

# GUIDING NEW PHYSICS SEARCHES WITH UNSUPERVISED LEARNING

[DS, Jacques - 1807.06038]

Andrea De Simone



# > New Physics ?

Searches for New Physics Beyond the Standard Model  
have been negative so far...

**MAYBE:**

## **1. New Physics (NP) is not accessible by LHC**

new particles are too light/heavy  
or interacting too weakly

## **2. We have not explored all the possibilities**

new physics may be buried under large bkg  
or hiding behind unusual signatures

## “Don’t want to miss a thing” (in data)

closer look at current data

get ready for upcoming data from next run

## Model-independent search

searches for specific models may be:

- time-consuming
- insensitive to unexpected/unknown processes

# > New Statistical Test

## Want a statistical test for NP which is:

### 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

## 1. Statistical test of dataset compatibility

- Nearest-Neighbors Two-Sample Test
- Identify Discrepancies
- Include Uncertainties

## 2. Applications to High-Energy Physics

## 1. Statistical test of dataset compatibility

- Nearest-Neighbors Two-Sample Test
- Identify Discrepancies
- Include Uncertainties

## 2. Applications to High-Energy Physics

# > Two-sample Test

Two sets:

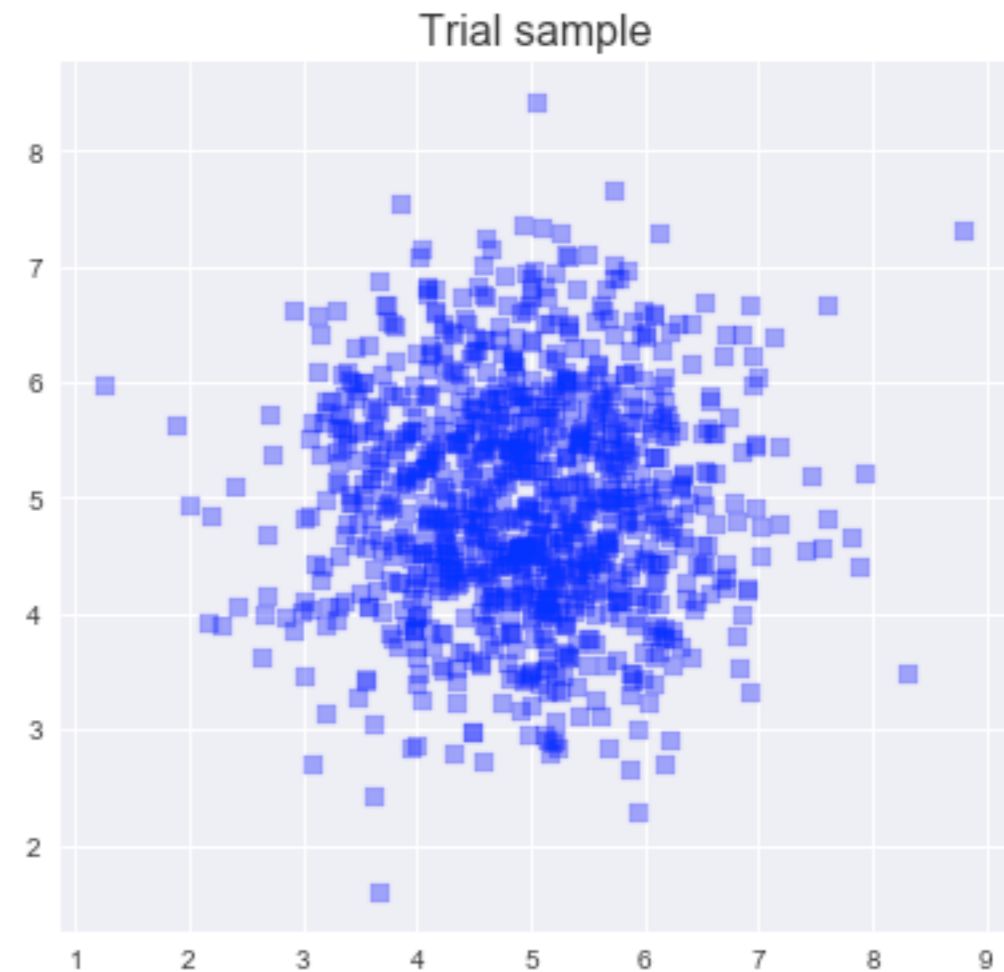
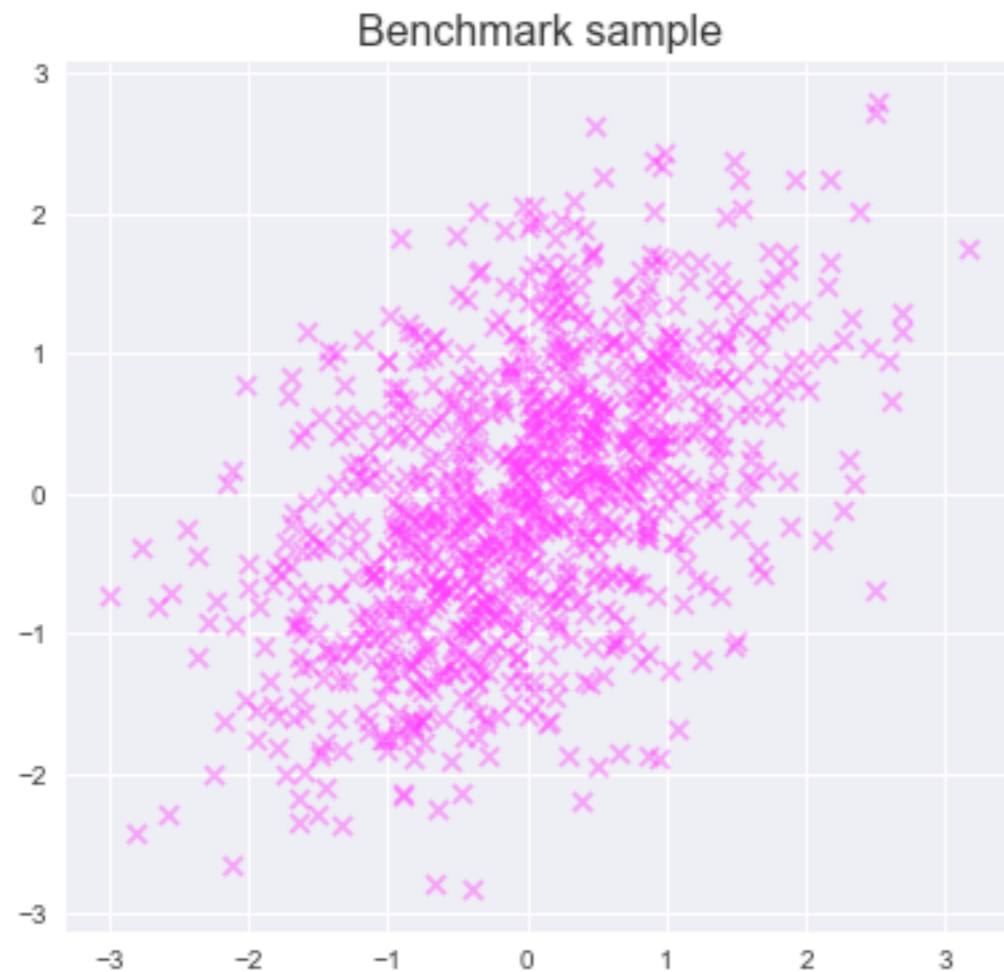
[a.k.a. “homogeneity test”]

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

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

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

probability distributions  $p_B, p_T$  unknown



e.g.: simulated SM bkg

real measured data

# > Two-sample Test

Two sets:

$$\begin{array}{ll} \text{Trial:} & \mathcal{T} = \{\mathbf{x}_1, \dots, \mathbf{x}_{N_T}\} \stackrel{\text{iid}}{\sim} p_T \\ \text{Benchmark:} & \mathcal{B} = \{\mathbf{x}'_1, \dots, \mathbf{x}'_{N_B}\} \stackrel{\text{iid}}{\sim} p_B \end{array} \quad \mathbf{x}_i, \mathbf{x}'_i \in \mathbb{R}^D$$

probability distributions  $p_B, p_T$  unknown

Are  $B, T$  drawn from the same prob. distribution?



easy...



