Two Sample Testing – Free Lunches & Comp-Stat Tradeoffs

Introduction

Thought Experiment: Red vs Blue Pill



4-week Experiment: Each day, give Neo either R or B pill, measure real-valued thyroxine concentration (X or Y respectively).Data: R: $X_1, ..., X_{14} \sim P$ (iid) and B: $Y_1, ..., Y_{14} \sim Q$ (iid). Mean-difference alternatives:

 $H_0: \mathbb{E}_P[X] = \mathbb{E}_Q[Y]$ vs. $H_1: \mathbb{E}_P[X] \neq \mathbb{E}_Q[Y]$

General alternatives:

 $H_0: P = Q$ vs. $H_1: P \neq Q$

Two-sample test: inputs data, outputs 0 or 1. "Stochastic proof by contradiction". Nonparametric: no assumptions about P, Q.

Real Experiment: Faces vs Houses

Question: Does brain region R differentiate faces and non-faces?

• Show someone a face:

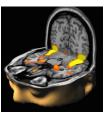


• Measure brain activity (in, say, 500 voxels)



• Repeat 200 times

- Show someone a house:
- Measure brain activity (in, say, 500 voxels)



- Repeat 200 times
- Data: $X_1, ..., X_{200} \sim P \in \mathbb{R}^{500}$ and $Y_1, ..., Y_{200} \sim Q \in \mathbb{R}^{500}$. Test: Mean-difference alternative or general alternative.

Errors, and power

Two ways that a two-sample test could be wrong:

. False Positive: When $P = Q(H_0)$, but the test returns 1. The type-1 error α is the probability of a false positive.

High α - false discoveries - dangerous! Control at (say) 0.05.

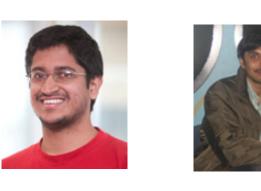
2. False Negative: When $P \neq Q$ (H_1), but the test returns 0. The type-2 error β is the probability of a false negative.

High β implies low power $\phi := 1 - \beta$ - very weak test incapable of detecting real differences that do exist.

A test is (classically) **consistent** if, while controlling α at any level (say 0.05), the power goes to 1 as the number of samples $n \to \infty$.

A test is (high-dim) **consistent** if, while controlling α at any level (say 0.05), the power goes to 1 as $(n,d) \rightarrow \infty$.

Aaditya Ramdas*, Sashank Reddi*, Barnabas Poczos, Aarti Singh, Larry Wasserman





Open Questions

Tests and Tradeoffs

- Parametric + Mean-Difference Alternative Eg: Threshold Hotelling's Statistic $(\bar{X} - \bar{Y})^T S^{-1} (\bar{X} - \bar{Y})$ Eg: Random Projections variant : Lopes-Jacob-Wainwright'12
- Nonparametric + Mean-Difference Alternative Eg: Diagonalize/drop S (SD, BS, CQ: Chen+Qin).
- Nonparametric + General Alternative Eg: Threshold the empirical Gaussian MMD (G- MMD^2) $\|x_{i} - x_{i}\|^{2}$ 1 $||y_{i}-y_{i}||^{2}$ o n n $||x_{i}-y_{i}||^{2}$

$$\frac{1}{\binom{n}{2}} \sum_{i \neq j} e^{-\frac{n - i - j - j - i}{\gamma^2}} + \frac{1}{\binom{n}{2}} \sum_{i \neq j} e^{-\frac{n - i - j - j - i}{\gamma^2}} - \frac{2}{n^2} \sum_{i=1} \sum_{j=1}^{n-1} e^{-\frac{n - i - j - j - i}{\gamma^2}}$$

Eg: Euclidean Energy Dist. (E-ED): $e^{-\frac{\|x_i - x_j\|^2}{\gamma^2}} \rightarrow \|x_i - x_j\|$ **Fact:** Population G-MMD² = 0 iff P = Q (also E-ED)

 $(\mathbf{Q1}: \mathsf{How} \mathsf{ do} \mathsf{ the} \mathsf{ latter} \mathsf{ tests} \mathsf{ perform} \mathsf{ in} \mathsf{ the} \mathsf{ former} \mathsf{ setting}?)$ **Q2:** What is the role of bandwidth γ on test power?

