

# Adaptive linear-time nonparametric two-sample testing

Zoltán Szabó (École Polytechnique)



Wittawat Jitkrittum



Kacper Chwialkowski



Arthur Gretton

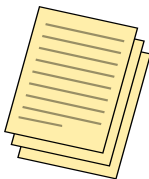
CMAP seminar, Palaiseau  
École Polytechnique November 22, 2016

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

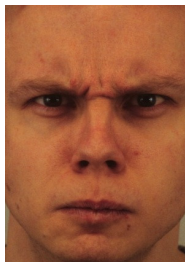
# Motivating examples

# 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:

- $X = \{\mathbf{x}_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}$ ,  $Y = \{\mathbf{y}_j\}_{j=1}^n \stackrel{i.i.d.}{\sim} \mathbb{Q}$ .
- Example:  $\mathbf{x}_i = i^{th}$  happy face,  $\mathbf{y}_j = j^{th}$  sad face.

# What is a two-sample test?

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

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

$$H_1 : \mathbb{P} \neq \mathbb{Q}.$$

# What is a two-sample test?

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

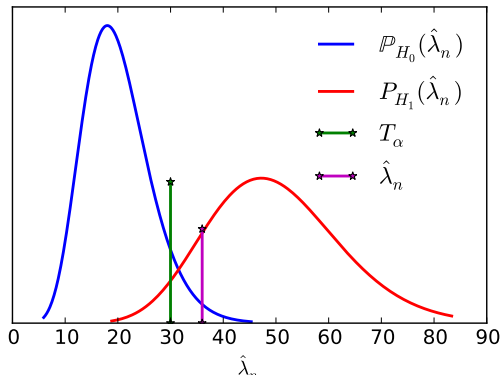
$$H_0 : \mathbb{P} = \mathbb{Q}, \text{ vs}$$

$$H_1 : \mathbb{P} \neq \mathbb{Q}.$$

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

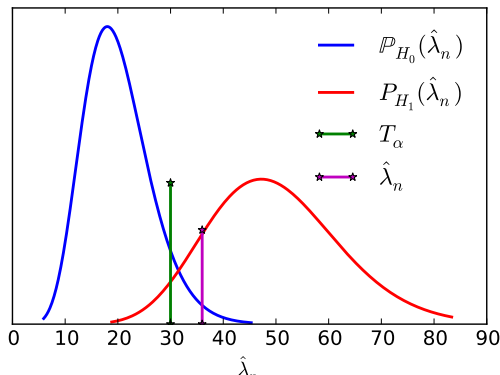
# Ingredients of two-sample test

- Test statistic:  $\hat{\lambda}_n = \hat{\lambda}_n(X, Y)$ , random.
- Significance level:  $\alpha = 0.01$ .
- Under  $H_0$ :  $P_{H_0}(\underbrace{\hat{\lambda}_n \leq T_\alpha}_{\text{correctly accepting } H_0}) = 1 - \alpha$ .



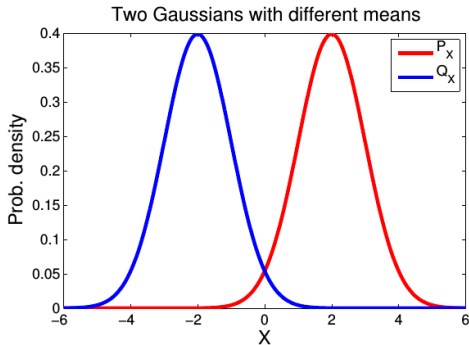
# Ingredients of two-sample test

- Test statistic:  $\hat{\lambda}_n = \hat{\lambda}_n(X, Y)$ , random.
- Significance level:  $\alpha = 0.01$ .
- Under  $H_0$ :  $P_{H_0}(\underbrace{\hat{\lambda}_n \leq T_\alpha}_{\text{correctly accepting } H_0}) = 1 - \alpha$ .
- Under  $H_1$ :  $P_{H_1}(T_\alpha < \hat{\lambda}_n) = P(\text{correctly rejecting } H_0) =: \text{power}$ .



