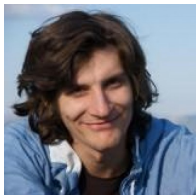


# Big Hypothesis Testing with Kernel Embeddings

Dino Sejdinovic

Department of Statistics  
University of Oxford

9 January 2015  
UCL Workshop on the Theory of Big Data



Heiko Strathmann



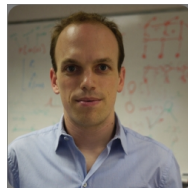
Soumyajit De



Wojciech Zaremba



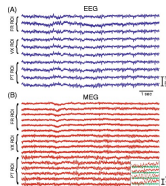
Matthew Blaschko



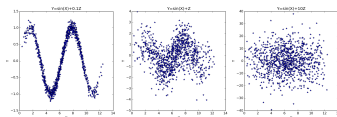
Arthur Gretton

# Making Hard Inference Possible

- many dimensions



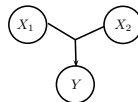
- low signal-to-noise ratio



- highly non-linear associations

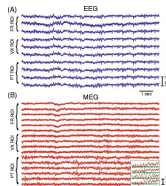


- higher-order interactions

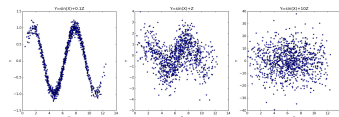


# Making Hard Inference Possible

- many dimensions



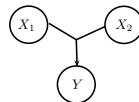
- low signal-to-noise ratio



- highly non-linear associations



- higher-order interactions



need an expressive model **and** a very large number of observations

# Making Hard Inference Possible

- many dimensions



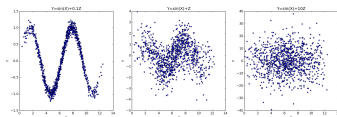
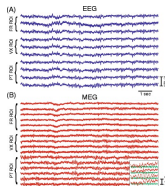
- highly non-linear associations



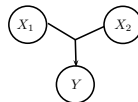
need an expressive model **and** a very large number of observations

cannot use batch algorithms

- low signal-to-noise ratio



- higher-order interactions



# Overview

- 1 Kernel Embeddings and MMD
- 2 Scaling up Kernel Tests
- 3 Experiments

# Outline

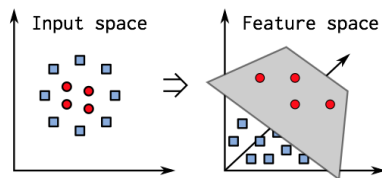
1 Kernel Embeddings and MMD

2 Scaling up Kernel Tests

3 Experiments

# Kernel Embedding

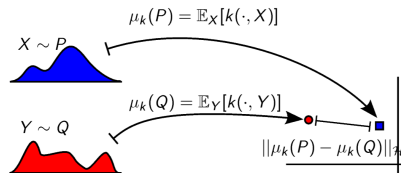
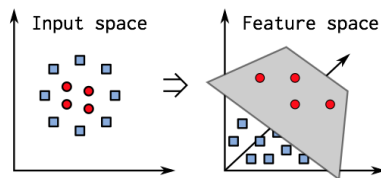
- **feature map:**  $x \mapsto k(\cdot, x) \in \mathcal{H}_k$   
instead of  
 $x \mapsto (\varphi_1(x), \dots, \varphi_s(x)) \in \mathbb{R}^s$
- $\langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}_k} = k(x, y)$   
inner products easily **computed**





# Kernel Embedding

- **feature map:**  $x \mapsto k(\cdot, x) \in \mathcal{H}_k$   
instead of  
 $x \mapsto (\varphi_1(x), \dots, \varphi_s(x)) \in \mathbb{R}^s$
- $\langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}_k} = k(x, y)$   
inner products easily **computed**
- **embedding:**  
 $P \mapsto \mu_k(P) = \mathbb{E}_{X \sim P} k(\cdot, X) \in \mathcal{H}_k$   
instead of  
 $P \mapsto (\mathbb{E} \varphi_1(X), \dots, \mathbb{E} \varphi_s(X)) \in \mathbb{R}^s$
- $\langle \mu_k(P), \mu_k(Q) \rangle_{\mathcal{H}_k} = \mathbb{E}_{X, Y} k(X, Y)$   
inner products easily **estimated**

