

Josua's Discussion

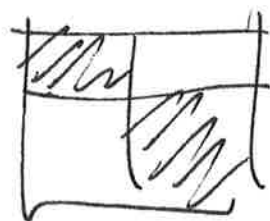
$$A = \{x_1, x_2\}$$

$$B = \{x_3, x_4\}$$

$$\tau = -|x_1 - x_2|^2 - |x_3 - x_4|^2 + 2|x_1 - x_3|^2 + 2|x_2 - x_4|^2$$

$$X = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix}$$

$$\text{we want } \tau(x, z) \stackrel{?}{=} \text{Tr}(z X X^T z)$$



$\uparrow z$

$$X^T X^T$$



$\uparrow z$

quadratic
function

$$\tau^* = \max_z \tau(x, z)$$

something
like this.

(Jan 22, 2017)

Write T in a quadratic form in terms of a matrix Z that partitions the data points into the right classes. (1)

Goal: $T^* = \max_Z \text{Tr}(Z^T D Z)$.

Thus test statistic \sim clustering.

Paper: Energy statistics, Székely, Rizzo 2013

ϵ -stats: function of distances between statistical observations.

Typically more general and powerful than classical alternatives such as Correlation, F-stats, ...

$$\text{Ex.: } \begin{cases} V_n = \frac{1}{n^2} \sum_{i,j=1}^n h(X_i, X_j) \\ h(X_i, X_j) = h(\|X_i - X_j\|) \text{ func. of Euclidean } d\text{-dimensional distance.} \end{cases}$$

Székely proposed

Rotational invariant

$$2 \int_{-\infty}^{\infty} (F(x) - G(x))^2 dx = 2 \mathbb{E} \|X - X'\| - \mathbb{E} \|X - X'\| - \mathbb{E} \|Y - Y'\|$$

where $X \sim F$
 $Y \sim G$

X' is a copy of X , $X' \stackrel{\text{iid}}{\sim} F$

Y' " Y , $Y' \stackrel{\text{iid}}{\sim} G$

It can be shown that this quantity is ≥ 0 and $= 0$ iff $X, Y \sim F$, identically distributed.

(2)

Energy distance:
$$\varepsilon(x, y) = 2\mathbb{E}|x - y| - \mathbb{E}|x - x'| - \mathbb{E}|y - y'|$$

Prop:
$$\varepsilon(x, y) = \frac{1}{c_d} \int_{\mathbb{R}^d} \frac{|\hat{f}(t) - \hat{g}(t)|^2}{1 + |t|^{\frac{d+1}{2}}} dt$$

$$c_d = \frac{\pi^{(d+1)/2}}{\Gamma(\frac{d+1}{2})}$$

$\Rightarrow \varepsilon(x, y) \geq 0$ with equality iff x, y are identically distributed.

Proof. Use characteristic function, and Parseval-Plancherel.

If x, y live in a space with metric $\delta(x, y)$ then,

$$\varepsilon(x, y) \equiv 2\mathbb{E}\delta(x, y) - \mathbb{E}\delta(x, x') - \mathbb{E}\delta(y, y')$$

\uparrow careful with the previous Prop. which may not hold.

Why special?

L_2 distance (weighted) is ε under assumption of rotation and scaling invariance.

- Look in the literature for: Maximum Mean Discrepancy (Kernel methods).
- Klebanov (2005) N -distances and their Applications.

Testing equal distributions

(3)

H_0 : X and Y have the same distr.
 H_1 : They don't.

$\{x_1, \dots, x_n\} \rightarrow$ sample from X

$\{y_1, \dots, y_m\} \rightarrow$ sample from Y

$$A = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \|x_i - y_j\|$$

$$B = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \|x_i - x_j\|$$

$$C = \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \|y_i - y_j\|$$

E-statistic energy

$$\boxed{\mathcal{E}(X, Y) \equiv 2A - B - C}$$

T-statistic test.

$$\boxed{T = \frac{n \cdot m}{n+m} \mathcal{E}(X, Y)}$$

Under H_0 , $T \xrightarrow{d}$ quadratic form independent of ?
 Normal r.v.

