

# Change detection in multi-dimensional datasets and time series

**Andrea De Simone**

[andrea.desimone@sissa.it](mailto:andrea.desimone@sissa.it)



Univ. Camerino, 2019-02-26

[DS, Jacques – arXiv:1807.06038]

# > Outline

- 1 Two-Sample Test: Intro & Motivation
- 2 Nearest Neighbors Two-Sample Test (NN2ST)
- 3 Gaussian Examples
- 4 Outlook: Time Series Data

## Two-Sample Test

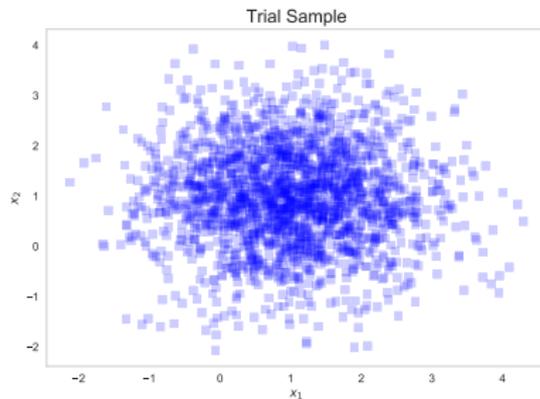
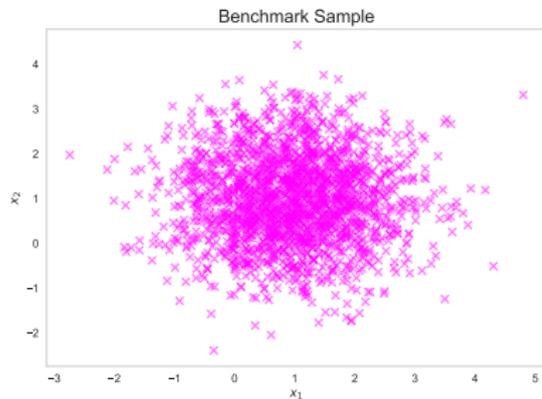
Two sets:

$$\text{Trial: } \mathcal{T} \equiv \{\mathbf{x}_1, \dots, \mathbf{x}_{N_T}\} \stackrel{\text{iid}}{\sim} p_T,$$

$$\text{Benchmark: } \mathcal{B} \equiv \{\mathbf{x}'_1, \dots, \mathbf{x}'_{N_B}\} \stackrel{\text{iid}}{\sim} p_B.$$

$$\mathbf{x}_i, \mathbf{x}'_i \in \mathbb{R}^D$$

$p_B, p_T$  unknown



## Two-Sample Test

Two sets:

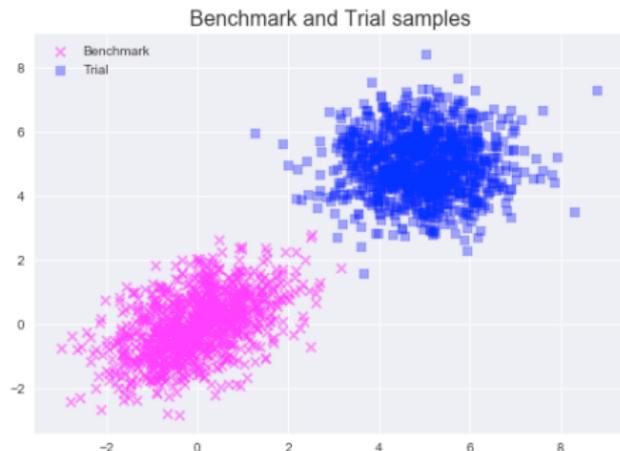
$$\text{Trial: } \mathcal{T} \equiv \{\mathbf{x}_1, \dots, \mathbf{x}_{N_T}\} \stackrel{\text{iid}}{\sim} p_T,$$

$$\text{Benchmark: } \mathcal{B} \equiv \{\mathbf{x}'_1, \dots, \mathbf{x}'_{N_B}\} \stackrel{\text{iid}}{\sim} p_B.$$

$$\mathbf{x}_i, \mathbf{x}'_i \in \mathbb{R}^D$$

$p_B, p_T$  unknown

« Are  $\mathcal{B}, \mathcal{T}$  drawn from the same probability distribution? »



easy...