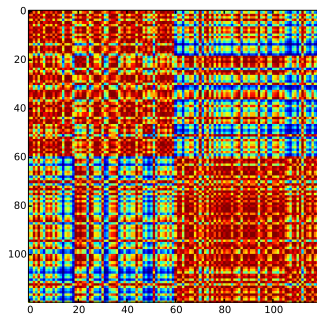
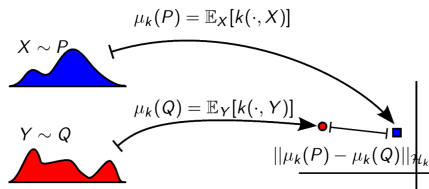


# Kernel MMD (1)

## Definition

**Kernel metric (MMD)** between  $P$  and  $Q$ :

$$\begin{aligned} \text{MMD}_k^2(P, Q) &= \|\mathbb{E}k(\cdot, X) - \mathbb{E}k(\cdot, Y)\|_{\mathcal{H}_k}^2 \\ &= \mathbb{E}_{XX'}k(X, X') + \mathbb{E}_{YY'}k(Y, Y') - 2\mathbb{E}_{XY}k(X, Y) \end{aligned}$$



## Kernel MMD (2)

- A polynomial kernel  $k(z, z') = (1 + z^\top z')^s$  captures the difference in first  $s$  moments only
- For a certain family of kernels (**characteristic**):  $\text{MMD}_k(P, Q) = 0$  if and only if  $P = Q$ : Gaussian  $\exp(-\frac{1}{2\sigma^2} \|z - z'\|_2^2)$ , Laplacian, inverse multiquadratics,  $B_{2n+1}$ -splines...
- Under mild assumptions,  $k$ -MMD metrizes weak\* topology on probability measures (Sriperumbudur, 2010):

$$\text{MMD}_k(P_n, P) \rightarrow 0 \Leftrightarrow P_n \rightsquigarrow P$$

# Nonparametric two-sample tests

- Testing  $\mathbf{H}_0 : \mathbf{P} = \mathbf{Q}$  vs.  $\mathbf{H}_A : \mathbf{P} \neq \mathbf{Q}$   
based on samples  $\{x_i\}_{i=1}^{n_x} \sim \mathbf{P}$ ,  $\{y_i\}_{i=1}^{n_y} \sim \mathbf{Q}$ .
- Test statistic is an estimate of  
 $\text{MMD}_k^2(\mathbf{P}, \mathbf{Q}) = \mathbb{E}_{\mathbf{X}\mathbf{X}'} k(\mathbf{X}, \mathbf{X}') + \mathbb{E}_{\mathbf{Y}\mathbf{Y}'} k(\mathbf{Y}, \mathbf{Y}') - 2\mathbb{E}_{\mathbf{X}\mathbf{Y}} k(\mathbf{X}, \mathbf{Y})$ :

$$\widehat{\text{MMD}} = \frac{1}{n_x(n_x - 1)} \sum_{i \neq j} k(x_i, x_j) + \frac{1}{n_y(n_y - 1)} \sum_{i \neq j} k(y_i, y_j) - \frac{2}{n_x n_y} \sum_{i,j} k(x_i, y_j).$$

- $O(n^2)$  to compute ( $n = n_x + n_y$ )
- Degenerate U-statistic:  $\frac{1}{\sqrt{n}}$ -convergence to MMD under  $\mathbf{H}_A$ ,  
 $\frac{1}{n}$ -convergence to 0 under  $\mathbf{H}_0$ .