Under H_1 , $T \xrightarrow{d} \infty$

$$H = \frac{\mathcal{E}^2(X, Y)}{2\mathbb{E}\|X - Y\|} = \frac{2\mathbb{E}\|X - Y\| - \mathbb{E}\|X - X'\| - \mathbb{E}\|Y - Y'\|}{2\mathbb{E}\|X - Y\|}$$

$H=0$ when $X \sim Y$.

This is normality test

From the paper: $H_0: F = F_0$, $H_1: F \neq F_0$ goodness-of-fit

$$\mathcal{E}(X, F_0) = \frac{2}{n} \sum_{i=1}^n \mathbb{E}|x_i - X| - \mathbb{E}|X - X'|$$

$$- \frac{1}{n^2} \sum_{l=1}^n \sum_{m=1}^n |x_l - x_m|$$

$n \mathcal{E} \rightarrow \infty$ for H_1

$n \mathcal{E} \rightarrow$ some distr for H_0 .

$\alpha \in (0, 2)$ generalization

(4)

$$(1) \quad \varepsilon^{(1)}(X, Y) = 2\mathbb{E}|X - Y|^\alpha - \mathbb{E}|X - X'|^\alpha - \mathbb{E}|Y - Y'|^\alpha$$

$$(2) \quad \varepsilon^{(2)}(X, Y) = 2|\mathbb{E}X - \mathbb{E}Y|^2$$

$\varepsilon^{(1)}(X, Y) \geq 0$ with equality $\Leftrightarrow X, Y \sim F$.
This does not hold for (2), since it gives 0 if $\mathbb{E}X = \mathbb{E}Y$.

Conditionally negative function:

$$z_{ij} \equiv x_i - y_j. \quad \sum_{i=1}^n \sum_{j=1}^m c_i \bar{c}_j \phi(z_{ij}) \leq 0$$

$$\forall c_i, c_j \in \mathbb{R}, \text{ whenever } \sum_i c_i = 0.$$

In this case we can replace $|x - y| \rightarrow \phi(x, y)$,
it's better if ϕ is strictly negative.

Testing for equal distributions

$X = \{x_1, x_2, \dots, x_{n_1}\}$ iid random samples.

$Y = \{y_1, \dots, y_{n_2}\}$

$$\begin{aligned} \varepsilon_{n_1, n_2}(X, Y) &= \frac{2}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{m=1}^{n_2} |x_i - y_m| \quad \rightarrow 2A \\ &\quad - \frac{1}{n_1^2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_1} |x_i - x_j| \quad \rightarrow B \\ &\quad - \frac{1}{n_2^2} \sum_{\ell=1}^{n_2} \sum_{m=1}^{n_2} |y_\ell - y_m| \quad \rightarrow C \end{aligned} \quad = 2A - B - C$$

$$T_{n_1, n_2} = \frac{n_1 n_2}{n_1 + n_2} E_{n_1, n_2}$$

(5)

$H_0: X, Y \sim F$, small T

$H_1: X, Y$ diff. distn, large T

Distance Components (DISCO)

$H_0: F_1 = F_2 = \dots = F_k \quad k \geq 2.$

two samples $A = \{a_1, \dots, a_{n_1}\}$
 $B = \{b_1, \dots, b_{n_2}\}$

let $g_\alpha(A, B) \equiv \frac{1}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{m=1}^{n_2} |a_i - b_m|^\alpha \quad 0 \leq \alpha \leq 2.$

Let A_1, A_2, \dots, A_k be samples of sizes n_1, n_2, \dots, n_k
 and $\sum_{j=1}^k n_j = N.$

→ total dispersion

$$T_\alpha = T_\alpha(A_1, \dots, A_k) = \frac{N}{2} g_\alpha(A, A)$$

A is the pooled sample of size N . I guess this means putting all the data from every sample i into A .