## Two-Sample Test

Two sets:

$$\text{Trial: } \mathcal{T} \equiv \{\mathbf{x}_1, \dots, \mathbf{x}_{N_T}\} \stackrel{\text{iid}}{\sim} p_T,$$

$$\text{Benchmark: } \mathcal{B} \equiv \{\mathbf{x}'_1, \dots, \mathbf{x}'_{N_B}\} \stackrel{\text{iid}}{\sim} p_B.$$

$$\mathbf{x}_i, \mathbf{x}'_i \in \mathbb{R}^D$$

$p_B, p_T$  unknown

« Are  $\mathcal{B}, \mathcal{T}$  drawn from the same probability distribution? »



... hard

## > Two-Sample Test

### Why is it important?

- detect departures from benchmark
- find anomalous points (outliers)
- check if observed data are compatible with expectations
- detect changes in underlying distributions
- real-time detect events/shifts in time series

## > Two-Sample Test

### *Desiderata for a statistical test*

(1) **model-independent**

no assumption about underlying physical model to interpret data

→ more general

(2) **non-parametric**

compare two samples as a whole (not just their means, etc.)

→ fewer assumptions, no max likelihood estim.

(3) **un-binned**

high-dim feature space partitioned without rectangular bins

→ retain full multi-dim info of data

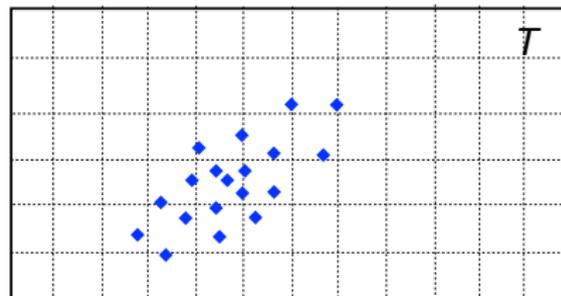
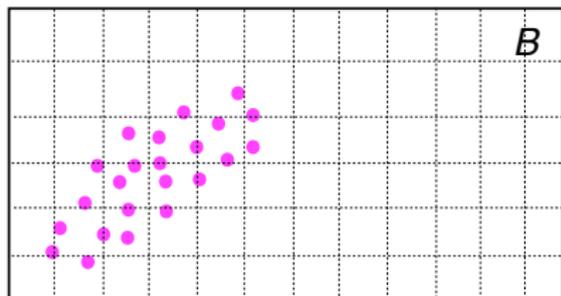
# Two-Sample Test

## Recipe

- (1) **Density Estimator**  
→ reconstruct PDF from samples
- (2) **Test Statistic (TS)**  
→ “measure distance” between PDFs
- (3) **TS distribution**  
→ associate probabilities to TS  
under null hypothesis  $H_0 : p_B = p_T$
- (4)  **$p$ -value**  
→ if  $p < \alpha$  then reject  $H_0$

Let's build the **Nearest Neighbors Two-Sample Test (NN2ST)**

# > 1. Density Estimator



Divide space in square bins?

- ✓ easy
- ✓ can use simple statistics (e.g.  $\chi^2$ )
- ✗ hard/slow/impossible in high- $D$

**Need un-binned,  
multi-variate approach**

Find PDF *estimators*  $\hat{p}_B, \hat{p}_T$ ,  
e.g. based on density of points

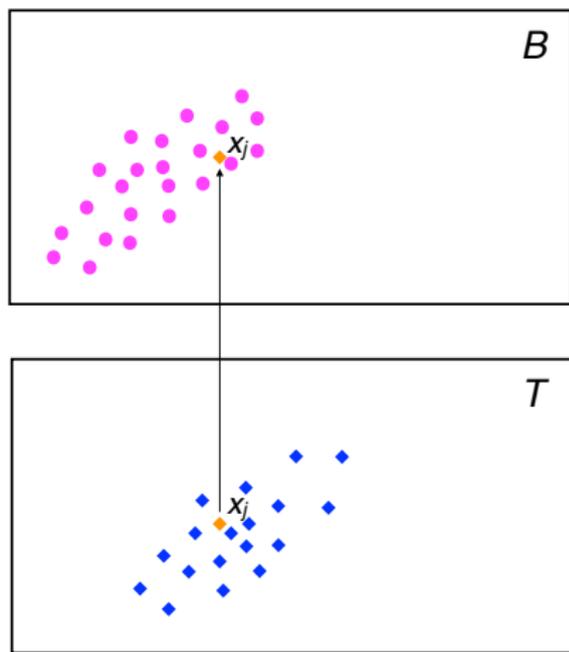
$$\hat{p}_{B,T}(\mathbf{x}) = \frac{\rho_{B,T}(\mathbf{x})}{N_{B,T}}$$