# Nonparametric independence tests

- $H_0 : X \perp\!\!\!\perp Y$
- $H_A : X \not\perp\!\!\!\perp Y$

# Nonparametric independence tests

- $H_0 : X \perp\!\!\!\perp Y \Leftrightarrow \mathbf{P}_{XY} = \mathbf{P}_X \mathbf{P}_Y$
- $H_A : X \not\perp\!\!\!\perp Y \Leftrightarrow \mathbf{P}_{XY} \neq \mathbf{P}_X \mathbf{P}_Y$

- Test statistic:

$$\text{HSIC}(X, Y) = \left\| \mu_{\kappa}(\hat{P}_{XY}) - \mu_{\kappa}(\hat{P}_X \hat{P}_Y) \right\|_{\mathcal{H}_{\kappa}}^2,$$

with  $\kappa = k \otimes l$

Gretton et al (2005, 2008); Smola et al (2007)

$$\begin{array}{c} k(\boxed{1}, \boxed{2}) \quad l(\boxed{1}, \boxed{2}) \\ \downarrow \\ \kappa(\boxed{1}, \boxed{1}, \boxed{2}, \boxed{2}) = \\ k(\boxed{1}, \boxed{2}) \times l(\boxed{1}, \boxed{2}) \end{array}$$

# Nonparametric independence tests

- $H_0 : X \perp\!\!\!\perp Y \Leftrightarrow \mathbf{P}_{XY} = \mathbf{P}_X \mathbf{P}_Y$
- $H_A : X \not\perp\!\!\!\perp Y \Leftrightarrow \mathbf{P}_{XY} \neq \mathbf{P}_X \mathbf{P}_Y$

- Test statistic:

$$\text{HSIC}(X, Y) = \left\| \mu_{\kappa}(\hat{P}_{XY}) - \mu_{\kappa}(\hat{P}_X \hat{P}_Y) \right\|_{\mathcal{H}_{\kappa}}^2,$$

with  $\kappa = k \otimes l$

Gretton et al (2005, 2008); Smola et al (2007)

$$k(\boxed{1}, \boxed{2}) \quad l(\boxed{1}, \boxed{2})$$

↓

$$\kappa(\boxed{1}, \boxed{1}, \boxed{2}, \boxed{2}) = k(\boxed{1}, \boxed{2}) \times l(\boxed{1}, \boxed{2})$$

Extensions: conditional independence testing (Fukumizu, Gretton, Sun and Schölkopf, 2008; Zhang, Peters, Janzing and Schölkopf, 2011), three-variable interaction (DS, Gretton and Bergsma, 2013)

# Outline

- 1 Kernel Embeddings and MMD
- 2 Scaling up Kernel Tests
- 3 Experiments

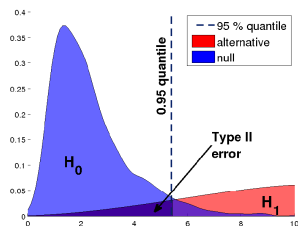


# Test threshold

- under  $H_0$  :  $P = Q$ :

$$\frac{n_x n_y}{n_x + n_y} \widehat{\text{MMD}}_k \rightsquigarrow \sum_{r=1}^{\infty} \lambda_r (Z_r^2 - 1), \quad \{Z_r\} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$$

- $\{\lambda_r\}$  depend on both  $k$  and  $P$ .
- expensive threshold computation:
  - Estimate leading  $\lambda_r$ 's (requires eigendecomposition of the kernel matrix):  $O(n^3)$
  - Permutation test:  $\# \text{shuffles} \times O(n^2)$



