is wrong.

testing composite hypotheses is possible

Quite unusual: best method for estimating parameters: MoM (method of moments) Example:

 $X_{1}, \dots, X_{n} \text{ iid } N(\mu, \sigma)$   $\mu = E[X_{1}] \simeq \frac{1}{n} \sum_{i=1}^{n} X_{i}$   $\sigma^{2} = Var(X_{1}) = E[X_{1}^{2}] - E[X_{1}]^{2}$   $E[X_{1}^{2}] \simeq \frac{1}{n} \sum_{i=1}^{n} X_{i}^{2}$   $\sigma^{2} \simeq \frac{1}{n} \sum_{i=1}^{n} X_{i}^{2} - \left(\frac{1}{n} \sum_{i=1}^{n} X_{i}\right)^{2}$ 

#### And many more...

- Tests based on moments
- Tests specific for a distribution (Normal: more than 20 tests)
- A good place to start: "Comparing Distributions", Olivier Thais, Springer

### **Multidimensional Data**

 $\chi^2$  tests: Curse of Dimensionality (R. Bellman)

Example:  $H_0: (X_1, ..., X_d) \sim U[0,1]^d$ We want  $E \ge 5$  and we want 10 bins in each dimension. What n do we need?  $d=1: E = n/_{10} \cong 5 \rightarrow n \cong 50$  $d=2: E = n/_{10^2} \cong 5 \rightarrow n \cong 500$  $d=3: E = n/_{10^3} \cong 5 \rightarrow n \cong 5000$ ...

d=10:  $E = n/10^{10} \cong 5 \rightarrow n \cong 50$  billion

### **EDF Tests**

EDF needs an ordering.

What comes first in  $R^2$ : (0,1) or (1,0)?

One can impose an ordering but in  $R^d$  there are  $2^d - 1$  ways to do it!

Also, Probability Integral Transform no longer works, so KS test is no longer distribution free. But one can always use MC to find null distribution. GOF tests beyond 2 or 3 dimensions unlikely to be very useful.

At the very least will require gigantic data sets

Still a wide open problem!

## **Special Cases**

Often data has special features that need to be taken into account

Example: High Energy Physics

- 1) Data is truncated
- 2) Sample size is random
- 3) Data is binned

#### **Truncated Data**

Data in High Energy Physics is always truncated to a finite interval.

• Care needs to be taken with normalization (aka  $\int_{-\infty}^{\infty} f(x) dx = 1$ )

### Sample Size

In HEP experiments sample size is not fixed apriori but is a consequence of the run time  $n \sim Poisson(\lambda)$ 

If n is fixed:  $(N_1, ..., N_k) \sim Multinomial(n, p_1, ..., p_k)$ 

But if n is Poisson  $N_i \sim Poisson(\lambda p_i)$  and  $N_1, \dots, N_k$  independent! (Theory of Marked Poisson processes) Consequence:  $X^2 \sim \chi^2(k)$  (not k-1) Not an issue if null distribution is found via MC

#### **Binned Data**

Data in HEP is often already binned for various reasons, for example detector resolution

Still need to consider rebinning for chi square tests.

How about Kolmogorov–Smirnov?  $KS = max \left\{ \left| \frac{i}{n} - F(X_{(i)}) \right|, \left| F(X_{(i)}) - \frac{i-1}{n} \right| \right\}$ But we only know  $b_i < X_{(i)} < b_{i+1}$ 

Obvious answer: 
$$x_i = \frac{b_i + b_{i+1}}{2}$$
 midpoint

Better answer: spread out  $N_i$  points in  $(b_i, b_{i+1})$  according to F.

Can be quite slow (requires finding quantiles of F, solve many non-linear equations), in practise spreading them uniformly almost as good.

#### 2-Sample Problem

Say we have

$$X_1, \ldots, X_n \sim F$$
 and  $Y_1, \ldots, Y_m \sim G$ 

and we want to test

$$H_0: F = G vs H_0: F \neq G$$

At first this seems a very different problem, but fairly generally a method for one can be turned into a method for the other.

#### Example: Kolmogorov-Smirnov

GOF: 
$$KS = max\{|\hat{F}(x) - F(x)|; x\}$$

2-Sample: 
$$KS = max\{|\widehat{F}(x) - \widehat{G}(x)|; x\}$$

### **Permutation Tests**

```
Under H_0: F = G
so
```

$$X_1, \ldots, X_n, Y_1, \ldots, Y_m \sim F$$

samples are independent, so any reordering is equally likely.

Basic permutation test:

Find random permutation of data, split in vectors of size n and m, calculate test statistic, repeat many times. Compare to actual data.