The within-sample dispersion is

$$W_\alpha = W_\alpha(A_1, \dots, A_k) = \sum_{j=1}^k \frac{n_j}{2} g_\alpha(A_j, A_j)$$

Between-sample energy statistics:

(6)

$$S_{m,\alpha} = \sum_{1 \leq j < k \leq K} \left(\frac{n_j + n_k}{2N} \right) \left[\frac{n_j n_k}{n_j + n_k} E_{n_j, n_k}^{(\alpha)}(A_j, A_k) \right]$$

$$\uparrow \sum_{j=1}^K \sum_{\substack{k=1 \\ j < k}}^K$$

$$S_{m,\alpha} = \sum_{j=1}^K \sum_{\substack{k=1 \\ j < k}}^K \frac{n_j n_k}{2N} \left(2g_\alpha(A_j, A_k) - g_\alpha(A_j, A_j) - g_\alpha(A_k, A_k) \right)$$

For $0 < \alpha \leq 2$ we have $T_\alpha = S_\alpha + V_\alpha$

If $0 < \alpha < 2$ we have a statistically consistent test of equality of distr.

If $\alpha = 2$, ~~the test~~ E can be zero if the means of the distributions are identical.

Paper: Equivalence of Distance-Based and RKHS-Based Statistics in Hypothesis Testing, Ann. Stat (2013)

energy distance
dist. covariance
(statistics)

maximum mean discrepancy (MMD)

this is a distance between embeddings of distr. to reproducing kernel Hilbert spaces (RKHS). (Mach. Learning)

To any positive definite kernel, MMD is "kinda" equivalent to Energy distance.

(*) They show that E-dist. most commonly used in stats. is just one member of a parametric family of kernels, and other choices can yield more powerful tests. (7)

E-stats

(Statistics)

two-sample testing in Euclidean space

Szekely, Rizzo 2004, 2005
Borngens, Franz 2004

Dependence measure: dist covariance



Rkhs

(Mach. Learning)

Embeddings of prob. distr. into reproducing kernel Hilbert spaces.
Test stat: ~~dist. between~~ difference between embeddings.
Maximum Discrepancy (MMD).
Gretton et al. (2007, 2012)

Seems to be more general and recover E-stats in some limit?

two-sample testing: E-dist is a maximum mean discrepancy.

E-dist arise from a particular choice of kernel!

Semi-Metric of Negative Type:

Z is a set. $\rho: Z \times Z \rightarrow [0, \infty)$ st. $\rho(z, z') = 0 \Leftrightarrow z = z'$, and $\rho(z, z') = \rho(z', z)$. (Z, ρ) is a semi-metric space.
Negative type: $\sum_{i=1}^n \alpha_i = 0$, $\sum_{i,j=1}^n \alpha_i \alpha_j \rho(z_i, z_j) \leq 0$

Prop. If ρ is of negative type, so is ρ^q for $q \in (0, 1)$. ρ is a semimetric of neg. type $\Leftrightarrow \exists$ a Hilbert space \mathcal{H} and $\varphi: \mathcal{Z} \rightarrow \mathcal{H}$ (i-jjective) s.t. $\rho(z, z') = \|\varphi(z) - \varphi(z')\|_{\mathcal{H}}^2$.
 So $\rho^{1/2}$ is a metric, even though ρ is a semimetric. (8)

Energy Distance: measure of stat. distance between two prob. measures P and Q on \mathbb{R}^d :

$$D_E(P, Q) = 2 \mathbb{E} \|Z - W\|_2 - \mathbb{E} \|Z - Z'\|_2 - \mathbb{E} \|W - W'\|_2$$

$$\begin{array}{l} Z, Z' \overset{iid}{\sim} P \\ W, W' \overset{iid}{\sim} Q \end{array} \quad \begin{array}{l} (i) D_E(P, Q) \geq 0 \\ (ii) D_E(P, Q) > 0 \text{ if } P \neq Q. \end{array}$$

This can be generalized to

$$D_{E, \rho}(P, Q) = 2 \mathbb{E} \rho(Z, W) - \mathbb{E} \rho(Z, Z') - \mathbb{E} \rho(W, W')$$

ρ must satisfy certain conditions so $D_{E, \rho}$ is a metric.