**Nearest Neighbors!**

[Schilling 1986, Henze 1988]

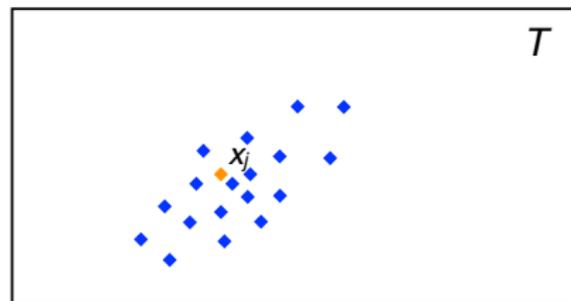
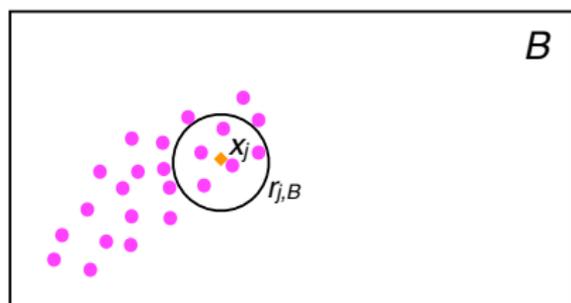
[Wang et al. 2005-2006, Perez-Cruz. 2008]

## > 1. Density Estimator



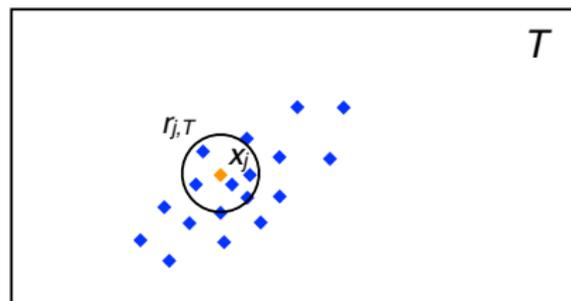
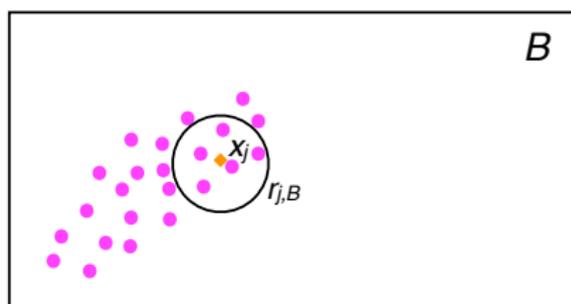
- Fix integer  $K$ .
- Choose query point  $x_j$  in  $\mathcal{T}$  and draw it in  $\mathcal{B}$ .

## 1. Density Estimator



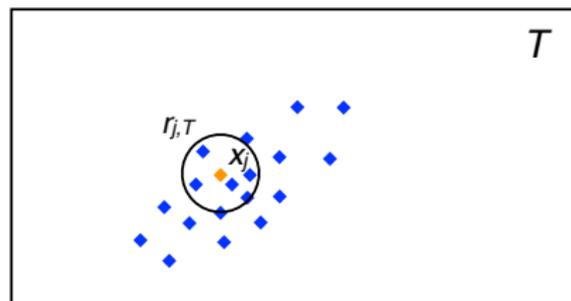
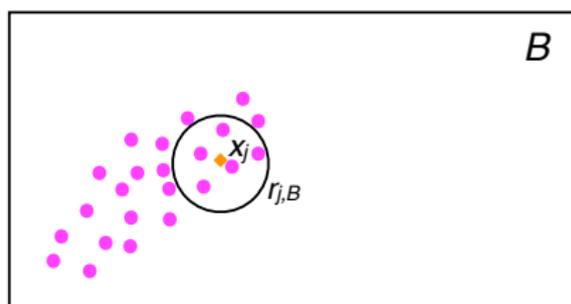
- Fix integer  $K$ .
- Choose query point  $x_j$  in  $\mathcal{T}$  and draw it in  $\mathcal{B}$ .
- Find the distance  $r_{j,B}$  of the  $K^{\text{th}}$ -NN of  $x_j$  in  $\mathcal{B}$ .

