# Distinguishing Distributions with Maximum Testing Power

#### Zoltán Szabó (Gatsby Unit, UCL)



Wittawat Jitkrittum



Kacper Chwialkowski



Arthur Gretton

Realeyes, Budapest August 24, 2016

Zoltán Szabó Distinguishing Distributions with Maximum Testing Power

- Motivating examples: NLP, computer vision.
- Two-sample test: t-test  $\rightarrow$  distribution features.
- Linear-time, interpretable, high-power, nonparametric t-test.
- Numerical illustrations.

# Motivating examples

#### Motivating example-1: NLP

- Given: two categories of documents (Bayesian inference, neuroscience).
- Task:
  - test their distinguishability,
  - most discriminative words  $\rightarrow$  interpretability.



#### Motivating example-2: computer vision





- Given: two sets of faces (happy, angry).
- Task:
  - check if they are different,
  - determine the most discriminative features/regions.

Contribution:

- We propose a nonparametric t-test.
- It gives a reason why  $H_0$  is rejected.
- It has high test power.
- It runs in linear time.

Contribution:

- We propose a nonparametric t-test.
- It gives a reason why  $H_0$  is rejected.
- It has high test power.
- It runs in linear time.

Dissemination, code:

- NIPS-2016 [Jitkrittum et al., 2016]: full oral = top 1.84%.
- https://github.com/wittawatj/interpretable-test.

# Two-sample test, distribution features

#### What is a two-sample test?

• Given:

• 
$$X = {\mathbf{x}_i}_{i=1}^n \overset{i.i.d.}{\sim} \mathbb{P}, \ \mathbf{Y} = {\mathbf{y}_j}_{j=1}^n \overset{i.i.d.}{\sim} \mathbb{Q}.$$

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

• Given:

• 
$$X = {\mathbf{x}_i}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}, \ \mathbf{Y} = {\mathbf{y}_j}_{j=1}^n \stackrel{i.i.d.}{\sim} \mathbb{Q}.$$

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

• Problem: using X, Y test

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

• 
$$X = {\mathbf{x}_i}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}, \ \mathbf{Y} = {\mathbf{y}_j}_{j=1}^n \stackrel{i.i.d.}{\sim} \mathbb{Q}.$$

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

• Problem: using X, Y test

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

• Assume  $X, Y \subset \mathbb{R}^d$ .

#### Ingredients of two-sample test

- Test statistic:  $\hat{\lambda}_n = \hat{\lambda}_n(X, Y)$ , random.
- Significance level:  $\alpha = 0.01$ .

• Under 
$$H_0$$
:  $P_{H_0}(\hat{\lambda}_n \leq T_{\alpha}) = 1 - \alpha$ .

correctly accepting  $H_0$ 



#### Ingredients of two-sample test



#### Towards representations of distributions: $\mathbb{E}X$

- Given: 2 Gaussians with different means.
- Solution: *t*-test.



## Towards representations of distributions: $\mathbb{E}X^2$

- Setup: 2 Gaussians; same means, different variances.
- Idea: look at 2nd-order features of RVs.



#### Towards representations of distributions: $\mathbb{E}X^2$

- Setup: 2 Gaussians; same means, different variances.
- Idea: look at 2nd-order features of RVs.

• 
$$\varphi_x = x^2 \Rightarrow \text{difference in } \mathbb{E}X^2$$



#### Towards representations of distributions: further moments

- Setup: a Gaussian and a Laplacian distribution.
- Challenge: their means and variances are the same.
- Idea: look at higher-order features.



#### Let us consider feature/distribution representations!

Zoltán Szabó Distinguishing Distributions with Maximum Testing Power

#### Kernel: similarity between features

• Given:  $\mathbf{x}$  and  $\mathbf{x}'$  objects (images or texts).

#### Kernel: similarity between features

- Given: **x** and **x**' objects (images or texts).
- Question: how similar they are?

#### Kernel: similarity between features

- Given:  $\mathbf{x}$  and  $\mathbf{x}'$  objects (images or texts).
- Question: how similar they are?
- Define features of the objects:

 $\varphi_{\mathbf{x}}$ : features of  $\mathbf{x}$ ,  $\varphi_{\mathbf{x}'}$ : features of  $\mathbf{x}'$ .