## Limited data, unlimited time

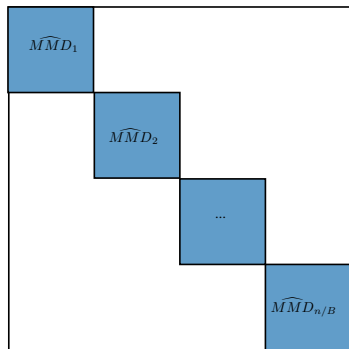
$$\text{MMD}_k^2(P, Q) = \mathbb{E}_{\mathbf{x}\mathbf{x}'} k(\mathbf{x}, \mathbf{x}') + \mathbb{E}_{\mathbf{y}\mathbf{y}'} k(\mathbf{y}, \mathbf{y}') - 2\mathbb{E}_{\mathbf{x}\mathbf{y}} k(\mathbf{x}, \mathbf{y})$$

- Estimate with

$$\widehat{\text{MMD}} = \frac{1}{n_x(n_x - 1)} \sum_{i \neq j} k(x_i, x_j) + \frac{1}{n_y(n_y - 1)} \sum_{i \neq j} k(y_i, y_j) - \frac{2}{n_x n_y} \sum_{i,j} k(x_i, y_j).$$

- Complexity:  $O(n^2)$ .

# Limited time, unlimited data



- Process mini-batches of size  $B$  at a time:  

$$\hat{\eta}_k = \frac{B}{n} \sum_{b=1}^{n/B} \widehat{MMD}_{k,b}$$
- Complexity:  $O(nB)$ .
- Provided  $B/n \rightarrow 0$ :  
 $\frac{1}{\sqrt{n}}$ -convergence to MMD if  $\text{MMD} \neq 0$ ,  $\frac{1}{\sqrt{nB}}$ -convergence to 0 under  $\mathbf{H}_0$ .

- A.Gretton, B.Sriperumbudur, D.S., H.Strathmann, S.Balakrishnan, M.Pontil and K.Fukumizu, **Optimal kernel choice for large-scale two-sample tests**, *NIPS* 2012.  
 - W.Zaremba, A.Gretton, M.Blaschko, **B-test: A Non-Parametric, Low Variance Kernel Two-Sample Test**, *NIPS* 2013.

## Full statistic vs. mini-batch statistic

	$U$ -statistic	mini-batch
time	$O(n^2)$	$O(nB)$
storage	$O(n^2)$	$O(B^2)$
null distribution	infinite sum of chi-squares	normal
computing p-value	$O(n^3)$ or #shuffles $\times O(n^2)$	$O(nB)$
convergence rate	$1/n$	$1/\sqrt{nB}$

- $\frac{n_x n_y}{(n_x + n_y)^{3/2}} \sqrt{B} \hat{\eta}_k \rightsquigarrow \mathcal{N}(0, \sigma_k^2)$  under  $\mathbf{H}_0$
- $\sigma_k^2$  (depends on  $k$  and  $\mathbf{P}$ ) can be unbiasedly estimated on each block in  $O(B^2)$  time

# Asymptotic efficiency criterion

- A. Gretton, B. Sriperumbudur, DS, H. Strathmann, S. Balakrishnan, M. Pontil and K. Fukumizu, **Optimal kernel choice for large-scale two-sample tests**, *NIPS* 2012.

## Proposition

*For given  $P$  and  $Q$ . Let  $\eta_k = \text{MMD}_k^2(P, Q)$ , and let  $\sigma_k^2$  be the asymptotic variance of the linear-time statistic  $\hat{\eta}_k$ . Then*

$$k_* = \arg \max_{k \in \mathcal{K}} \eta_k / \sigma_k$$

*minimizes the asymptotic Type II error probability on  $\mathcal{K}$ .*

# Asymptotic efficiency criterion

- A. Gretton, B. Sriperumbudur, D.S. H. Strathmann, S. Balakrishnan, M. Pontil and K. Fukumizu, **Optimal kernel choice for large-scale two-sample tests**, *NIPS* 2012.

