

# Locally-Adaptive Kernel Tests\*

Zoltán Szabó, Éric Moulines

CMAP, École Polytechnique

## Quick summary

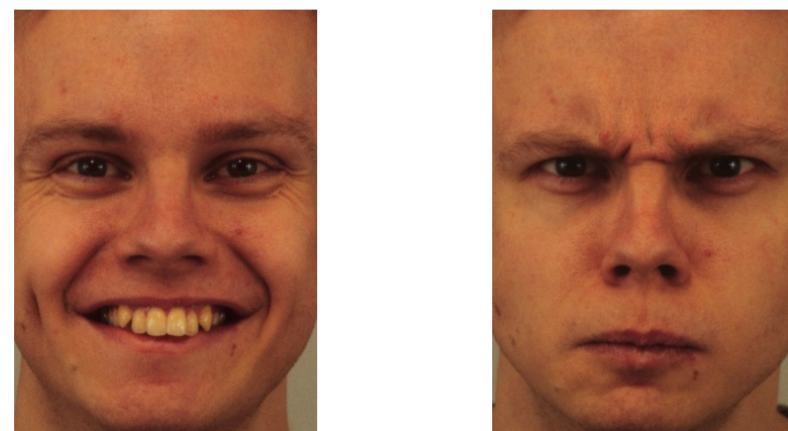
- Focus:
  - hypothesis testing  $\xrightarrow{\text{example}}$  2-sample testing.
  - Check if  $\mathbb{P} = \mathbb{Q}$  based on samples.
- Fast analytic kernel-based tests [1]:
  - Adaptivity [3]: parameters optimized for 'power'.  $\mathbb{P}, \mathbb{Q}$ : fixed.
- Challenge:
  - $\mathbb{P}$ : fixed. Sequence of alternatives:  $\mathbb{Q}_n \xrightarrow{n \rightarrow \infty} \mathbb{P}$ .
  - Adaptivity: realizable? Objective function?

## Two-sample testing

- Given:  $X = \{\mathbf{x}_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}, Y = \{\mathbf{y}_j\}_{j=1}^n \stackrel{i.i.d.}{\sim} \mathbb{Q}$ .
- Task: using  $X, Y$  test

$$H_0 : \mathbb{P} = \mathbb{Q}, \text{ vs } H_1 : \mathbb{P} \neq \mathbb{Q}.$$

- Example:  $\mathbf{x}_i = i^{\text{th}}$  happy face,  $\mathbf{y}_j = j^{\text{th}}$  sad face.



- Challenge (intuition): difference in emotions  $\rightarrow 0$ .

## Test power

- Test statistic:  $\hat{\lambda}_n = \hat{\lambda}_n(X, Y)$ , random.
- Decision:  $H_0$  is rejected if  $\hat{\lambda}_n$  is 'large'.
- $P(\text{we say } H_1 | H_0) \leq 1 - \alpha$ . Typically  $\alpha = 0.01$ .
- Power =  $1 - \underbrace{P(H_0 \text{ is accepted} | H_1)}_{\text{Type-II error}} \rightarrow \max$ .

## Representation of distributions

- Mean embedding:

$$\mu_{\mathbb{P}} := \underbrace{\int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x)}_{\text{Bochner integral}} \in \mathcal{H}_k.$$

- Maximum mean discrepancy:

$$MMD(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}.$$

In two-sample testing [2]:

- $\hat{\lambda}_n = \widehat{MMD}(\mathbb{P}_n, \mathbb{Q}_n)$ . Computation:  $\mathcal{O}(n^2)$ .

## Metric with analytic kernels

- Replace  $\|\cdot\|_{\mathcal{H}_k}$  in MMD with  $\|\cdot\|_{L^2(\mathcal{V})}$

$$\rho(\mathbb{P}, \mathbb{Q}) = \sqrt{\frac{1}{J} \sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2}.$$

$\mathcal{V} = \{\mathbf{v}_j\}_{j=1}^J$ : random locations. It is metric a.s. [1].