The Quadratic-time and Linear-time Statistics

Define (k is the Gaussian kernel or negative Euclidean distance) h := h(x, x', y, y') = k(x, x') + k(y, y') - k(x, y') - k(x', y)

$$MMD_u^2 = \frac{2}{n(n-1)} \sum_{i \neq j=1}^n h(x_i, x_j, y_i, y_j)$$

$$MMD_l^2 = \frac{1}{n/2} \sum_{i=1}^{n/2} h(x_{2i-1}, x_{2i}, y_{2i-1}, y_{2i})$$

It is easy to define a sub-quadratic time block-based variant that looks at n^{2x} entries for $0.5 \le x \le 1$.

(Q3: What is the tradeoff of computation vs power for these tests?

How did we land here? Misconceptions! **Q4:** How do $G-MMD^2$ and E-ED perform in high-dimensions?

- 1. Szekely¹: "The power of dCor test for independence is very good especially for high dimensions p,q". No proof!
- 2. Gretton et al²: "Estimation of MMD^2 is independent of d". True, but...Power \neq estimation error! MMD² is 1/poly(d)!
- Misleading Experiments 2, P and Q are mean-separated Gaussians with means (0, 0, ..., 0) and (1, 1, ..., 1).
- 4. Misleading Experiments³ Table 1^3 shows that unbiased dCov estimates dCov well (dCov hovers near 0 under independence).



Some Details...

Test based on MMD_1^2

Define $V = 2 \operatorname{Var}[h]$. By CLT

$$\frac{\sqrt{n}(\mathrm{MMD}_l^2 - \mathrm{MMD}^2)}{\sqrt{V}} \rightsquigarrow N(0, 1)$$

Let v be the empirical counterpart of V. Define the test

Reject when
$$\sqrt{n}rac{\mathrm{MMD}_l^2}{\sqrt{v}} > z_d$$

where $z_{\alpha} = \Phi^{-1}(1 - \alpha)$ for standard Gaussian cdf Φ , i.e. $P(Z > z_{\alpha}) = \alpha$ for standard Gaussian random variable Z.

- The type-1 error and power of this test changes with
- 1. Number of points n
- 2. Underlying dimensionality d of x, y
- 3. The signal-to-noise ratio $\Psi = \|\mathbb{E}_P[X] \mathbb{E}_Q[X]\|^2/\sigma^2$
- 4. The bandwidth of the Gaussian kernel γ

The "Classical" Power of MMD_1^2

If \Pr denotes the probability under H_1 , and Φ is the standard normal cdf, the power is

$$\Pr\left(\frac{\sqrt{n} \text{MMD}_{l}^{2}}{\sqrt{v}} > z_{\alpha}\right)$$

$$= \Pr\left(\frac{\sqrt{n} (\text{MMD}_{l}^{2} - \text{MMD}^{2})}{\sqrt{V}} > \sqrt{\frac{v}{V}} z_{\alpha} - \frac{\sqrt{n} \text{MMD}^{2}}{\sqrt{V}}\right)$$

$$\xrightarrow{n \to \infty} \Pr\left(Z > z_{\alpha} - \frac{\sqrt{n} \text{MMD}^{2}}{\sqrt{V}}\right)$$

$$= 1 - \Phi\left(z_{\alpha} - \frac{\sqrt{n} \text{MMD}^{2}}{\sqrt{V}}\right)$$

$$= \Phi\left(\frac{\sqrt{n} \text{MMD}^{2}}{\sqrt{V}} - z_{\alpha}\right)$$

This behaves like $\Phi(\sqrt{n})$ since the population MMD² and V are constants that are both independent of n.

Challenges in the high dimensional setting

- 1. MMD^2 depends on dimension
- 2. V depends on dimension
- 3. $(n,d) \rightarrow \infty$ at any rate
- 4. Does $v/V \to 1$ even if $(n, d) \to \infty$?