# Towards representations of distributions: EX

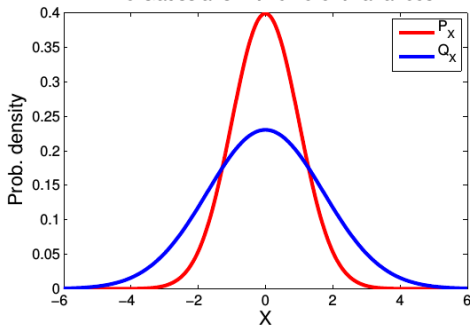
- Given: 2 Gaussians with (possibly) different means.
- Solution: *t*-test.



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

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

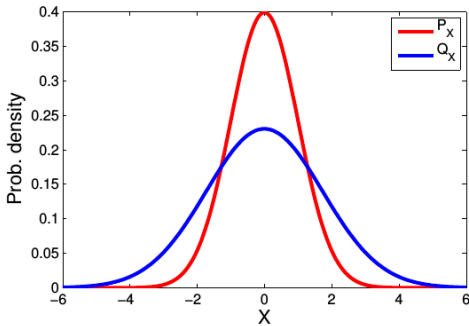
Two Gaussians with different variances



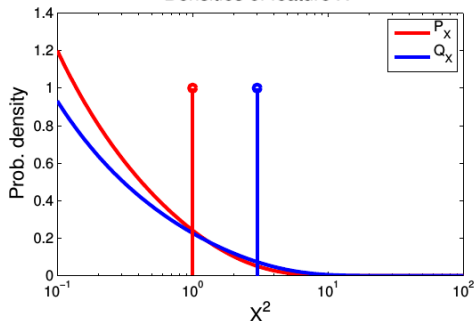
# Towards representations of distributions: $\mathbb{E}X^2$

- Setup: 2 Gaussians; same means, different variances.
- Idea: look at 2nd-order features of RVs.
- $\varphi_X = x^2 \Rightarrow$  difference in  $\mathbb{E}X^2$ .

Two Gaussians with different variances



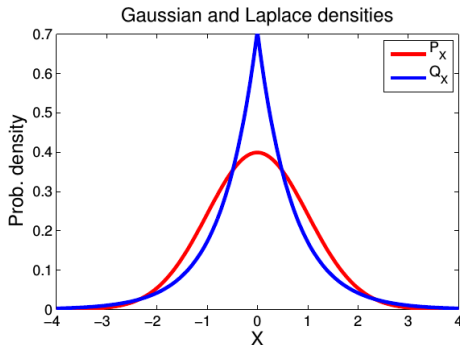
Densities of feature  $X^2$





# 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!

# Kernel: similarity between features

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

# Kernel: similarity between features

- Given:  $\mathbf{x}$  and  $\mathbf{x}'$  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

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

# Kernel examples on $\mathbb{R}^d$ ( $\gamma > 0, p \in \mathbb{Z}^+$ )

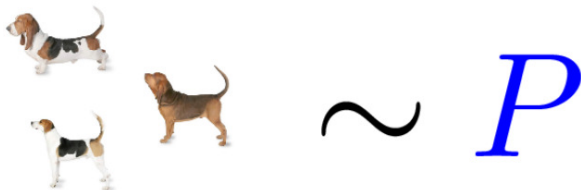
- Polynomial kernel:

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

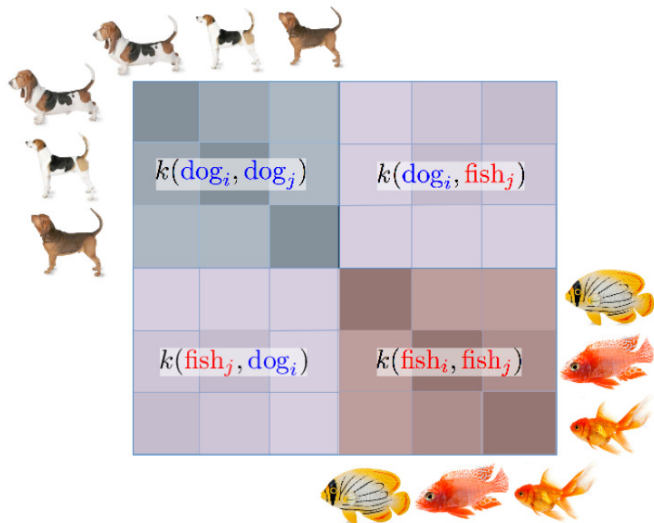
- Gaussian kernel:

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

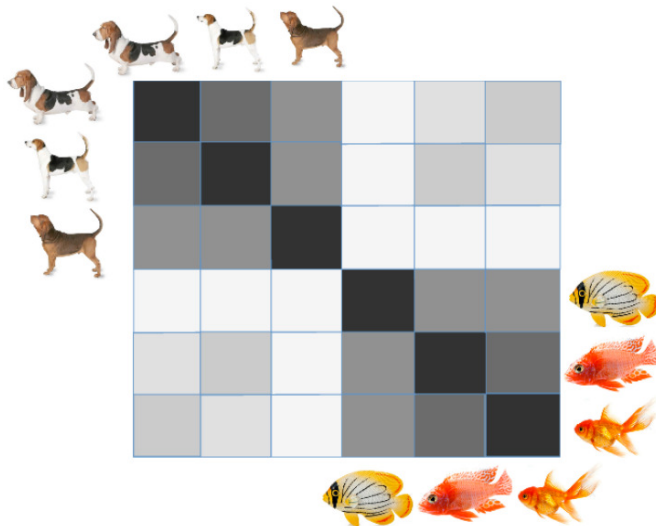
# Towards distribution features



# Towards distribution features

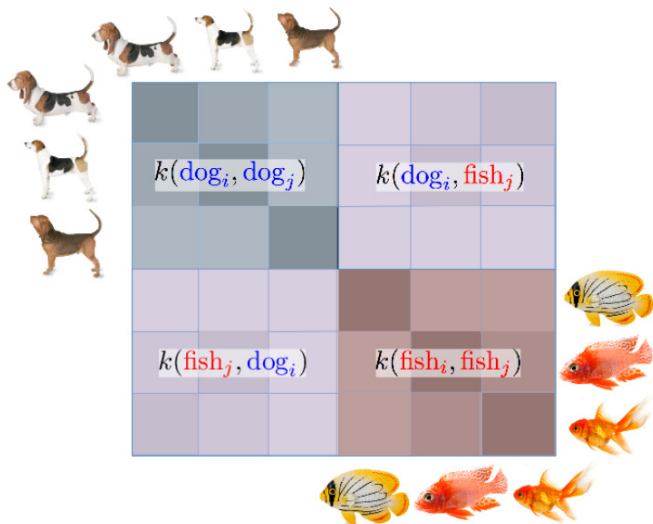


# Towards distribution features





# Towards distribution features



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

<sup>†</sup>  $\widehat{MMD}$  illustration credit: Arthur Gretton

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

# Kernel $\rightarrow$ distribution feature

- 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 $\rightarrow$ distribution feature

- 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 $\rightarrow$ distribution feature

- 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:
  - 1 Threshold choice: 'ugly' asymptotics of  $n\widehat{MMD^2}(\mathbb{P}, \mathbb{P})$ .
  - 2 Test statistic: quadratic time complexity.
  - 3 Witness  $\in \mathcal{H}(k)$ : can be hard to interpret.

# Linear-time tests

# Linear-time 2-sample test

- Recall:

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

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

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

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

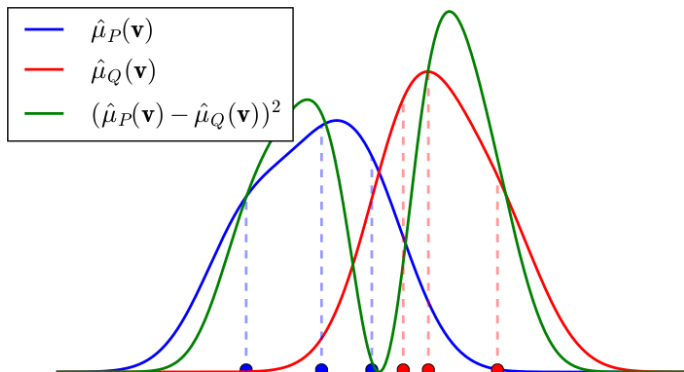
$\rho$  is a metric (a.s.). How do we estimate it? Distribution under  $H_0$ ?