• Kernel: inner product of these features

 $\boldsymbol{k}(\mathbf{x},\mathbf{x}') := \langle \varphi_{\mathbf{x}}, \varphi_{\mathbf{x}'} \rangle \,.$ 

• Polynomial kernel:

$$k(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + \gamma)^{p}.$$

• Gaussian kernel:

$$k(\mathbf{x},\mathbf{y}) = e^{-\gamma \|\mathbf{x}-\mathbf{y}\|_2^2}.$$



Zoltán Szabó Distinguishing Distributions with Maximum Testing Power







 $\widehat{MMD^2}(\mathbb{P},\mathbb{Q}) = \overline{K_{\mathbb{P},\mathbb{P}}} + \overline{K_{\mathbb{Q},\mathbb{Q}}} - 2\overline{K_{\mathbb{P},\mathbb{Q}}} \quad \text{(without diagonals in } \overline{K_{\mathbb{P},\mathbb{P}}}, \ \overline{K_{\mathbb{Q},\mathbb{Q}}})$ 

Zoltán Szabó

Distinguishing Distributions with Maximum Testing Power

• Kernel recall:  $k(\mathbf{x}, \mathbf{x}') = \langle \varphi_{\mathbf{x}}, \varphi_{\mathbf{x}'} \rangle$ .

- Kernel recall:  $k(\mathbf{x}, \mathbf{x}') = \langle \varphi_{\mathbf{x}}, \varphi_{\mathbf{x}'} \rangle$ .
- Feature of  $\mathbb{P}$  (mean embedding):

$$\mu_{\mathbb{P}} := \mathbb{E}_{\mathbf{x} \sim \mathbb{P}}[\varphi_{\mathbf{x}}].$$

- Kernel recall:  $k(\mathbf{x}, \mathbf{x}') = \langle \varphi_{\mathbf{x}}, \varphi_{\mathbf{x}'} \rangle$ .
- Feature of  $\mathbb{P}$  (mean embedding):

$$\mu_{\mathbb{P}} := \mathbb{E}_{\mathbf{x} \sim \mathbb{P}}[\varphi_{\mathbf{x}}].$$

• Previous quantity: unbiased estimate of

$$MMD^2(\mathbb{P},\mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|^2$$
.

- Kernel recall:  $k(\mathbf{x}, \mathbf{x}') = \langle \varphi_{\mathbf{x}}, \varphi_{\mathbf{x}'} \rangle$ .
- Feature of  $\mathbb{P}$  (mean embedding):

$$\mu_{\mathbb{P}} := \mathbb{E}_{\mathbf{x} \sim \mathbb{P}}[\varphi_{\mathbf{x}}].$$

• Previous quantity: unbiased estimate of

$$MMD^2(\mathbb{P},\mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|^2$$

- Valid test [Gretton et al., 2012]. Challenges:
  - Threshold choice: 'ugly' asymptotics of  $n\widehat{MMD^2}(\mathbb{P},\mathbb{P})$ .
  - Prest statistic: quadratic time complexity.

## Linear-time tests

Recall:

$$MMD^{2}(\mathbb{P},\mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}(k)}^{2}.$$

• Changing [Chwialkowski et al., 2015] this to

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

with random  $\{\mathbf{v}_j\}_{j=1}^J$  test locations

 $\rho$  is a metric (a.s.). How do we estimate it? Distribution under  $H_0$ ?