We will address these by

- Non-asymptotic, finite-sample Berry-Esseen theorem
- 2. Calculating MMD^2 , V explicitly by Taylor expansions
- Concentration bounds in terms of fourth moments

Explicit characterization of power as a function of $n, d, \Psi \approx KL(P, Q)$ in the high-dimensional setting as $(n,d) \rightarrow \infty$, for nonparametric P,Q differing in their means. 2. A clear and smooth computation-statistics tradeoff if computation scales as n^{2x} for $0.5 \le x \le 1$, then the power in the low SNR (low Ψ) regime is $\approx \Phi(n^x \Psi^2/\sqrt{d})$ 3. The power is independent of Gaussian kernel bandwidth,

as long as it is chosen large enough as $\Omega(\sqrt{d})$, which happens to be the choice made by the popular "median heuristic". **Energy Distance & Gaussian MMD have the same power** in this setting with a mean-difference between distributions. specialized tests that have been designed in the literature to test for mean differences (like Chen+Qin, Srivastava+Du).

5. Free Lunch! ED and GMMD have the same power as

(Linear) The right-hand side of the Berry Esseen lemma $10 \frac{\xi_3}{V^{3/2} \sqrt{n}}$ is actually $\leq 20/\sqrt{n}$, independent of d! **Implication:** Null, alternate distributions are always Gaussian.

Answers!

Summary of Some Results

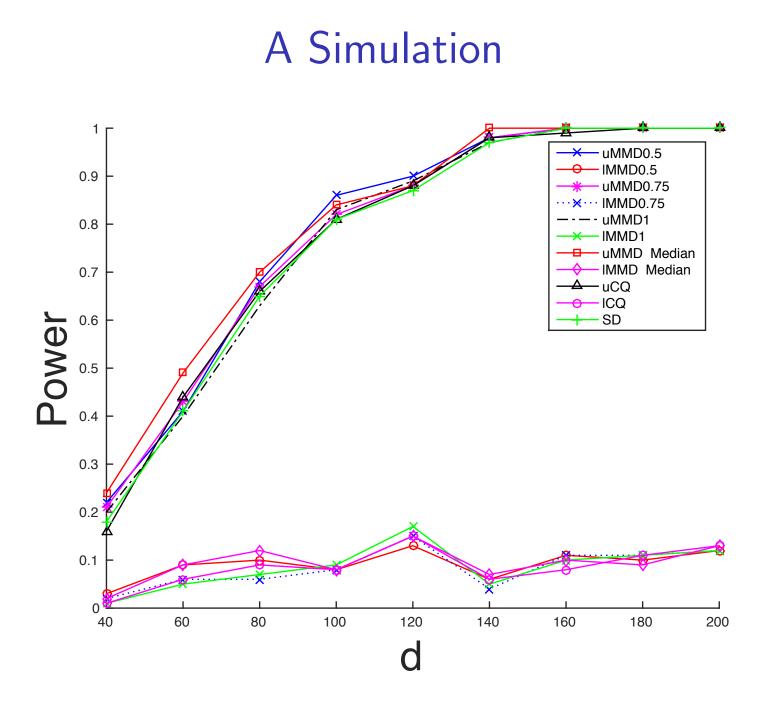


Figure : Parameters: P,Q Gaussians, $d = 40, 60, ..., 200, n = d, \Psi = 1$. **Top set**: U-statistics (G-MMD for many bandwidths, E-ED, CQ, SD). **Bottom set**: Linear-time statistics.

Summary of Some Techniques

1. G-MMD² $\approx \|\delta\|^2/\gamma^2$. Recall $\gamma = \Omega(\sqrt{d})$. **Implication:** This is why estimation error alone is misleading.

2. Variance V of test statistic also decays with d, but slower. **Implication:** This is why power decays with d.

Ratio of $v/V \to 1$ as $n \to \infty$ independent of how d grows. **Implication:** Studentization works fine.

(U-statistic) Martingale central limit theorem works out. **Implication:** Null, alternate distributions are Gaussian, *not* infinite sums of weighted chi-squared distributions (fixed d).

¹W'shop on Nonparam. Measures of Dependence (Columbia Univ, May'14) ²A Kernel Two Sample Test, Gretton et al. (JMLR'12)

³The dist. corr. t-test of indep. in high dim., Szekely et al (JMA'14)