

STATISTICS FOR DATA ANALYSIS

EMPHASIS: METHODS THAT ARE INSENSITIVE TO HOW DATA OR ERRORS ARE DISTRIBUTED.

I. FITTING

- LINEAR LEAST SQUARES (ADVANTAGES & LIMITATIONS)
- MINIMUM ABSOLUTE DEVIATION; NONLINEAR FITTING (LATER WE'LL TOUCH ON A NONPARAMETRIC ALTERNATIVE)

II. CONFIDENCE INTERVALS

- χ^2 (BRIEFLY)
- STATISTICAL BOOTSTRAP
- INSTRUMENT FOCUS TEST EXAMPLE (FITTING W/ UNKNOWN ERRORS)

Today

III. NONPARAMETRIC (DISTRIBUTION FREE) METHODS

- EXAMPLE: MEDIAN VS. MEAN
- MEASURES OF CORRELATION
- NONPARAMETRIC FITTING
- KOLMOGOROV-SMIRNOV TEST
- EXAMPLE: DISTRIBUTION OF NEARBY STARS

WE TAKE K-S TEST AS AN EXAMPLE OF FREQUENTIST HYPOTHESIS TESTING.

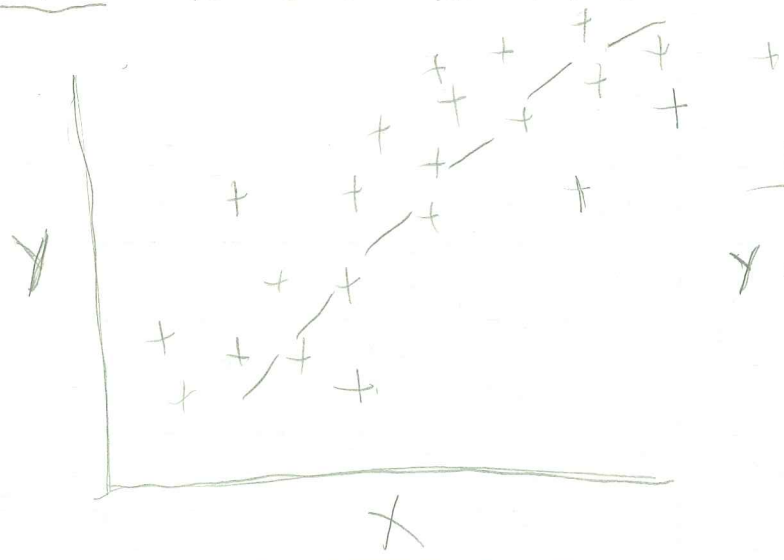
[examples/Octave/statistical_analysis/descriptive_stats_examples/gauss_example.m](https://www.mathworks.com/help/examples/Octave/statistical_analysis/descriptive_stats_examples/gauss_example.m) & [lorentzian_example.m](https://www.mathworks.com/help/examples/Octave/statistical_analysis/descriptive_stats_examples/lorentzian_example.m)

CORRELATION

A SUBSET OF STATISTICAL DEPENDENCE; APPLIES ONLY TO ORDINAL OR CONTINUOUS VARIABLES

GIVEN $\{X_1, \dots, X_N\}$ AND $\{Y_1, \dots, Y_N\}$

CORRELATION MEANS A MONOTONIC TREND IN Y VS. X



E.g.

$$Y = mX + b + \epsilon$$

↑
ADDITIVE NOISE

LINEAR CORRELATION COEFFICIENT:

$$r = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}}$$

NORMALIZATION STUFF →

r ASKS: WHEN $X_i > \bar{x}$, IS $Y_i > \bar{y}$?
 WHEN $X_i < \bar{x}$, IS $Y_i < \bar{y}$?