## Estimation

#### Estimate

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

where  $\hat{\mu}_{\mathbb{P}}(\mathbf{v}) = \frac{1}{n} \sum_{i=1}^{n} k(\mathbf{x}_i, \mathbf{v})$ . Using  $k(\mathbf{x}, \mathbf{v}) = e^{-\frac{\|\mathbf{x}-\mathbf{v}\|^2}{2\sigma^2}}$ ,



Zoltán Szabó Distinguishing Distributions with Maximum Testing Power

#### Estimation – continued

$$\widehat{\rho^{2}(\mathbb{P},\mathbb{Q})} = \frac{1}{J} \sum_{j=1}^{J} [\widehat{\mu}_{\mathbb{P}}(\mathbf{v}_{j}) - \widehat{\mu}_{\mathbb{Q}}(\mathbf{v}_{j})]^{2}$$

$$= \frac{1}{J} \sum_{j=1}^{J} \left[ \frac{1}{n} \sum_{i=1}^{n} k(\mathbf{x}_{i}, \mathbf{v}_{j}) - \frac{1}{n} \sum_{i=1}^{n} k(\mathbf{y}_{i}, \mathbf{v}_{j}) \right]^{2} = \frac{1}{J} \sum_{j=1}^{J} (\overline{z}_{n})_{j}^{2} = \frac{1}{J} \overline{z}_{n}^{T} \overline{z}_{n},$$
where  $\overline{z}_{n} = \frac{1}{n} \sum_{i=1}^{n} \underbrace{[k(\mathbf{x}_{i}, \mathbf{v}_{j}) - k(\mathbf{y}_{i}, \mathbf{v}_{j})]_{j=1}^{J}}_{=:\overline{z}_{i}} \in \mathbb{R}^{J}.$ 

• Good news: estimation is linear in n!

• Bad news: intractable null distr. =  $\sqrt{n\rho^2(\mathbb{P},\mathbb{P})} \xrightarrow{w}$  sum of J correlated  $\chi^2$ .

Modified test statistic:

$$\hat{\lambda}_n = n \bar{\mathbf{z}}_n^T \boldsymbol{\Sigma}_n^{-1} \bar{\mathbf{z}}_n,$$

where  $\Sigma_n = cov(\{\mathbf{z}_i\}_i)$ .

- Under H<sub>0</sub>:
  - $\hat{\lambda}_n \xrightarrow{w} \chi^2(J)$ .  $\Rightarrow$  Easy to get the  $(1 \alpha)$ -quantile!

# Our idea

- Until this point: test locations ( $\mathcal{V}$ ) are fixed.
- Instead: choose  $\boldsymbol{\theta} = \{\mathcal{V}, \sigma\}$  to

maximize lower bound on the test power.

- Until this point: test locations  $(\mathcal{V})$  are fixed.
- Instead: choose  $\theta = \{\mathcal{V}, \sigma\}$  to

maximize lower bound on the test power.

Theorem (Lower bound on power)

For large n, test power  $\geq L(\lambda_n)$ ; L: explicit function, increasing.

• Here,

• 
$$\lambda_n = n\mu^T \Sigma^{-1}\mu$$
: population version of  $\hat{\lambda}_n$ .  
•  $\mu = \mathbb{E}_{xy}[z_1], \Sigma = \mathbb{E}_{xy}[(z_1 - \mu)(z_1 - \mu)^T].$ 

Training objective  $\hat{\lambda}_n(X_{tr}, Y_{tr})$  converges to  $\lambda_n$ .

- But  $\lambda_n$  is unknown.
- Split (X, Y) into  $(X_{tr}, Y_{tr})$  and  $(X_{te}, Y_{te})$ . Use  $\hat{\lambda}_n(X_{tr}, Y_{tr}) \approx \lambda_n$ .

Theorem (Guarantee on objective approximation)

$$\left|\sup_{\mathcal{V},\mathcal{K}} \bar{\mathbf{z}}_n^T (\boldsymbol{\Sigma}_n + \gamma_n)^{-1} \bar{\mathbf{z}}_n - \sup_{\mathcal{V},\mathcal{K}} \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right| = \mathcal{O}\left(n^{-\frac{1}{4}}\right)$$

Training objective 
$$\hat{\lambda}_n(X_{tr}, Y_{tr})$$
 converges to  $\lambda_n$ .

- But  $\lambda_n$  is unknown.
- Split (X, Y) into  $(X_{tr}, Y_{tr})$  and  $(X_{te}, Y_{te})$ . Use  $\hat{\lambda}_n(X_{tr}, Y_{tr}) \approx \lambda_n$ .