## > 1. Density Estimator



- Fix integer  $K$ .
- Choose query point  $x_j$  in  $\mathcal{T}$  and draw it in  $\mathcal{B}$ .
- Find the distance  $r_{j,B}$  of the  $K^{\text{th}}$ -NN of  $x_j$  in  $\mathcal{B}$ .
- Find the distance  $r_{j,T}$  of the  $K^{\text{th}}$ -NN of  $x_j$  in  $\mathcal{T}$ .

## > 1. Density Estimator



- Fix integer  $K$ .
- Choose query point  $\mathbf{x}_j$  in  $\mathcal{T}$  and draw it in  $\mathcal{B}$ .
- Find the distance  $r_{j,B}$  of the  $K^{\text{th}}$ -NN of  $\mathbf{x}_j$  in  $\mathcal{B}$ .
- Find the distance  $r_{j,T}$  of the  $K^{\text{th}}$ -NN of  $\mathbf{x}_j$  in  $\mathcal{T}$ .
- Estimate PDFs:

$$\hat{p}_B(\mathbf{x}_j) = \frac{K}{N_B} \frac{1}{\omega_D r_{j,B}^D}$$

$$\hat{p}_T(\mathbf{x}_j) = \frac{K}{N_T - 1} \frac{1}{\omega_D r_{j,T}^D}$$

## 2. Test Statistic

- Measure the “distance” between 2 PDFs
- Define **Test Statistic** (to detect under-/over-densities)

$$\text{TS}(\mathcal{T}) \equiv \frac{1}{N_T} \sum_{j=1}^{N_T} \log \frac{\hat{p}_T(\mathbf{x}_j)}{\hat{p}_B(\mathbf{x}_j)}$$

- Form NN-estimated PDFs:

$$\text{TS}(\mathcal{T}) = \frac{D}{N_T} \sum_{j=1}^{N_T} \log \frac{r_{j,B}}{r_{j,T}} + \log \frac{N_B}{N_T - 1}$$

- Related to *Kullback-Leibler* divergence as:  $\text{TS}(\mathcal{T}) = \hat{D}_{\text{KL}}(\hat{p}_T || \hat{p}_B)$   
 $[D_{\text{KL}}(p||q) \equiv \int_{\mathbb{R}^D} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}]$

- **Theorem:**

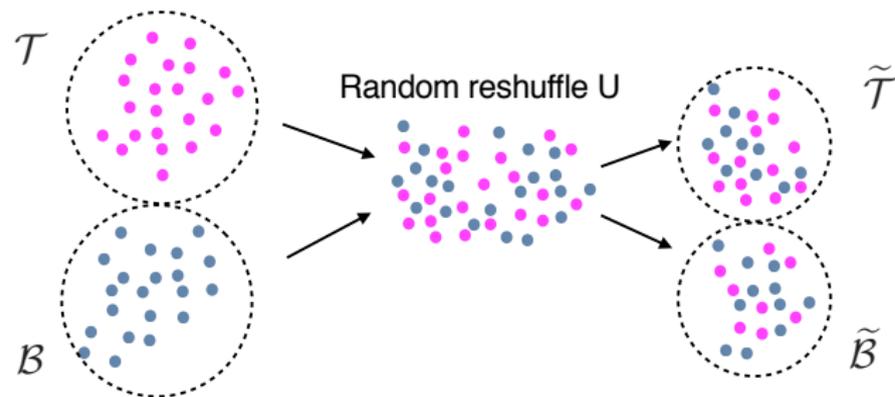
this estimator converges to  $D_{\text{KL}}(p_B || p_T)$ , in the large sample limit

[Wang et al. – 2005, 2006]

### 3. Test Statistic Distribution

How is TS distributed? **Permutation test!**

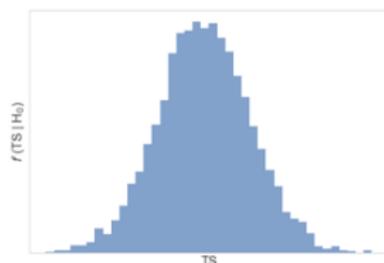
Assume  $p_B = p_T$ . Union set  $\mathcal{U} = \mathcal{T} \cup \mathcal{B}$ .