IF $y = mx + b$ EXACTLY,

$$r = \begin{cases} 1 & m > 0 \\ -1 & m < 0 \end{cases}$$

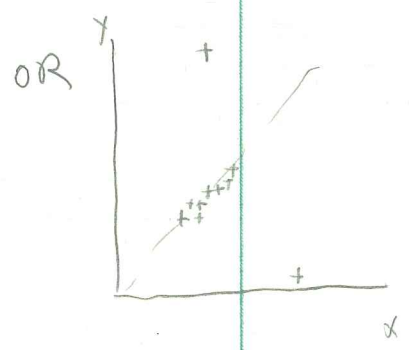
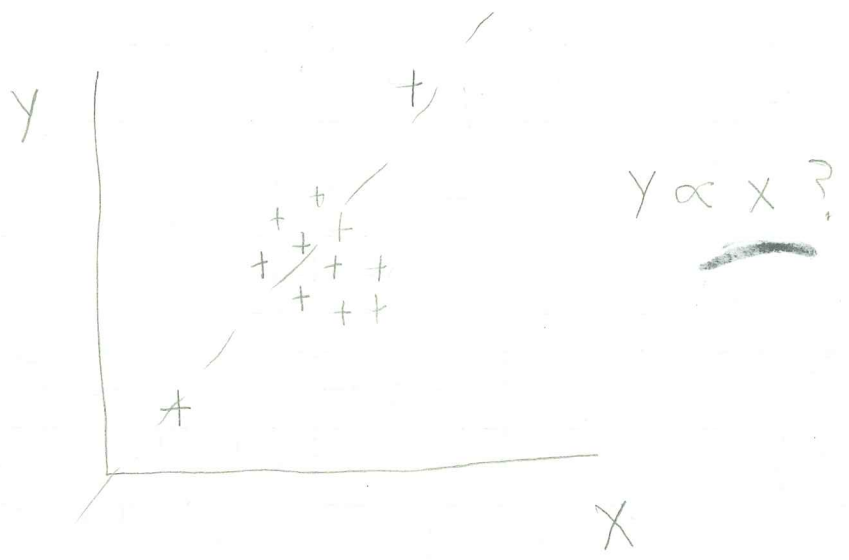
LINEAR CORRELATION COEFF. "r" IS DESIGNED FOR DATA (X, Y) WITH ADDITIVE GAUSSIAN NOISE.

PERFECT WORLD: $y' = mx' + b$

NOISED: $Y = Y' + \epsilon_1$, $P(\epsilon_1) \propto e^{-\epsilon_1^2 / 2\sigma_1^2}$

$X = X' + \epsilon_2$, --- "

BUT "r" IS EASILY BROKEN BY A FEW "OUTLIERS:"



THIS HAPPENS MOST OFTEN WITH NOISE THAT HAS FAT TAILS, E.G.

$P(\epsilon) \propto \frac{1}{1+x^2}$

r SIMPLY "WEIGHTS" OUTLIERS TOO HEAVILY.

NONPARAMETRIC ^(RANK) CORRELATION STATISTICS (NR § 14.6)

GENERAL IDEA:

IF I SORT MY DATA (X, Y) BY ASCENDING X, HOW CLOSE IS THAT TO SORTING BY ASCENDING Y?

MOTIVATION:

MUCH LIKE THE MEDIAN, ALL THAT MATTERS IS RANK ORDER (NOT MAGNITUDE). ∴ INSENSITIVE TO OUTLIERS.

SPEARMAN RANK-ORDER CORRELATION COEFF.

- ① REPLACE (X_i, Y_i) WITH THEIR RANKS (R_i, S_i)
- ② "r_s" IS THE LINEAR CORRELATION COEFFICIENT USING (R, S) IN PLACE OF (X, Y) :

$$r_s = \frac{\sum_{i=1}^N (R_i - \frac{N}{2})(S_i - \frac{N}{2})}{\sqrt{\sum_{i=1}^N (R_i - \frac{N}{2})^2} \sqrt{\sum_{i=1}^N (S_i - \frac{N}{2})^2}}$$

$-1 \leq r_s \leq 1$

(NOTE THAT $\frac{N}{2} = \bar{R} = \bar{S}$)

SPEARMAN IS FAST (LIMITED BY SORTING, $N \log N$ FOR HEAPSORT)

NR ROUTINE `spear()` GIVEN (X_i, Y_i) CALCULATES r_s AND SIGNIFICANCE.