Kernel-Based Approach:

RKHS. \mathcal{H} is a Hilbert space of real valued functions defined on a set \mathcal{Z} . $K: \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ is called a reproducing kernel of \mathcal{H} if:

$$1. K(\cdot, z) \in \mathcal{H}, \forall z \in \mathcal{Z}.$$

$$2. \langle f, K(\cdot, z) \rangle_{\mathcal{H}} = f(z), \forall z \in \mathcal{Z}, \forall f \in \mathcal{H}.$$

If \mathcal{H} has such a K then \mathcal{H} is a reproducing kernel Hilbert space.

Theo. (Moore-Aronszajn) To every symmetric, positive definite function $k: \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$, there is an RKHS \mathcal{H}_k of real valued functions. (9)
 $\varphi: \mathcal{Z} \rightarrow \mathcal{H}_k$, $\varphi(z) \mapsto k(\cdot, z)$ is the canonical feature map.

Def.: Let k be a kernel on \mathcal{Z} , and ν a prob measure on \mathcal{Z} . The kernel embedding of ν into RKHS is $\mu_k(\nu) \in \mathcal{H}_k$ s.t.
 $\int f(z) d\nu(z) = \langle f, \mu_k(\nu) \rangle_{\mathcal{H}_k} \quad \forall f \in \mathcal{H}_k.$

Def.: (MMD) Let k be a kernel on \mathcal{Z} . P, Q prob. measures on \mathcal{Z} . The MMD γ_k between P and Q is $\gamma_k(P, Q) = \|\mu_k(P) - \mu_k(Q)\|_{\mathcal{H}_k}$
 Usefull formula:

$$\gamma_k^2(P, Q) = \mathbb{E} k(z, z') + \mathbb{E} k(w, w') - 2\mathbb{E} k(z, w)$$

$z, z' \stackrel{\text{iid}}{\sim} P$
 $w, w' \stackrel{\text{iid}}{\sim} Q.$

if the restriction of μ_k to some prob space is well defined and injective γ_k is a metric.

Lemma Let $\rho: \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ be a semimetric. Let $z_0 \in \mathcal{Z}$, and denote $k(z, z') = \rho(z, z_0) + \rho(z', z_0) - \rho(z, z')$

Then k is positive definit iff ρ is conditionally negative, or a semimetric of negative type.

Distance-induced kernel:

(10)

$$k(z, z') = \frac{1}{2} \{ \rho(z, z_0) + \rho(z', z_0) - \rho(z, z') \}$$

also just called distance kernel.

Prop. 1. $\rho(z, z') = k(z, z) + k(z', z') - 2k(z, z')$
 $= \|k(\cdot, z) - k(\cdot, z')\|_{\mathcal{H}_k}^2$

2. $z \mapsto k(\cdot, z)$ is injective

Ex.: $z \subseteq \mathbb{R}^d$, $\rho_q(z, z') = \|z - z'\|^q$ $0 < q \leq 2$.

$$k_q(z, z') = \frac{1}{2} \{ \|z - z_0\|^q + \|z' - z_0\|^q - \|z - z'\|^q \}$$

$$\text{if } z_0 = 0 \Rightarrow k_q(z, z') = \frac{1}{2} \{ \|z\|^q + \|z'\|^q - \|z - z'\|^q \}$$

} let k be a nondegenerate kernel on Z . The
 $\rho(z, z') = k(z, z') + k(z', z) - 2k(z, z')$
defines a valid semimetric ρ of negative type.
 k generates ρ .

if k and k' generates the same ρ , k and k' are equivalent.