#### Power

$$F = N(0,1) \ G = N(\mu, 1)$$
  
in example: n=50

Compare to classic 2-sample t test:

$$TS = \frac{\bar{X} - \bar{Y}}{s_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2,\alpha/2}$$



#### GOFer Online Goodness-of-Fit Testing

- No knowledge of R required
- Can handle continuous or binned data
- Can handle model specified via expression (density) or via bin probabilities
- Model expressions can be specified via R or C++
- Finds a variety of standard tests (Chisquare variants and EDF tests)
- Allows for "optimal" binning

## How to run App

- App is online at <u>https://drrolke.shinyapps.io/GOFer/</u>
- You can install all you need to run the app on your computer in less than 10 minutes. To see how go here <u>http://academic.uprm.edu/wrolke/GOFer/</u>
- Page also has detailed explanations on how to use app.

#### GOFer

#### For a detailed explanation of the app go here

Run MC!	Data is	Model is	MC runs	Sample size is
	Continuous •	Continuous	▼ 1000	fixed •
Source of Data?	Select Data Set		Left End Point	Right End Point
Use Included Examples	Uniform y=1	¥	0	1
Choose Probability Model to Test	Enter R expressi	on for density and hit Go	Go	
Enter R expression for density	1			
Select EDF Tests	Select Chi Squar	e Tests	Number of Bins (0=Default Formula)	Type of bins
Kolmogorov-Smirnov	✓ Pearson χ <sup>2</sup>		0	Foual Probability
Anderson-Darling	χ² λ, p			Equarriobability
Cramer-vonMises	$\ \ \square \ \chi^2 \ \lambda,m$			
Type I error rate for envelope test				
α = 5% <b>•</b>				
Do you want to find the Power of the Tests?				
No				

Choose Probability Model to Test	Enter R expression for density and hit Go	Go	
Enter R expression for density	exp(-x)		
Select EDF Tests	Select Chi Square Tests	Number of Bins (0=Default Formula)	Type of bins
Kolmogorov-Smirnov	📝 Pearson χ <sup>2</sup>	0	Equal Probability
Anderson-Darling	χ² λ, ρ		
Cramer-vonMises	🔲 χ² λ,m		

#### Type I error rate for envelope test



#### Do you want to find the Power of the Tests?

No



▼

•

Chi Square Tests. Number of bins: 32				
e				
1				
5				
ō				
3				
5				

Choose Probability Model to Test	Enter R expression for density and hit Go	Go		
Enter R expression for density	exp(-x)			
Select EDF Tests	Select Chi Square Tests	Number of Bins (0=Default Formula)	Type of bins	
Kolmogorov-Smirnov	√ Pearson χ <sup>2</sup>		Faual Drobability	
Anderson-Darling	χ² λ,ρ	0		
Cramer-vonMises	$\chi^2 \lambda, m$			
Type I error rate for envelope test				
α = 5% <b>•</b>				
Do you want to find the Power of the	Enter R expression for alternative density			
Tests?	1-x/1.3			

Power of the Tests				
χ <sup>2</sup> Tests				
Method	Bin Type	Power		
Pearson $\chi^2$	Equal Probability	27.9 %		
	EDF Type Tests			
Kolmogorov-Smirnov		61.3 %		
Anderson-Darling		69.4 %		
Cramer-vonMises		59 %		

Result of GOF Tests					
Chi Square Tests. Number of bins: 32					
Method	Bin Type	p-value			
Pearson $\chi^2$	Equal Probability	p = 0.561			
	EDF Type Tests				
Kolmogorov-Smirnov		p = 0.825			
Anderson-Darling		p = 0.593			
Cramer-vonMises		p = 0.966			



Yes

### File Upload

#### R

```
Density
exp(-param[1]*x)
Distribution
1 - \exp(-param[1]*x)
Alternative
1 - x/5
Estimator
xbar < -mean(x)
new<-param[1]
repeat {
     old<-new
     new < -old - (1/old - xbar - exp(-old)/(1 - vbar - exp(-old)))
exp(-old)))/(-1/old^2+exp(-old)/(1-exp(-old))))/(-1/old^2+exp(-old)))
old))^2)
     if(abs(old-new)<0.0001) break
}
return(new)
```

#### C++

```
Density
for(int i=0;i<n;++i) y[i] = exp(-
param[0]*x[i]);
Distribution
for(int i=0; i < n; +i) y[i] = 1.0-exp(-i)
param[0]*x[i]);
Alternative
for(int i=0;i<n;++i) y[i] = 1.0-x[i]/5.0;
Estimator
double pold, pnew, xbar;
xbar=0:
for(int i=0; i < n; ++i) xbar + = x[i];
xbar = xbar/n;
pnew=param[0];
while(abs(pold-pnew)>0.0001) {
    pold=pnew;
    pnew=pold-(1.0/pold-xbar-exp(-
pold)/(1.0-exp(-pold)))/(-
1.0/(pold*pold)+exp(-pold)/(1.0-exp(-
pold))/(1.0-exp(-pold)));
p[0]=pnew;
```

## Thanks!