# Estimation

Compute

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

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





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

$$\begin{aligned}\widehat{\rho^2(\mathbb{P}, \mathbb{Q})} &= \frac{1}{J} \sum_{j=1}^J [\hat{\mu}_{\mathbb{P}}(\mathbf{v}_j) - \hat{\mu}_{\mathbb{Q}}(\mathbf{v}_j)]^2 \\ &= \frac{1}{J} \sum_{j=1}^J \left[ \frac{1}{n} \sum_{i=1}^n k(\mathbf{x}_i, \mathbf{v}_j) - \frac{1}{n} \sum_{i=1}^n k(\mathbf{y}_i, \mathbf{v}_j) \right]^2\end{aligned}$$

$$\begin{aligned}\widehat{\rho^2(\mathbb{P}, \mathbb{Q})} &= \frac{1}{J} \sum_{j=1}^J [\hat{\mu}_{\mathbb{P}}(\mathbf{v}_j) - \hat{\mu}_{\mathbb{Q}}(\mathbf{v}_j)]^2 \\ &= \frac{1}{J} \sum_{j=1}^J \left[ \frac{1}{n} \sum_{i=1}^n k(\mathbf{x}_i, \mathbf{v}_j) - \frac{1}{n} \sum_{i=1}^n k(\mathbf{y}_i, \mathbf{v}_j) \right]^2 = \frac{1}{J} \sum_{j=1}^J (\bar{\mathbf{z}}_n)_j^2 = \frac{1}{J} \bar{\mathbf{z}}_n^T \bar{\mathbf{z}}_n,\end{aligned}$$

$$\text{where } \bar{\mathbf{z}}_n = \frac{1}{n} \sum_{i=1}^n \underbrace{[k(\mathbf{x}_i, \mathbf{v}_j) - k(\mathbf{y}_i, \mathbf{v}_j)]}_{=:\mathbf{z}_i} \Big|_{j=1}^J \in \mathbb{R}^J.$$

# Estimation – continued

$$\begin{aligned}\widehat{\rho^2(\mathbb{P}, \mathbb{Q})} &= \frac{1}{J} \sum_{j=1}^J [\hat{\mu}_{\mathbb{P}}(\mathbf{v}_j) - \hat{\mu}_{\mathbb{Q}}(\mathbf{v}_j)]^2 \\ &= \frac{1}{J} \sum_{j=1}^J \left[ \frac{1}{n} \sum_{i=1}^n k(\mathbf{x}_i, \mathbf{v}_j) - \frac{1}{n} \sum_{i=1}^n k(\mathbf{y}_i, \mathbf{v}_j) \right]^2 = \frac{1}{J} \sum_{j=1}^J (\bar{\mathbf{z}}_n)_j^2 = \frac{1}{J} \bar{\mathbf{z}}_n^T \bar{\mathbf{z}}_n,\end{aligned}$$

where  $\bar{\mathbf{z}}_n = \frac{1}{n} \sum_{i=1}^n \underbrace{[k(\mathbf{x}_i, \mathbf{v}_j) - k(\mathbf{y}_i, \mathbf{v}_j)]}_{=: \mathbf{z}_i} \Big|_{j=1}^J \in \mathbb{R}^J$ .

- Good news: estimation is linear in  $n$ !
- Bad news: intractable null distr.  $= \sqrt{n} \widehat{\rho^2(\mathbb{P}, \mathbb{P})} \xrightarrow{w} \text{sum of } J \text{ correlated } \chi^2$ .

# Normalized version gives tractable null

- Modified test statistic:

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

where  $\Sigma_n = \text{cov}(\{\mathbf{z}_i\}_{i=1}^n)$ .