Compute the test statistic  $\text{TS}_n$  on  $(\tilde{\mathcal{B}}, \tilde{\mathcal{T}})$ .

Repeat many times.

Distribution of TS under  $H_0 : f(\text{TS}|H_0) \leftarrow \{\text{TS}_n\}$   
[asymptotically normal with zero mean]



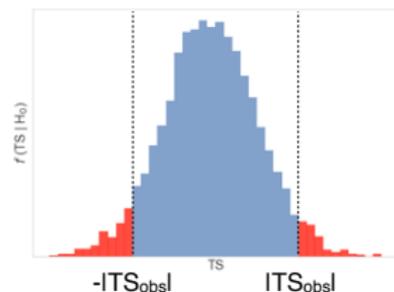
## > 4. $p$ -value

- Find  $\hat{\mu}, \hat{\sigma}$ : mean, variance of  $f(\text{TS}|H_0)$
- Standardize the TS:

$$\text{TS} \rightarrow \text{TS}' \equiv \frac{\text{TS} - \hat{\mu}}{\hat{\sigma}}$$

- $\text{TS}'$  distributed according to  $f'(\text{TS}'|H_0) = \hat{\sigma} f(\hat{\mu} + \hat{\sigma}\text{TS}'|H_0)$
- Two-sided  $p$ -value

$$p = 2 \int_{|\text{TS}_{\text{obs}}|}^{\infty} f'(\text{TS}'|H_0) d\text{TS}'$$



## > NN2ST: Summary

### INPUT:

Trial sample:	$\mathcal{T} \equiv \{\mathbf{x}_1, \dots, \mathbf{x}_{N_T}\} \stackrel{\text{iid}}{\sim} p_T,$	$\mathbf{x}_i, \mathbf{x}'_i \in \mathbb{R}^D$ $p_B, p_T$ unknown
Benchmark sample:	$\mathcal{B} \equiv \{\mathbf{x}'_1, \dots, \mathbf{x}'_{N_B}\} \stackrel{\text{iid}}{\sim} p_B$	
$K$ :	number of nearest neighbors	
$N_{\text{perm}}$ :	number of permutations	

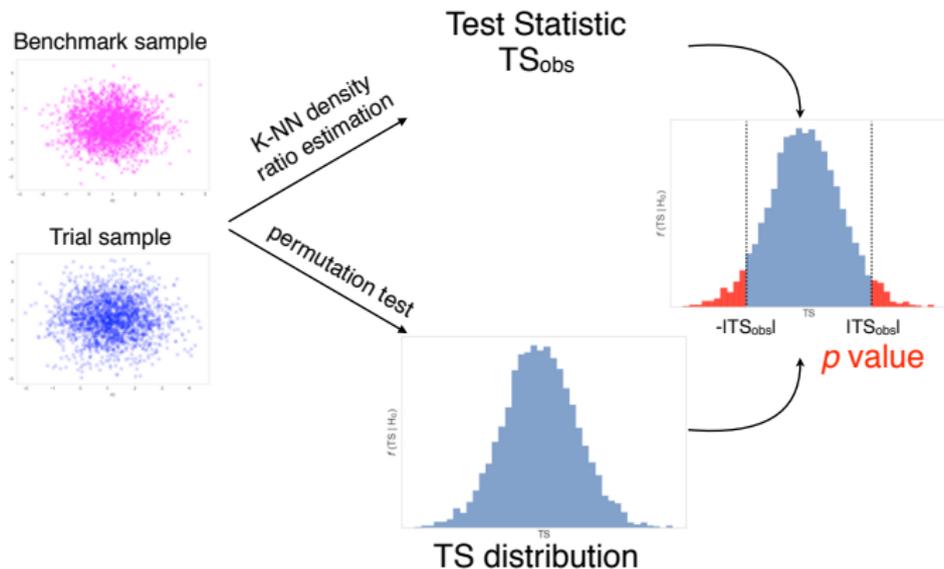
### OUTPUT:

**$p$ -value of the null hypothesis  $H_0 : p_B = p_T$**

[check compatibility between 2 samples]

[detect changes in underlying distributions]

# > NN2ST: Summary



Python code:

[github.com/de-simone/NN2ST](https://github.com/de-simone/NN2ST)

[DS, Jacques – arXiv:1807.06038]

## > NN2ST: Summary

- ✓ general, model-independent
- ✓ solid math foundations
- ✓ fast, no optimization
- ✓ sensitive to unspecified signals
- ✗ need to run for each sample pair
- ✗ permutation test is bottleneck

# > NN2ST on Gaussian Samples

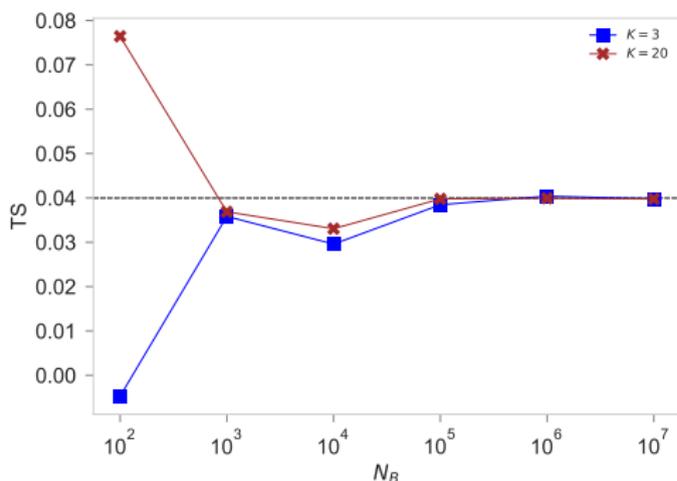
Random samples from  
 $D$ -dimensional Gaussians

$$D = 2,$$

$$p_B = \mathcal{N}(\boldsymbol{\mu}_B, \Sigma_B),$$
$$p_T = \mathcal{N}(\boldsymbol{\mu}_T, \Sigma_T).$$

$$\boldsymbol{\mu}_B = \begin{pmatrix} 1.0 \\ 1.0 \end{pmatrix}, \quad \boldsymbol{\mu}_T = \begin{pmatrix} 1.2 \\ 1.2 \end{pmatrix},$$

$$\Sigma_B = \Sigma_T = \mathbf{I}_2.$$



Convergence to exact  
KL divergence

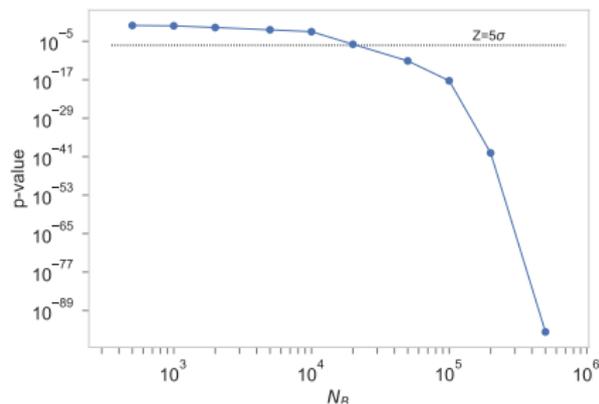
# > NN2ST on Gaussian Samples

Dataset	$\mu$	$\Sigma$
$\mathcal{B}$	$1_D$	$I_D$
$\mathcal{T}_{G0}$	$1_D$	$I_D$
$\mathcal{T}_{G1}$	$1.12_D$	$I_D$
$\mathcal{T}_{G2}$	$1_D$	$\begin{pmatrix} 0.95 & 0.1 & \mathbf{0} \\ 0.1 & 0.8 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & I_{D-2} \end{pmatrix}$
$\mathcal{T}_{G3}$	$1.15_D$	$I_D$