KENDALL'S TAU (MORE NONPARAMETRIC / ROBUST / COMP. INTENSIVE THAN SPEARMAN)

- N DATA POINTS (X_i, Y_i)
- EXAMINE ALL PAIRS (X_i, Y_i) vs. (X_j, Y_j)
 (THERE ARE $\frac{1}{2}N[N-1]$ OF THEM!) $\Rightarrow \mathcal{O}(N^2)$ COMPUTATIONS.
- CLASSIFY PAIRS AND COUNT HOW MANY ARE:
 - N_c , CONCORDANT: $(X_j - X_i)$ SAME SIGN AS $(Y_j - Y_i)$
 - N_D , DISCORDANT: $(X_j - X_i)$ OPPOSITE SIGN AS $(Y_j - Y_i)$
 - $N_{=X}$, EXTRA-X: $X_i = X_j$ (BUT $Y_i \neq Y_j$)
 - $N_{=Y}$, EXTRA-Y: $Y_i = Y_j$ (BUT $X_i \neq X_j$)

$$\tau \equiv \frac{N_c - N_D}{\sqrt{N_c + N_D + N_{=Y}} \cdot \sqrt{N_c + N_D + N_{=X}}}$$

$$-1 < \tau < 1$$

UNDER NULL HYPOTHESIS OF NO ASSOCIATION X vs. Y,
 τ IS APPROX. NORMALLY DISTRIBUTED, WITH $\bar{\tau} = 0$

AND

$$\sigma_{\tau} = \sqrt{\frac{4N + 10}{9N(N-1)}} \rightarrow \sqrt{\frac{2}{3N}} \text{ FOR LARGE } N$$

NOTE: τ CAN BE SPEED UP FOR BINNED DATA
 (E.G. ORDINAL CONTINGENCY TABLE)

NR IMPLEMENTATIONS OF KENDALL'S τ :

kendr1() --- FOR 2D CONTINUOUS DATA
GIVEN (x_i, y_i)
CALCULATES τ , SIGNIFICANCE

kendr2() --- FOR 2D ORDINAL (BINNED) DATA
GIVEN A TABLE T_{ij} OF COUNTS
CALCULATES τ , SIGNIFICANCE

BOTTOM LINE: ① FOR UNBINNED DATA (x_i, y_i) ,
SPEARMAN IS FASTER. KENDALL OFFERS NO
GREAT ADVANTAGE.

② FOR BINNED DATA, τ IS QUICKER AND MORE
NATURAL. (IF # BINS IS MUCH SMALLER THAN
DATA POINTS).

A ROBUST ALTERNATIVE TO LINEAR LEAST SQUARES:

APPLICATION:

PORTER & KLIMCHUK 1995 APJ 454, 499
OR

KLIMCHUK & PORTER 1995 NATURE 377, 131

FITTING $Y'_i = mX'_i + b$

WE MEASURE $Y_i = Y'_i + \epsilon_i$, $X_i = X'_i$

GIVEN NO PRIOR KNOWLEDGE OF $P(\epsilon)$ [BUT ASSUME INDEPENDENT]
→ EXCEPT WE ASSUME MEDIAN $(\epsilon) = 0$.
(BUT YOU CAN BET $P(\epsilon)$ IS NO WUSSY GAUSSIAN!)

APPROACH

OBVIOUSLY, WE EXPECT Y_i CORRELATED WITH X_i . CONSIDER:

① $Y_i^* \equiv Y_i - m^* X_i = b + \epsilon_i + (m - m^*) X_i$

Y^* WILL BE UNCORRELATED (IN FACT INDEPENDENT) OF X

IF $m^* = m$. THUS,

② MINIMIZE: $\tau(X, Y^*)$ W.R.T. m^* .

(THE RESULT IS $m^* \approx m$.)

③ LET $b^* = \text{MEDIAN}(Y^*)$ --- ASSUMES MEDIAN $(\epsilon_i) = 0$.

ALTERNATIVELY, CONSIDER:

$$X_i = \frac{Y'_i}{m} - \frac{b}{m} = \frac{1}{m}(Y_i - \epsilon_i - b)$$

LET $X_i^* \equiv X_i - \frac{1}{m^*} Y_i = \frac{1}{m}(\epsilon_i + b) + (\frac{1}{m} - \frac{1}{m^*}) Y_i$

CAN WE MINIMIZE $\tau(X^*, Y)$ W.R.T. m^*

AND GET THE SAME RESULT AS BEFORE?