## Proposition

*For given  $P$  and  $Q$ . Let  $\eta_k = \text{MMD}_k^2(P, Q)$ , and let  $\sigma_k^2$  be the asymptotic variance of the linear-time statistic  $\hat{\eta}_k$ . Then*

$$k_* = \arg \max_{k \in \mathcal{K}} \eta_k / \sigma_k$$

*minimizes the asymptotic Type II error probability on  $\mathcal{K}$ .*

- We only have estimates of  $\eta_k$  and  $\sigma_k$  !
- Will the kernel optimization using plug-in estimates be consistent?
- Over what families of kernels can we perform such optimization *efficiently*?

# Asymptotic efficiency criterion

- A. Gretton, B. Sriperumbudur, D.S., H. Strathmann, S. Balakrishnan, M. Pontil and K. Fukumizu, **Optimal kernel choice for large-scale two-sample tests**, *NIPS* 2012.

## Proposition

*For given  $P$  and  $Q$ . Let  $\eta_k = \text{MMD}_k^2(P, Q)$ , and let  $\sigma_k^2$  be the asymptotic variance of the linear-time statistic  $\hat{\eta}_k$ . Then*

$$k_* = \arg \max_{k \in \mathcal{K}} \eta_k / \sigma_k$$

*minimizes the asymptotic Type II error probability on  $\mathcal{K}$ .*

- We only have estimates of  $\eta_k$  and  $\sigma_k$  !
- Will the kernel optimization using plug-in estimates be consistent? **yes!**
- Over what families of kernels can we perform such optimization *efficiently*? **linear combinations (MKL)**

# Outline

- 1 Kernel Embeddings and MMD
- 2 Scaling up Kernel Tests
- 3 Experiments



# Hard-to-detect differences: Gaussian blobs

**Difficult problems:** lengthscale of the *difference* in distributions not the same as that of the distributions.

# Hard-to-detect differences: Gaussian blobs

**Difficult problems:** lengthscale of the *difference* in distributions not the same as that of the distributions.

We distinguish grids of Gaussian blobs with different covariances.

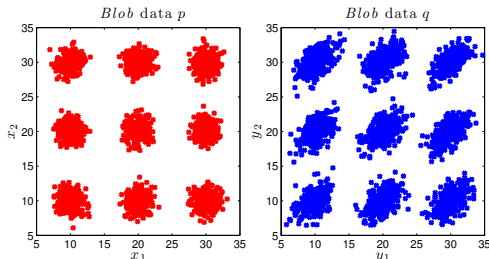


Figure :  $3 \times 3$  blobs, ratio  $\varepsilon = 3.2$  of largest-to-smallest eigenvalues of blobs in  $Q$ .

## Gaussian blobs (2)

$12 \times 12$  blobs with  $\varepsilon = 1.4$ . Linear time statistic vs. Quadratic time statistic. Fixed kernel.

## Gaussian blobs (2)

$12 \times 12$  blobs with  $\varepsilon = 1.4$ . Linear time statistic vs. Quadratic time statistic. Fixed kernel.

	<i>m</i> per trial	Type II error	Trials
Quadratic	5,000	[0.7996, 0.8516]	820
	10,000	[0.5161, 0.6175]	367
	> 10,000	Buy more RAM!	

## Gaussian blobs (2)

$12 \times 12$  blobs with  $\varepsilon = 1.4$ . Linear time statistic vs. Quadratic time statistic. Fixed kernel.

	<i>m</i> per trial	Type II error	Trials
Quadratic	5,000	[0.7996, 0.8516]	820
	10,000	[0.5161, 0.6175]	367
	> 10,000	Buy more RAM!	
Linear	$\sim 100,000,000$	[0.2250, 0.3049]	468
	$\sim 200,000,000$	[0.1873, 0.2829]	302
	$\vdots$	$\vdots$	$\vdots$
	$\sim 500,000,000$	<b><math>0.0270 \pm 0.0302</math></b>	111