- Under  $H_0$ :
  - $\hat{\lambda}_n \xrightarrow{w} \chi^2(J)$ .  $\Rightarrow$  Easy to get the  $(1 - \alpha)$ -quantile!

# Our idea

- Until this point: test locations ( $\mathcal{V}$ ) are fixed.
- Instead: choose  $\theta = \{\mathcal{V}, \sigma\}$  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  $\hat{\lambda}_n$ .
  - $\mu = \mathbb{E}_{\mathbf{xy}}[\mathbf{z}_1]$ ,  $\Sigma = \mathbb{E}_{\mathbf{xy}}[(\mathbf{z}_1 - \mu)(\mathbf{z}_1 - \mu)^T]$ .



# Convergence of the $\lambda_n$ estimator

But  $\lambda_n$  is **unknown**. Split  $(X, Y)$  into  $(X_{tr}, Y_{tr})$  and  $(X_{te}, Y_{te})$ .

- ① Locations, kernel parameter:  $\hat{\theta} = \arg \max_{\theta} \hat{\lambda}_{\frac{n}{2}}^{tr}(\theta)$ .

# Convergence of the $\lambda_n$ estimator

But  $\lambda_n$  is **unknown**. Split  $(X, Y)$  into  $(X_{tr}, Y_{tr})$  and  $(X_{te}, Y_{te})$ .

- ① Locations, kernel parameter:  $\hat{\theta} = \arg \max_{\theta} \hat{\lambda}_{\frac{n}{2}}^{tr}(\theta)$ .
- ② Test statistic:  $\hat{\lambda}_{\frac{n}{2}}^{te}(\hat{\theta})$ .

# Convergence of the $\lambda_n$ estimator

Theorem (Guarantee on objective approximation,  $\gamma_n \rightarrow 0$ )

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

# Convergence of the $\lambda_n$ estimator

Theorem (Guarantee on objective approximation,  $\gamma_n \rightarrow 0$ )

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

Examples:

$$\mathcal{K} = \left\{ k_{\sigma}(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2\sigma^2}} : \sigma > 0 \right\},$$

$$\mathcal{K} = \left\{ k_{\mathbf{A}}(\mathbf{x}, \mathbf{y}) = e^{-(\mathbf{x}-\mathbf{y})^T \mathbf{A} (\mathbf{x}-\mathbf{y})} : \mathbf{A} \succ 0 \right\}.$$

# Numerical demos

- 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  $\sigma$ . Random  $\mathcal{V}$  [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	$n^{te}$	ME-full	ME-grid	MMD-quad	MMD-lin
1. Bayes-Bayes	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 [ $\mathcal{O}(n^2)$ ].

# NLP: most/least discriminative words

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

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

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

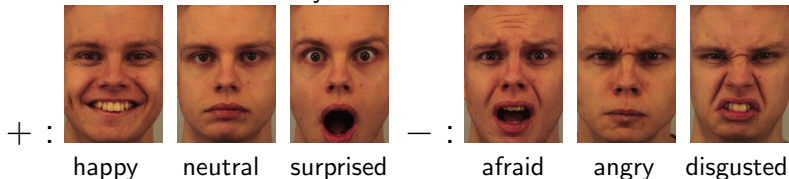


# NLP: most/least discriminative words

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



Problem	$n^{te}$	ME-full	ME-grid	MMD-quad	MMD-lin
$\pm$ vs. $\pm$	201	.010	.012	.018	.008
$+$ vs. $-$	201	.998	.656	1.00	.578



- Learned test location (averaged) =

- 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.