Prop. k and \tilde{k} are equivalent iff

$$\tilde{k}(z, z') = k(z, z') + f(z) + f(z')$$

for some shift function $f: Z \rightarrow \mathbb{R}$.

Equivalence of MMD and Energy distance

(11)

For every ρ , $D_{E,\rho}$ is related to the MMD associated to a kernel k that generates ρ .

Theo. (\mathcal{Z}, ρ) is a semimetric space of negative type and k is a kernel that generates ρ . Then

$$D_{E,\rho}(P, Q) = 2 \gamma_k^2(P, Q)$$

Recall $\gamma_k^2(P, Q) = \mathbb{E} k(z, z') + \mathbb{E} k(w, w') - 2 \mathbb{E} k(z, w)$

$$\mathcal{Z} = \{z_i\}_{i=1}^m \stackrel{\text{iid}}{\sim} P$$

$$\mathcal{W} = \{w_i\}_{i=1}^n \stackrel{\text{iid}}{\sim} Q \quad \bullet \text{ The empirical estimate is}$$

$$\begin{aligned} \hat{\gamma}_k(z, w) &= \gamma_k^2\left(\frac{1}{m} \sum_{i=1}^m \delta_{z_i}, \frac{1}{n} \sum_{j=1}^n \delta_{w_j}\right) \\ &= \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m k(z_i, z_j) + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(w_i, w_j) \\ &\quad - \frac{2}{m \cdot n} \sum_{i=1}^m \sum_{j=1}^n k(z_i, w_j) \end{aligned}$$

Theo. $\mathcal{Z} = \{z_i\}_{i=1}^m, \mathcal{W} = \{w_i\}_{i=1}^n$ are two iid samples from P . Assume $S_{\tilde{H}_P}$ is trace class

$$\frac{m}{2} \hat{\gamma}_k^2(\mathcal{Z}, \mathcal{W}) \xrightarrow{m} \sum_{i=1}^{\infty} d_i N_i^2$$

where $N_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$ and d_i are the eigenvalues of $S_{\tilde{H}_P}$.

Two-sample Experiment

(12)

- (1) Two Multivariate Gaussians, where the means differ in 1D only and all variances are equal.
- (2) Two Gaussians but with identical means, but variance that differ in a single dimension.
- (3) One distr. is a Gaussian in 1D, and the other is a Gaussian in 1D with a sinusoidal perturbation of increasing frequency.

Exponent $q < 2$ seems to give good results in (3).

Conclusion: Energy distance is a particular case of a larger class of discrepancy measures.

kernel based methods can be applied to data that do not lie in E^n .

A New test for Multivariate Normality
Székely, Rizzo (2005)

New class of rotation invariant and consistent goodness-of-fit tests for multivariate distributions based on Euclidean distance between sample elements.

The most widely applied tests of normality are based on Mardia's multivariate generalization of skewness and kurtosis.

(13)

Let X_1, \dots, X_n be a d -dimensional sample.

T_n is affine invariant if $T_n(A(X_1), \dots, A(X_n)) = T_n(X_1, \dots, X_n)$

for every affine transformation $A: \mathbb{R}^d \rightarrow \mathbb{R}^d$.

A goodness-of-fit test of $H_0: F \in \mathcal{F}$ versus $H_1: F \notin \mathcal{F}$ is consistent against all fixed alternatives if the probability of rejecting the null hypothesis $\rightarrow 1$ when $n \rightarrow \infty$ and the actual distribution of the sampled population is not in \mathcal{F} .