## Gaussian blobs (3)

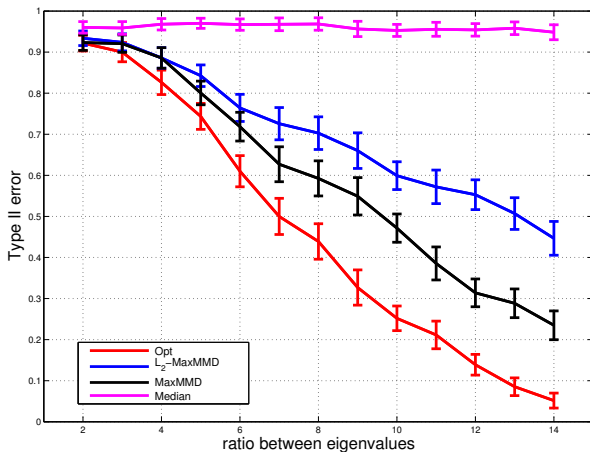


Figure :  $m = 10,000$ ; family generated by gaussian kernels with bandwidths  $\{2^{-5}, \dots, 2^{15}\}$ .

# Hard-to-detect differences: UCI HIGGS

- P. Baldi, P. Sadowski, and D. Whiteson. **Searching for Exotic Particles in High-energy Physics with Deep Learning**. *Nature Communications* 5, 2014.

- benchmark dataset for distinguishing a signature of Higgs boson vs. background
- joint distributions of the azimuthal angular momenta  $\varphi$  for four particle jets: low-signal, low-level features
- Do *joint* angular momenta carry any discriminating information?

# Hard-to-detect differences: UCI HIGGS

- P. Baldi, P. Sadowski, and D. Whiteson. **Searching for Exotic Particles in High-energy Physics with Deep Learning**. *Nature Communications* 5, 2014.

- benchmark dataset for distinguishing a signature of Higgs boson vs. background
- joint distributions of the azimuthal angular momenta  $\varphi$  for four particle jets: low-signal, low-level features
- Do *joint* angular momenta carry any discriminating information?

sample size:	1e4	5e4	1e5	5e5	1e6
p-value (gauss-med):	.757	.217	.475	.391	.074



# Hard-to-detect differences: UCI HIGGS

- P. Baldi, P. Sadowski, and D. Whiteson. **Searching for Exotic Particles in High-energy Physics with Deep Learning**. *Nature Communications* 5, 2014.

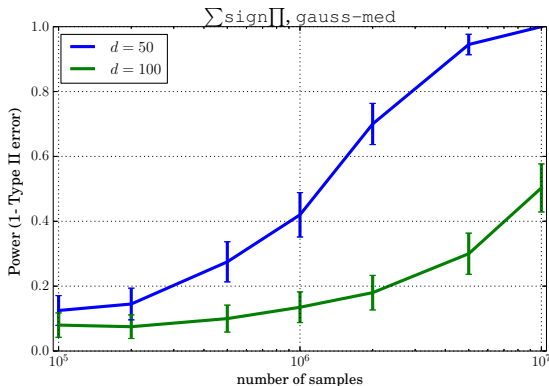
- benchmark dataset for distinguishing a signature of Higgs boson vs. background
- joint distributions of the azimuthal angular momenta  $\varphi$  for four particle jets: low-signal, low-level features
- Do *joint* angular momenta carry any discriminating information?

sample size:	1e4	5e4	1e5	5e5	1e6
p-value (gauss-med):	.757	.217	.475	.391	.074

train/test size:	2e3/8e3	1e4/4e4	2e4/8e4	1e5/4e5	2e5/8e5
p-value (gauss-opt):	.139	.476	.035	<b>6.12e-5</b>	<b>1.02e-18</b>

# Experiment: Independence Test ( $\sum \text{sign} \Pi$ )

- $X \sim \mathcal{N}(0, I_d)$ ,  
 $Y = \sqrt{\frac{2}{d}} \sum_{j=1}^{d/2} \text{sign}(X_{2j-1} X_{2j}) |Z_j| + Z_{\frac{d}{2}+1}$ , where  $Z \sim \mathcal{N}(0, I_{\frac{d}{2}+1})$