NO: FOR $m = m^*$, $X_i^* = \frac{1}{m}(\epsilon_i + b)$. X^* IS

STILL CORRELATED WITH Y BECAUSE ϵ IS CORRELATED WITH Y !
(THOUGH ϵ IS UNCORRELATED WITH Y' !)

⊗ ASSUME $P(X', Y', \epsilon) = P(X', Y)P(\epsilon)$.

NONPARAMETRIC (RANK) CORRELATION IS A

(4)

SUBSET OF "ROBUST ESTIMATION" METHODS (NR §15.7) *

VERY BROAD FIELD, EXPLODED C. 1960, STILL UNDER RESEARCH!

OBJECTIVE	A TRADITIONAL APPROACH	ROBUST METHOD(S)
CENTRAL TENDENCY OF $\{X_i\}$	MEAN	MEDIAN
COMPARE 2 DISTRIBUTIONS	MOMENTS	KOLMOGOROV-SMIRNOV STATISTIC
CORRELATION	LINEAR COEFF. "r"	NONPARAMETRIC CORRELATION E.G. KENDALL'S τ
FITTING TRENDS $Y = mx + b$	LINEAR LEAST SQUARES	METHOD OF KLIMCHUK & PORTER; MIN ABS DEV
MULTI-PARAMETER FITTING $Y(x)$	NONLINEAR LEAST SQUARES	LOESS NONPARAMETRIC REGRESSION; MIN ABS DEV
ESTIMATING CONFIDENCE LEVELS	INTEGRATE PDF OF SOME STATISTIC, ASSUMING GAUSSIAN ADDITIVE NOISE IN THE DATA	MONTÉ CARLO (E.G. BOOTSTRAP) * NR §15.6

IN GENERAL ...

<p>↑</p> <p><u>FASTER</u></p> <p>GREATER STATISTICAL POWER (IF ALL ASSUMPTIONS MET)</p>	<p>↑</p> <p><u>MORE RELIABLE</u></p> <p>(FEWER ASSUMPTIONS)</p>
---	---

* SEE ALSO EFRON & TIBSHIRANI, SCIENCE 253, 390 (1991)

KOLMOGOROV-SMIRNOV TEST

DEPENDS ON EVALUATING AND COMPARING

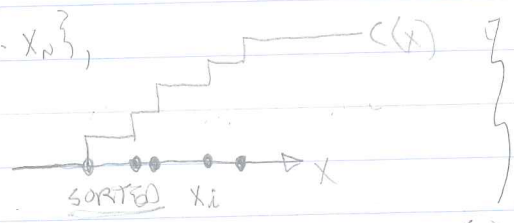
MY NO. 1 STANDARD
INSTAT FOR FOR
CUMULATIVE DIST.

• CUMULATIVE DISTRIBUTION FUNCTIONS:

$$C(x) = \int_{-\infty}^x P(x') dx'$$

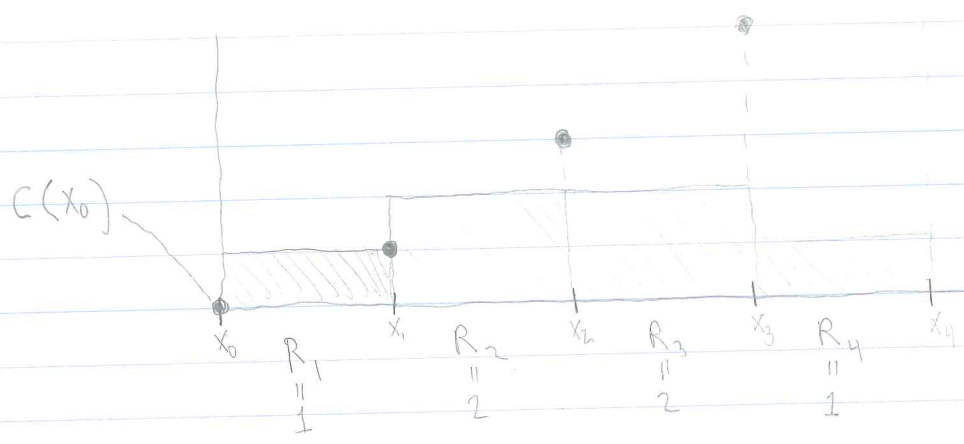
• WITH SAMPLED DATA $\{x_1, \dots, x_n\}$,

$$C(x) = \sum_{x_i \leq x} \frac{1}{N}$$



WE HAVE
C(x) AT
2N POINTS.