Theorem (Guarantee on objective approximation)

$$\left|\sup_{\mathcal{V},\mathcal{K}} \bar{\mathsf{z}}_n^T (\boldsymbol{\Sigma}_n + \gamma_n)^{-1} \bar{\mathsf{z}}_n - \sup_{\mathcal{V},\mathcal{K}} \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}\right| = \mathcal{O}\left(n^{-\frac{1}{4}}\right)$$

Examples:

$$\begin{split} \mathcal{K} &= \{k_{\sigma}(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x} - \mathbf{y}\|^2} : \sigma > 0\},\\ \mathcal{K} &= \{k_{\mathbf{A}}(\mathbf{x}, \mathbf{y}) = e^{-(\mathbf{x} - \mathbf{y})^T \mathbf{A}(\mathbf{x} - \mathbf{y})} : \mathbf{A} > 0\}. \end{split}$$

# Numerical demos

#### Parameter settings

- Gaussian kernel ( $\sigma$ ).  $\alpha = 0.01$ . J = 1. Repeat 500 trials.
- Report

$$\mathbb{P}(\text{reject } H_0) \approx \frac{\#\text{times } \hat{\lambda}_n > T_\alpha \text{ holds}}{\#\text{trials}}$$

- Compare 4 methods
  - **ME-full**: Optimize  $\mathcal{V}$  and Gaussian bandwidth  $\sigma$ .
  - **ME-grid**: Optimize  $\sigma$ . Fix  $\mathcal{V}$  [Chwialkowski et al., 2015].
  - MMD-quad: Test with quadratic-time MMD [Gretton et al., 2012].
  - MMD-lin: Test with linear-time MMD [Gretton et al., 2012].
- Optimize kernels to power in MMD-lin, MMD-quad.

#### NLP: discrimination of document categories

- 5903 NIPS papers (1988-2015).
- Keyword-based category assignment into 4 groups:
  - Bayesian inference, Deep learning, Learning theory, Neuroscience
- d = 2000 nouns. TF-IDF representation.

Problem	n <sup>te</sup>	ME-full	ME-grid	MMD-quad	MMD-lin
<ol> <li>Bayes-Bayes</li> </ol>	215	.012	.018	.022	.008
2. Bayes-Deep	216	.954	.034	.906	.262
3. Bayes-Learn	138	.990	.774	1.00	.238
4. Bayes-Neuro	394	1.00	.300	.952	.972
5. Learn-Deep	149	.956	.052	.876	.500
6. Learn-Neuro	146	.960	.572	1.00	.538

• Performance of ME-full  $[\mathcal{O}(n)]$  is comparable to MMD-quad  $[\mathcal{O}(n^2)]$ .

- Aggregating over trials; example: 'Bayes-Neuro'.
- Most discriminative words:

spike, markov, cortex, dropout, recurr, iii, gibb.

- learned test locations: highly interpretable,
- 'markov', 'gibb' (< Gibbs): Bayesian inference,
- 'spike', 'cortex': key terms in neuroscience.

• Aggregating over trials; example: 'Bayes-Neuro'.

• Least dicriminative ones:

circumfer, bra, dominiqu, rhino, mitra, kid, impostor.

## Distinguish positive/negative emotions

- Karolinska Directed Emotional Faces (KDEF) [Lundqvist et al., 1998].
- 70 actors = 35 females and 35 males.
- $d = 48 \times 34 = 1632$ . Grayscale. Pixel features.



happy

neutral



surprised





angry



afraid

disgusted

Problem	n <sup>te</sup>	ME-full	ME-grid	MMD-quad	MMD-lin			
$\pm$ vs. $\pm$	201	.010	.012	.018	.008			
+ vs. $-$	201	.998	.656	1.00	.578			
- AND								

Learned test location (averaged) = ٩

- We proposed a nonparametric t-test:
  - linear time,
  - high-power ( $\approx$  'MMD-quad'),
- 2 demos: discriminating
  - documents of different categories,
  - positive/negative emotions.

# Thank you for the attention!



**Acknowledgements**: This work was supported by the Gatsby Charitable Foundation.

- Non-convexity, informative features.
- Number of locations (J).
- Computational complexity.
- Estimation of *MMD*<sup>2</sup>.