# Experiment: Independence Test (sine $\sum$ )

- $X_1, X_2 \stackrel{i.i.d.}{\sim} \text{Unif}[0, 2\pi]$ ,  
 $Y = \sin(X_1 + X_2) + 10Z$ , with  $Z \sim \mathcal{N}(0, 1)$ .

sine $\sum$ , $B = 100$	brown: $q = 1$	brown: opt	gauss: med	gauss: opt
$N = 5e5$ , 1-Type II	.277 $\pm$ .059	<b>.675 <math>\pm</math> .065</b>	.190 $\pm$ .054	<b>.740 <math>\pm</math> .061</b>
Type I	.035 $\pm$ .025	.025 $\pm$ .022	.085 $\pm$ .039	.040 $\pm$ .027
$N = 1e6$ , 1-Type II	.460 $\pm$ .069	<b>.915 <math>\pm</math> .039</b>	.325 $\pm$ .065	<b>.905 <math>\pm</math> .041</b>
Type I	.055 $\pm$ .032	.050 $\pm$ .030	.025 $\pm$ .022	.060 $\pm$ .033

# Shogun

Home About Documentation Contact Blog News Events

將軍 The Shogun Machine Learning Toolbox

**SHOGUN 3.2.0**  
In \$DEITY we trust all others bring data.

DOWNLOAD NOW

What's New: SHOGUN Release version 3.2.0 (libshogun 16.0, data 0.8, parameter 1)

Follow us on Twitter

WORKSHOP VIDEOS! The recordings of the presentations from the SHOGUN machine learning workshop 2014 are now available.

c-base

**SHOGUN FEATURES**

A large scale machine learning toolbox. SHOGUN is designed for unified large-scale learning for a broad range of feature types and learning settings, like classification, regression, or explorative data analysis.

**SHOGUN TALKS**

SHOGUN was presented in July at the [EuroPython conference](#) in Berlin and in August at the [Open Machine Learning Workshop](#) in New York, both times by Heiko Strathmann. These talks are of interest both to people that do not know about SHOGUN and are interested in using it or just getting an introduction, as well as for people with prior experience using SHOGUN.

**WHAT'S NEW**

Feb. 17, 2014 → [SHOGUN 3.2.0](#)  
Jan. 6, 2014 → [SHOGUN 3.1.1](#)  
Jan. 8, 2014 → [SHOGUN 3.1.0](#)  
Oct. 28, 2013 → [SHOGUN 3.0.0](#)  
March 17, 2013 → [SHOGUN 2.1.0](#)  
Sept. 1, 2012 → [SHOGUN 2.0.0](#)  
Dec. 1, 2011 → [SHOGUN 1.1.0](#)

- Written in C++ with interfaces to Python, Matlab, Java, R.
- Google Summer of Code (2012, 2014).

# Summary

- Hypothesis testing based on kernel embeddings reveals hard-to-detect differences between distributions and non-linear low-signal associations.

# Summary

- Hypothesis testing based on kernel embeddings reveals hard-to-detect differences between distributions and non-linear low-signal associations.
- A simple mini-batch procedure allows us to run the tests on large-scale problems and on streaming data.

# Summary

- Hypothesis testing based on kernel embeddings reveals hard-to-detect differences between distributions and non-linear low-signal associations.
- A simple mini-batch procedure allows us to run the tests on large-scale problems and on streaming data.
- Can select kernel parameters on-the-fly in order to explicitly maximise test power.

# Summary

- Hypothesis testing based on kernel embeddings reveals hard-to-detect differences between distributions and non-linear low-signal associations.
- A simple mini-batch procedure allows us to run the tests on large-scale problems and on streaming data.
- Can select kernel parameters on-the-fly in order to explicitly maximise test power.
- Both kernel selection and testing in  $O(n)$  time and  $O(1)$  storage (if  $B = \text{const}$ ).