- Plug-in estimate,  $\mathcal{O}(n)$ -time

$$\hat{\rho}(\mathbb{P}, \mathbb{Q}) = \frac{\bar{\mathbf{z}}_n^T \bar{\mathbf{z}}_n}{J}, \quad \bar{\mathbf{z}}_n = \frac{1}{n} \sum_{i=1}^n \underbrace{[k(\mathbf{x}_i, \mathbf{v}_j) - k(\mathbf{y}_i, \mathbf{v}_j)]}_{=: \mathbf{z}(\mathbf{x}_i, \mathbf{y}_i)}.$$

- Modified test statistic,  $\chi_J^2$  null:

$$\hat{\lambda}_n = n \bar{\mathbf{z}}_n^T \hat{\Sigma}_n^{-1} \bar{\mathbf{z}}_n, \quad \Sigma_n = \widehat{\text{cov}}(\{\mathbf{z}(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n).$$

## Optimize for test power

- Power proxy [3]:

$$\lambda_n = n \mathbf{m}^T \Sigma^{-1} \mathbf{m}, \quad \mathbf{m} = \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\mathbf{z}(\mathbf{x}, \mathbf{y})], \quad \Sigma = \text{cov}_{\mathbf{x}, \mathbf{y}} [\mathbf{z}(\mathbf{x}, \mathbf{y})].$$

- Objective function:

$$(k^*, \mathcal{V}^*) := \arg \max_{k, \mathcal{V}} \mathbf{m}^T \Sigma^{-1} \mathbf{m}.$$

## Locally-adaptive test

- Instead of fixed  $\mathbb{P}, \mathbb{Q}$ :

$$\mathbb{P}, \quad \mathbb{Q}_n = (1 - \alpha_n)\mathbb{Q} + \alpha_n\mathbb{P} \xrightarrow{n \rightarrow \infty} \mathbb{P}, \quad \alpha_n \rightarrow 1.$$

- By linearity: with  $\alpha_n = 1 - \frac{1}{\sqrt{n}}$

$$\mathbf{r}_{\mathbb{Q}_n} - \mathbf{r}_{\mathbb{P}} = \frac{\boldsymbol{\delta}}{\sqrt{n}}, \quad \boldsymbol{\delta} = \mathbf{r}_{\mathbb{Q}-\mathbb{P}}, \quad \mathbf{r}_{\mathbb{M}} := [\mu_{\mathbb{M}}(\mathbf{v}_j)]_{j=1}^J.$$

- Assume:  $\{\mathbf{x}_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}, \{\mathbf{y}_{i,n}\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{Q}_n$ . Let

$$\Sigma = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}, \mathbf{y} \sim \mathbb{P}} [\mathbf{z}(\mathbf{x}, \mathbf{y}) \mathbf{z}^T(\mathbf{x}, \mathbf{y})],$$

$$\hat{\Sigma}_n = \widehat{\text{cov}}(\mathbf{z}(\mathbf{x}_i, \mathbf{y}_{i,n})), \quad \bar{\mathbf{z}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{z}(\mathbf{x}_i, \mathbf{y}_{i,n}).$$

- After some algebra:  $\theta := (k, \mathcal{V})$

$$n \bar{\mathbf{z}}_n^T \hat{\Sigma}_n^{-1} \bar{\mathbf{z}}_n \xrightarrow{w} \chi_J^2 \left( \underbrace{\boldsymbol{\delta}^T \Sigma^{-1} \boldsymbol{\delta}}_{=: \lambda(\boldsymbol{\delta}, \theta)} \right).$$

- $\Rightarrow$  Power against  $(\mathbb{Q}_n)_{n=1}^{\infty}$  goes to

$$\beta(\boldsymbol{\delta}, \theta) = P(\chi_J^2(\lambda) \geq q) = M_{\frac{J}{2}}(\sqrt{\lambda}, q).$$

## Objective

Maximize the **worst-case local power**,

$$(k^*, \mathcal{V}^*) := \arg \max_{\theta} \min_{\boldsymbol{\delta}} \lambda(\boldsymbol{\delta}, \theta).$$

## References

- [1] K. Chwialkowski, A. Ramdas, D. Sejdinovic, and A. Gretton. Fast two-sample testing with analytic representations of probability measures. In *NIPS*, pages 1972–1980, 2015.
- [2] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *JMLR*, 13:723–773, 2012.
- [3] W. Jitkrittum, Z. Szabó, K. Chwialkowski, and A. Gretton. Interpretable distribution features with maximum testing power. In *NIPS*, pages 181–189, 2016.