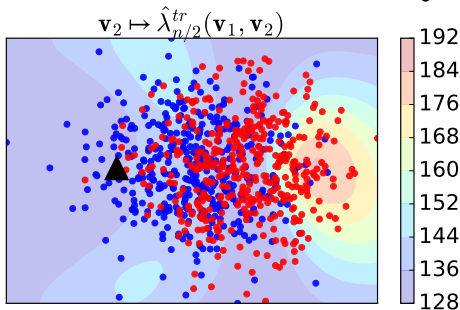
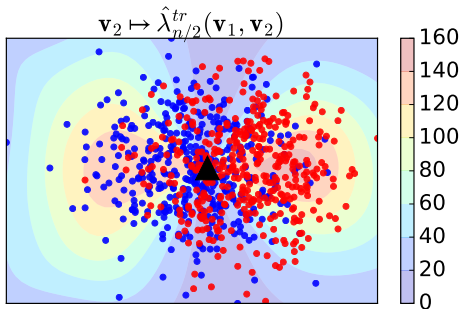
- Non-convexity, informative features.
- Number of locations ( $J$ ).
- MMD: IPM representation.
- Estimation of  $\text{MMD}^2$ .
- Proof idea.
- Computational complexity:  $(J, n, d)$ -dependence.

# Non-convexity, informative features

- 2D problem:

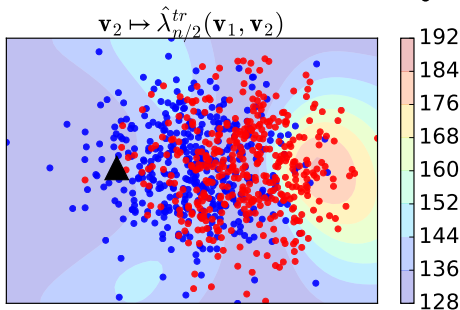
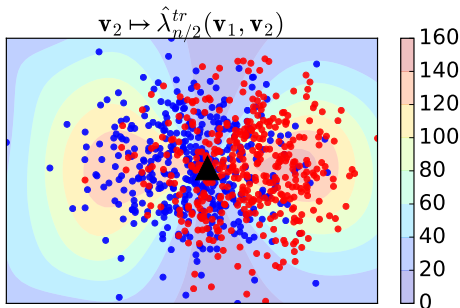
$$\mathbb{P} := \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \mathbb{Q} := \mathcal{N}(\mathbf{e}_1, \mathbf{I}).$$

- $\mathcal{V} = \{\mathbf{v}_1, \mathbf{v}_2\}$ . Fix  $\mathbf{v}_1$  to  $\blacktriangle$ .
- $\mathbf{v}_2 \mapsto \hat{\lambda}_n(\{\mathbf{v}_1, \mathbf{v}_2\})$ : contour plot.



# Non-convexity, informative features

- **Nearby locations:** do not increase discriminability.
- **Non-convexity:** reveals multiple ways to capture the difference.



# 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.



# Number of locations ( $J$ )

- **Very large  $J$ :**
  - test power need not increase monotonically in  $J$  (more locations  $\Rightarrow$  statistic can gain in variance).
  - defeats the purpose of a linear-time test.

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

$$MMD^2(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}(k)}^2 = \left[ \sup_{\|f\|_{\mathcal{H}(k)} \leq 1} \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, f \rangle_{\mathcal{H}(k)} \right]^2$$

# MMD: IPM representation

$$\begin{aligned} \text{MMD}^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}(k)}^2 = \left[ \sup_{\|f\|_{\mathcal{H}(k)} \leq 1} \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, f \rangle_{\mathcal{H}(k)} \right]^2 \\ &\stackrel{(*)}{=} \left[ \sup_{\|f\|_{\mathcal{H}(k)} \leq 1} \mathbb{E}_{x \sim \mathbb{P}} f(x) - \mathbb{E}_{y \sim \mathbb{Q}} f(y) \right]^2. \end{aligned}$$

$$\begin{aligned} \text{MMD}^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}(k)}^2 = \left[ \sup_{\|f\|_{\mathcal{H}(k)} \leq 1} \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, f \rangle_{\mathcal{H}(k)} \right]^2 \\ &\stackrel{(*)}{=} \left[ \sup_{\|f\|_{\mathcal{H}(k)} \leq 1} \mathbb{E}_{x \sim \mathbb{P}} f(x) - \mathbb{E}_{y \sim \mathbb{Q}} f(y) \right]^2. \end{aligned}$$

(\*) in details:

$$\langle \mu_{\mathbb{P}}, f \rangle_{\mathcal{H}(k)} = \left\langle \int k(\cdot, x) d\mathbb{P}(x), f \right\rangle_{\mathcal{H}(k)}$$