V-statistic with kernel:

$$h(x, y) = \mathbb{E} \|x - Y\| + \mathbb{E} \|y - Y'\| - \mathbb{E} \|Y - Y'\| - \|x - y\|$$

$$\|x\| = (x^T x)^{1/2} \text{ Euclidean}$$

X, X' iid F , Y is a d -dimensional random vector.

$$E_n = \frac{1}{n} \sum_{j,k=1}^n h(X_j, X_k). \quad \frac{E_n}{n} \text{ is a } V\text{-statistic.} \quad \text{degenerate}$$

If X, X', Y, Y' are independent then

$$\mathbb{E} h(X, Y) = 2 \mathbb{E} \|X - Y\| - \mathbb{E} \|X - X'\| - \mathbb{E} \|Y - Y'\| \geq 0$$

with equality iff $X, Y \sim F$.

S is a non-empty set. $\gamma: S \times S \rightarrow \mathbb{R}$ is negative definite (14)

definite if

$$\sum_{j,k=1}^n \gamma(x_j, x_k) r_j r_k \leq 0$$

where $\sum_{j=1}^n r_j = 0$. Continuous analogue is

$$\iint_S dQ(x) dQ(y) \gamma(x, y) r(x) r(y) \leq 0.$$

where $\int dQ(x) r(x) = 0$ for some measure Q .

Prop. If $\gamma(x, y) = \|x - y\|$ is the Euclidean dist. then it is strictly negative.

Theo. $2\mathbb{E} \gamma(X, Y) - \mathbb{E} \gamma(X, X') - \mathbb{E} \gamma(Y, Y') \geq 0$

iff γ is negative definite. If γ is strictly negative the equality holds iff $X \stackrel{d}{=} Y$.

Proof. Let $X \sim \mu$, $Y \sim \nu$. Let Q be an arbitrary prob measure. Define $r(x) = \frac{d\mu(x)}{dQ} - \frac{d\nu(x)}{dQ}$

We have

$$\begin{aligned} & \int \gamma(x, y) d\mu(x) d\nu(y) + \int \gamma(x, y) d\mu(x) d\nu(y) = - \frac{\int \gamma(x, y) r(x) r(y) dQ(x) dQ(y)}{\geq 0!} \\ & - \int \gamma(x, y) d\mu(x) d\mu(y) - \int \gamma(x, y) d\nu(x) d\nu(y) \end{aligned}$$

$$\left(\frac{d\mu(x)}{dQ} d\mu(y) - \frac{d\nu(x)}{dQ} d\nu(y) \right) = r(x) \cdot r(y)$$

$$\begin{aligned} r(x) \cdot r(y) &= \left(\frac{d\mu(x)}{dQ} - \frac{d\nu(x)}{dQ} \right) \left(\frac{d\mu(y)}{dQ} - \frac{d\nu(y)}{dQ} \right) \\ &= \frac{d\mu(x) d\mu(y)}{dQ dQ} - \frac{d\mu(x) d\nu(y)}{dQ dQ} - \frac{d\nu(x) d\mu(y)}{dQ dQ} + \frac{d\nu(x) d\nu(y)}{dQ dQ} \end{aligned}$$

we can use $\rho = \|\cdot\|_d$ in \mathbb{R}^d the equality (15) implies $X \stackrel{d}{=} Y$.

Suppose X_1, \dots, X_n is a random sample from F , and x_1, \dots, x_n are the observed values of the random sample. $H_0: F = F_0$ vs $H_1: F \neq F_0$.

$$E_n = n \left(\frac{2}{n} \sum_{j=1}^n \mathbb{E} \|X_j - X\| - \mathbb{E} \|X - X'\| - \frac{1}{n^2} \sum_{j,k=1}^n \|X_j - X_k\| \right)$$

$X, X' \stackrel{iid}{\sim} F_0$.

If $F_0 = N_d(\mu, \Sigma)$ then $y_j \equiv \Sigma^{-1/2}(x_j - \mu)$ and use

$$E_n = n \left(\frac{2}{n} \sum_{j=1}^n \mathbb{E} \|y_j - z\| - \mathbb{E} \|z - z'\| - \frac{1}{n^2} \sum_{j,k=1}^n \|y_j - y_k\| \right)$$

where $z, z' \stackrel{iid}{\sim} N_d(0, I)$.