• WITH BINNED DATA, WE CAN ONLY EVALUATE C(x) ON THE BORDERS BETWEEN BINS:



THE K-S STATISTIC:

$$D = \max | C(x) - C(y) |$$

• INVARIANT UNDER REPARAMETERIZATION OF X:

$$X \rightarrow \log X, \quad Y \rightarrow \log Y$$

NOTES ON K-S STATISTIC D:

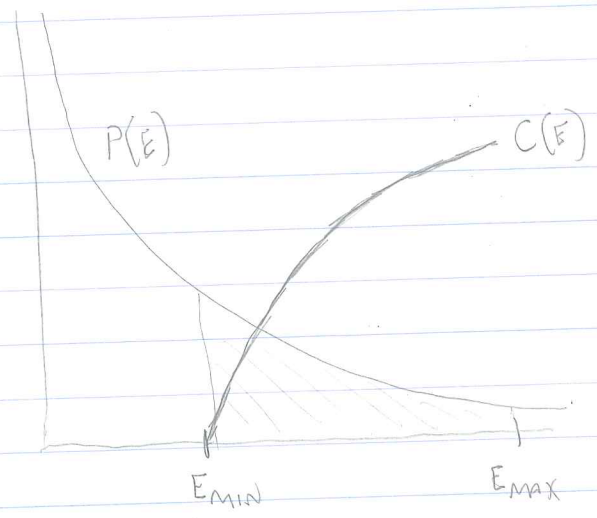
- VERY SENSITIVE TO MEDIAN \bar{x}_m .
- CONTAINS AN INTRINSIC NOTION OF FUNCTION WIDTH AND SHAPE THAT DOES NOT DEPEND ON THE EXISTENCE OF MOMENTS OF $P(x)$.

- WORKS FOR ANY $P(x)$ THAT CAN BE NORMALIZED
I.E. $\int_{-\infty}^{\infty} P(x) dx$ CONVERGES.

- CAN BE EXTENDED TO ANY $P(x) \geq 0$ ON A RESTRICTED INTERVAL $[a, b]$ AS LONG AS $\int_a^b P(x) dx$ CONVERGES.

E.g. POWER LAW DISTRIBUTION OF FLARE ENERGIES:

$$P(E) = E^{-\alpha}$$



$$C(E) = \int_{E_{min}}^E P(E') dE'$$

NULL HYPOTHESIS: $P(x) = P(y)$

K-S TEST:

$$\text{PROBABILITY } (D > \text{OBSERVED}) = Q_{KS} \left[D \left(\sqrt{N_e} + 0.12 + \frac{0.11}{\sqrt{N_e}} \right) \right]$$

WHERE: $Q_{KS}(\lambda) \equiv 2 \sum_{j=1}^{\infty} (-1)^{j-1} e^{-2j^2 \lambda^2}$ FOR $N_e \gg 1$

$(Q_{KS}(0) = 1, \quad Q_{KS}(\infty) = 0)$

★ → PLOT THIS! →

AND:

$$N_e \equiv \begin{cases} N & \text{IF COMPARING } \{X_1, \dots, X_N\} \text{ TO } P(x) \\ \frac{N_1 N_2}{N_1 + N_2} & \text{IF COMPARING } \{X_1, \dots, X_{N_1}\} \text{ TO } \{Y_1, \dots, Y_{N_2}\} \end{cases}$$

"EFFECTIVE NUMBER OF DATA POINTS"

FOR SMALL N_e , OR BINNED DATA, Q_{KS} MAY DIFFER.

POSSIBLE APPROACHES:

- DERIVE MORE GENERAL Q_{KS} ANALYTICALLY
- RUN A MONTE-CARLO TEST TO ESTIMATE STATISTICAL SIGNIFICANCE OF D