# MMD: IPM representation

$$\begin{aligned} \text{MMD}^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}(k)}^2 = \left[ \sup_{\|f\|_{\mathcal{H}(k)} \leq 1} \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, f \rangle_{\mathcal{H}(k)} \right]^2 \\ &\stackrel{(*)}{=} \left[ \sup_{\|f\|_{\mathcal{H}(k)} \leq 1} \mathbb{E}_{x \sim \mathbb{P}} f(x) - \mathbb{E}_{y \sim \mathbb{Q}} f(y) \right]^2. \end{aligned}$$

(\*) in details:

$$\langle \mu_{\mathbb{P}}, f \rangle_{\mathcal{H}(k)} = \left\langle \int k(\cdot, x) d\mathbb{P}(x), f \right\rangle_{\mathcal{H}(k)} = \int \underbrace{\langle k(\cdot, x), f \rangle_{\mathcal{H}(k)}}_{=f(x)} d\mathbb{P}(x)$$

# MMD: IPM representation

$$\begin{aligned} \text{MMD}^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}(k)}^2 = \left[ \sup_{\|f\|_{\mathcal{H}(k)} \leq 1} \langle \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}, f \rangle_{\mathcal{H}(k)} \right]^2 \\ &\stackrel{(*)}{=} \left[ \sup_{\|f\|_{\mathcal{H}(k)} \leq 1} \mathbb{E}_{x \sim \mathbb{P}} f(x) - \mathbb{E}_{y \sim \mathbb{Q}} f(y) \right]^2. \end{aligned}$$

(\*) in details:

$$\begin{aligned} \langle \mu_{\mathbb{P}}, f \rangle_{\mathcal{H}(k)} &= \left\langle \int k(\cdot, x) d\mathbb{P}(x), f \right\rangle_{\mathcal{H}(k)} = \int \underbrace{\langle k(\cdot, x), f \rangle_{\mathcal{H}(k)}}_{=f(x)} d\mathbb{P}(x) \\ &= \mathbb{E}_{x \sim \mathbb{P}} f(x). \end{aligned}$$

Squared difference between feature means:

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



# Estimation of $MMD^2$

Squared difference between feature means:

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

Unbiased empirical estimate for  $\{x_i\}_{i=1}^n \sim \mathbb{P}$ ,  $\{y_j\}_{j=1}^n \sim \mathbb{Q}$ :

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

① Lower bound on the test power:

- ①  $|\hat{\lambda}_n - \lambda_n| \lesssim \|\bar{\mathbf{z}}_n - \boldsymbol{\mu}\|_2 + \|\boldsymbol{\Sigma}_n - \boldsymbol{\Sigma}\|_F.$
- ② Bound the r.h.s. by Hoeffding inequality  $\Rightarrow P(|\hat{\lambda}_n - \lambda_n| \geq t).$
- ③ By reparameterization:  $P(\hat{\lambda}_n \geq T_\alpha)$  bound.

① Lower bound on the test power:

- ①  $|\hat{\lambda}_n - \lambda_n| \lesssim \|\bar{\mathbf{z}}_n - \boldsymbol{\mu}\|_2 + \|\boldsymbol{\Sigma}_n - \boldsymbol{\Sigma}\|_F.$
- ② Bound the r.h.s. by Hoeffding inequality  $\Rightarrow P(|\hat{\lambda}_n - \lambda_n| \geq t).$
- ③ By reparameterization:  $P(\hat{\lambda}_n \geq T_\alpha)$  bound.

② Uniformly  $\hat{\lambda}_n \approx \lambda_n$ :

- Reduction to bounding  $\sup_{\mathcal{V}, \mathcal{K}} \|\bar{\mathbf{z}}_n - \boldsymbol{\mu}\|_2, \sup_{\mathcal{V}, \mathcal{K}} \|\boldsymbol{\Sigma}_n - \boldsymbol{\Sigma}\|_F.$
- Empirical processes, Dudley entropy bound.

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



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)*.



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

The Karolinska directed emotional faces-KDEF.

Technical report, ISBN 91-630-7164-9.