It's possible to find formulas for $\mathbb{E}[\dots]$.

Then one estimates

$$\hat{\mu} = \bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

$$y_j = \hat{\Sigma}^{-1/2}(x_j - \hat{\mu})$$

The null hypothesis is rejected if E_n is large.

Projection Pursuit (PP): $X \in \mathbb{R}^d \sim N_d(\mu, \Sigma)$
 iff $a^T X \sim N_1(a^T \mu, a^T \Sigma a)$ for all $a \in \mathbb{R}^d$.
 The PP tests the "worst" d D projectible with
 a goodness-of-fit method.

(16)

$$C_n^* = \sup_{\substack{a \in \mathbb{R}^d \\ \|a\|=1}} C_n(a^T X_1, \dots, a^T X_n)$$

reject the H_0 if C_n^* is large.

C_n^* can be approximated by finding a
 finite set of projections determined by
 $\{a_1, a_2, \dots, a_m\}$ uniformly scattered in \mathbb{R}^d .

Testing For equal Distributions in High Dimensions

Szekely, Rizzo 2004

Nonparametric test for equality of two or
 more multivariate distr. based on Euclidean
 distance. Distr. are not specified. Test is
consistent. Prob. more sensitive than nearest
 neighbors and performs well in high dim.
 Computational complexity indep. of dimension
 and number of data points!

Kolmogorov-Smirnov

Cramer-von Mises

+ D

→ do not generalize
 to $D > 1$ well is
 a "distribution free"
 manner.

Suppose X_1, \dots, X_{n_1} and Y_1, \dots, Y_{n_2} are independent random samples of \mathbb{R}^d , according to F_1 or F_2 . (17)

two-sample prob: $H_0: F_1 = F_2$
 $H_1: F_1 \neq F_2$

n-samples prob: $H_0: F_1 = \dots = F_k$
 $H_1: F_i \neq F_j$ for some $i \neq j$

Empirically: ε -test way superior than nearest neighbor in high dimensions.

$A = \{a_1, \dots, a_{n_1}\}$ disjoint non-empty sets of \mathbb{R}^d
 $B = \{b_1, \dots, b_{n_2}\}$

$$\varepsilon(A, B) \equiv \frac{n_1 n_2}{n_1 + n_2} \left\{ \frac{2}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \|a_i - b_j\| \right. \\
\left. - \frac{1}{n_1^2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_1} \|a_i - a_j\| \right. \\
\left. - \frac{1}{n_2^2} \sum_{i=1}^{n_2} \sum_{j=1}^{n_2} \|b_i - b_j\| \right\} \geq 0$$

- See Székely, Rizzo (2003), Hierarchical clustering via joint between-within distances.

Two-sample probs

$$A = \{X_1, \dots, X_{n_1}\} \quad \mu_{AB} = E\|X - Y\|$$

$$B = \{Y_1, \dots, Y_{n_2}\} \quad \mu_A = E\|X_1 - X_2\|$$

$$\mu_B = E\|Y_1 - Y_2\|$$

$$E[\varepsilon(A, B)] = \frac{n_1 n_2}{n_1 + n_2} \left\{ \frac{2}{n_1 n_2} \frac{n_1 n_2}{\mu_{AB}} - \frac{1}{n_1^2} \frac{n_1(n_1-1)}{\mu_A} \right. \\
\left. - \frac{1}{n_2^2} \frac{n_2(n_2-1)}{\mu_B} \right\}$$