## Non-convexity, informative features

• 2D problem:

 $\mathbb{P} := \mathcal{N}([0;0], \mathbf{I}),$  $\mathbb{Q} := \mathcal{N}([1;0], \mathbf{I}).$ 

- $\mathcal{V} = \{ \mathbf{v}_1, \mathbf{v}_2 \}.$
- Fix v<sub>1</sub> to ▲.
- Contour plot of  $\mathbf{v}_2 \mapsto \hat{\lambda}_n(\{\mathbf{v}_1, \mathbf{v}_2\}).$



## Number of locations (J)

#### • Small *J*:

- often enough to detect the difference of  $\mathbb P$  &  $\mathbb Q.$
- few distinguishing regions to reject  $H_0$ .
- faster test.
- Very large J:
  - test power need not increase monotonically in J (more locations ⇒ statistic can gain in variance).
  - defeats the purpose of a linear-time test.

- Optimization & testing: linear in n.
- Testing:  $\mathcal{O}(ndJ + nJ^2 + J^3)$ .
- Optimization:  $\mathcal{O}\left(ndJ^2 + J^3\right)$  per gradient ascent.

Squared difference between feature means:

$$\begin{split} \mathcal{M}\mathcal{M}\mathcal{D}^{2}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}}^{2} = \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} \\ &= \langle \mu_{\mathbb{P}}, \mu_{\mathbb{P}} \rangle_{\mathcal{H}} + \langle \mu_{\mathbb{Q}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} - 2 \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} \\ &= \mathbb{E}_{\mathbb{P},\mathbb{P}}k(\mathbf{x},\mathbf{x}') + \mathbb{E}_{\mathbb{Q},\mathbb{Q}}k(\mathbf{y},\mathbf{y}') - 2\mathbb{E}_{\mathbb{P},\mathbb{Q}}k(\mathbf{x},\mathbf{y}). \end{split}$$

Squared difference between feature means:

$$\begin{split} \mathsf{MMD}^{2}(\mathbb{P},\mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}}^{2} = \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} \\ &= \langle \mu_{\mathbb{P}}, \mu_{\mathbb{P}} \rangle_{\mathcal{H}} + \langle \mu_{\mathbb{Q}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} - 2 \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}} \\ &= \mathbb{E}_{\mathbb{P},\mathbb{P}} k(\mathbf{x}, \mathbf{x}') + \mathbb{E}_{\mathbb{Q},\mathbb{Q}} k(\mathbf{y}, \mathbf{y}') - 2\mathbb{E}_{\mathbb{P},\mathbb{Q}} k(\mathbf{x}, \mathbf{y}). \end{split}$$

Unbiased empirical estimate for  $\{\mathbf{x}_i\}_{i=1}^n \sim \mathbb{P}$ ,  $\{\mathbf{y}_j\}_{j=1}^n \sim \mathbb{Q}$ :

$$\widehat{MMD}^2(\mathbb{P},\mathbb{Q}) = \overline{K_{\mathbb{P},\mathbb{P}}} + \overline{K_{\mathbb{Q},\mathbb{Q}}} - 2\overline{K_{\mathbb{P},\mathbb{Q}}}.$$

Chwialkowski, K., Ramdas, A., Sejdinovic, D., and Gretton, A. (2015).

Fast Two-Sample Testing with Analytic Representations of Probability Measures.

In *Neural Information Processing Systems (NIPS)*, pages 1981–1989.

Gretton, A., Borgwardt, K., Rasch, M., Schölkopf, B., and Smola, A. (2012).

A kernel two-sample test.

Journal of Machine Learning Research, 13:723–773.

Jitkrittum, W., Szabó, Z., Chwialkowski, K., and Gretton, A. (2016).
 Interpretable distribution features with maximum testing power.
 In Neural Information Processing Systems (NIPS).

(accepted).

Lundqvist, D., Flykt, A., and Öhman, A. (1998).

The Karolinska directed emotional faces-KDEF. Technical report, ISBN 91-630-7164-9.