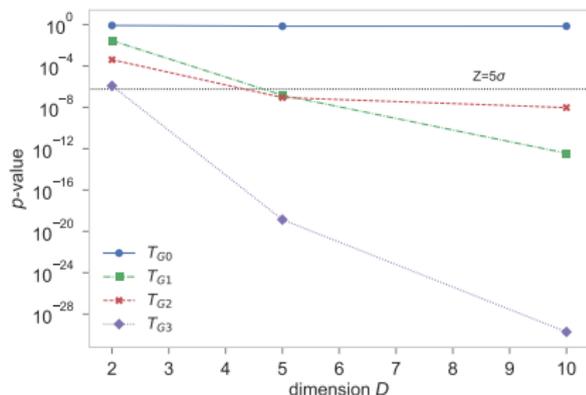
$$N_B = N_T = 20\,000$$

$$K = 5$$

$$N_{\text{perm}} = 1\,000$$



more data, more power



higher  $D$ , more power

## > Outlook: time series data

[Caveat Emptor: very preliminary!]

Real-time detection of changes in data streams:  
variation in underlying mechanism generating data.

$\mathcal{T}, \mathcal{B}$  samples: windows of time series data, ending at discrete times  $t, t'$

$$\begin{aligned}\mathcal{T}_t &= \{\mathbf{x}_{t-N+1}, \dots, \mathbf{x}_t\}, \\ \mathcal{B}_{t'} &= \{\mathbf{x}_{t'-N+1}, \dots, \mathbf{x}_{t'}\}, \quad (N_B = N_T \equiv N).\end{aligned}$$

**Trial** window sliding forward with time.

**Benchmark** window anchored or rolling.

- anchored  $\mathcal{B}$  window:  $t' = N \rightarrow \mathcal{B}_{t'} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$   
Captures cumulative changes over time.
- adjacent windows:  $t' = t - N \rightarrow \mathcal{B}_{t'} = \{\mathbf{x}_{t-2N+1}, \dots, \mathbf{x}_{t-N}\}$   
Captures “rate of change” at current time.

# > Outlook: time series data

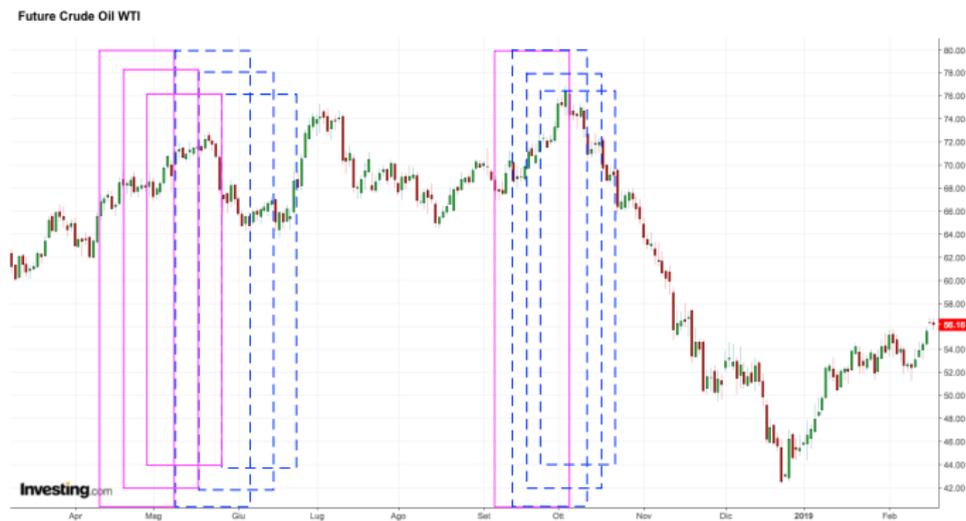


# > Outlook: time series data



adjacent vs. anchored windows

## › Outlook: time series data



- ▶ **Feature space** can be high-dimensional:  
prices (OHLC), prices of related markets, indicators, volumes, ...
- ▶ Reduce false alarms with **persistence factor**  $\gamma$  ( $\sim 1$ )%.

$H_0$  rejected  $\gamma \cdot N$  times in a row

→ **detected change in market conditions**

## > Take-Home Messages

- (1) **Proposed a new statistical test: NN2ST**
- (2) **Model-independent and suitable for high- $D$  data**
- (3) **Excellent results on static datasets**
- (4) **Promising applications for change detection in time series data**