(48)

$$E[E(A,B)] = \frac{n_1 n_2}{n_1 + n_2} \left\{ 2\mu_{AB} - \frac{n_1 - 1}{n_1} \mu_A - \frac{n_2 - 1}{n_2} \mu_B \right\}$$

$$= \frac{n_1 n_2}{n_1 + n_2} \left\{ 2\mu_{AB} - \mu_A - \mu_B \right\} + \frac{n_2}{n_1 + n_2} \mu_A + \frac{n_1}{n_1 + n_2} \mu_B$$

if $X \stackrel{D}{=} Y \Rightarrow \mu_{AB} = \mu_A = \mu_B$, $\stackrel{!}{=} 0 \Rightarrow E[E(A,B)] = \mu_{AB}$
 if $X \not\stackrel{D}{=} Y \Rightarrow \neq 0 \Rightarrow E[E(A,B)]$

$$E[E(A,B)] = \frac{n_1 n_2}{n} C + \underbrace{\left(\frac{n_2}{n} \mu_A \right)}_{C_2} + \underbrace{\left(\frac{n_1}{n} \mu_B \right)}_{C_1}$$

$$= C_1 C_2 n C$$

Fixed ratios

so $E[E(A,B)]$ is $\sim n$ (asymptotically)

as $n \rightarrow \infty$, under H_0 , $E[E(A,B)] \rightarrow \text{cte}$

under H_1 , $E[E(A,B)] \rightarrow \infty$

Not only $E[E]$ but E itself either converges or diverges. (in distr.)

$$E(A,B) = \frac{n_1 n_2}{n} C \quad \text{conditional on data}$$

$$= \frac{n_1}{n} \left(\frac{n_2}{n} \cdot n \cdot C \right)$$

fixed ratios (asymptotically).

$$C = 2\mu_{AB} - \mu_A - \mu_B = 0 \Leftrightarrow X \stackrel{D}{=} Y$$

Two-sample test stats.

$$E_{n_1, n_2} = \frac{n_1 n_2}{n_1 + n_2} \left(2 \sum_{i=1}^{n_1} \sum_{m=1}^{n_2} \|X_i - Y_m\| - \frac{1}{n_1^2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_1} \|X_i - X_j\| - \frac{1}{n_2^2} \sum_{l=1}^{n_2} \sum_{m=1}^{n_2} \|Y_l - Y_m\| \right)$$

(19)

under H_0 or H_1 , a random permutation of

$$W_1^{(\pi)}, \dots, W_n^{(\pi)} \quad n_1 + n_2 = n$$

of $X_1, \dots, X_{n_1}, Y_1, \dots, Y_{n_2}$, is equal in distribution to a random sample of size n from the mixture W where W samples from X with prob. $\frac{n_1}{n}$ and from Y with prob. $\frac{n_2}{n}$.

$$\lim_{n \rightarrow \infty} P(E_n > C_\alpha) = \alpha \in (0, 1)$$

Reject H_0 if $E_n > C_\alpha$

Implementation

- Permutation test Approach, see Efron 1993, chap 15
An introduction to the bootstrap.

$$\begin{matrix} A_1, A_2, \dots, A_K \\ \sum \\ F_1, F_2, \dots, F_K \end{matrix} \xrightarrow[n = \sum n_i]{\text{pooled sample}} \{W_1, \dots, W_n\} = A_1 \cup \dots \cup A_K$$

under H_0 , W_1, \dots, W_n are iid, with distr. F .
If the desired significance level is α ,

resample (without replacement) from W_1, \dots, W_n
 B samples of size n so that $(B+1)\alpha$ is
 an integer.

$$m_j = \sum_{i=1}^j n_i \quad m_0 = 0.$$

$\{W_1^{(b)}, \dots, W_n^{(b)}\} \rightarrow$ one bootstrap sample.
 $b = 1 \dots B$
 Compute $E_n^{(b)}$ from $A_i^{(b)} = \{W_{m_{i-1}+1}^{(b)}, \dots, W_{m_i}^{(b)}\}$
 $i = 1 \dots k$

The bootstrap estimate of $P_n(\cdot E_n \leq t)$ is

$$\frac{1}{B} \sum_{b=1}^B \mathbb{I}(E_n^{(b)} \leq t)$$

Reject H_0 if the observed E_n exceeds
 $100(1-\alpha)\%$ of the replicates $E_n^{(b)}$