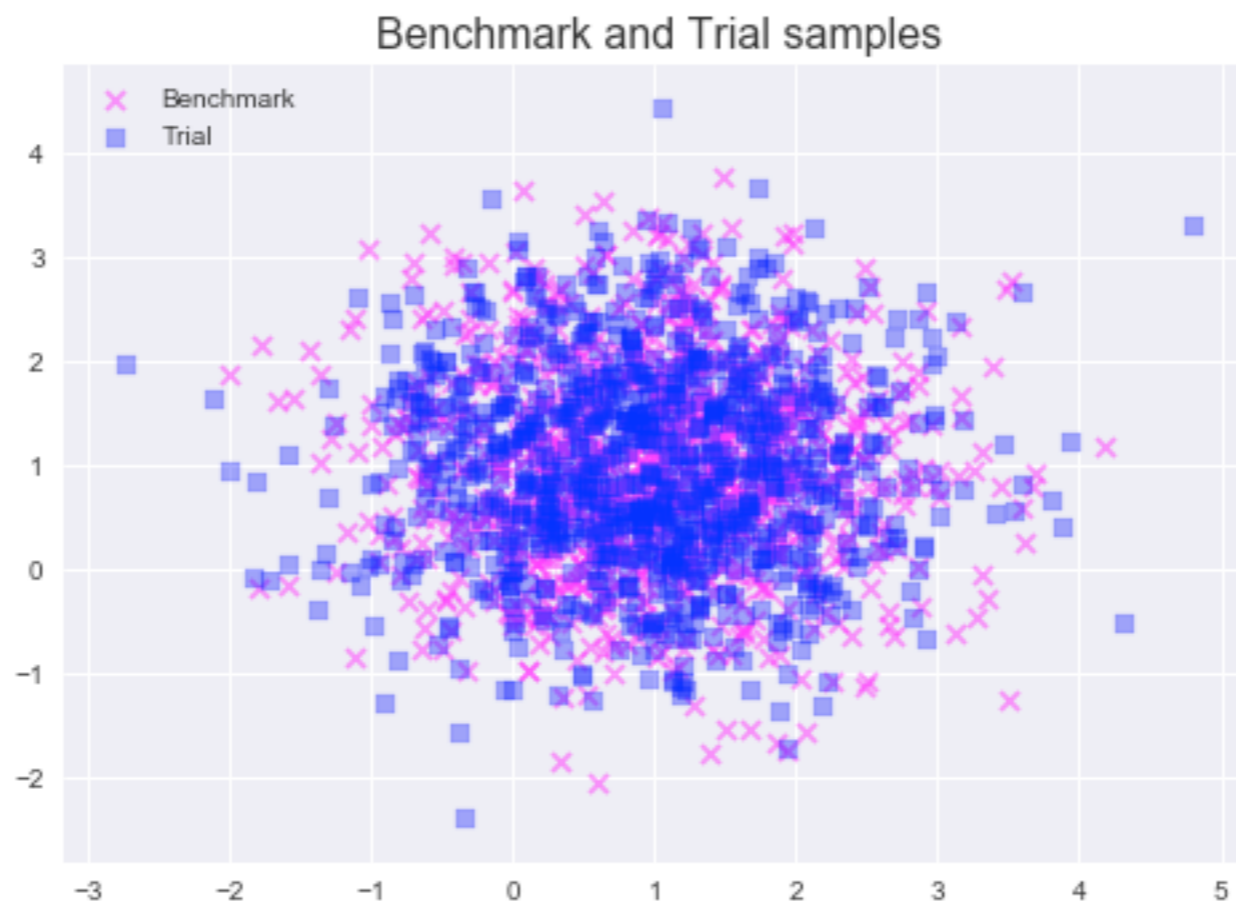
# > Two-sample Test

Two sets:

$$\begin{aligned} \text{Trial:} \quad & \mathcal{T} = \{\mathbf{x}_1, \dots, \mathbf{x}_{N_T}\} \stackrel{\text{iid}}{\sim} p_T \\ \text{Benchmark:} \quad & \mathcal{B} = \{\mathbf{x}'_1, \dots, \mathbf{x}'_{N_B}\} \stackrel{\text{iid}}{\sim} p_B \end{aligned} \quad \mathbf{x}_i, \mathbf{x}'_i \in \mathbb{R}^D$$

probability distributions  $p_B, p_T$  unknown

Are  $B, T$  drawn from the same prob. distribution?



... hard!

# > Two-sample Test

## RECIPE:

### 1. Density Estimator

—————→ reconstruct PDFs 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

—————→ accept/reject  $H_0$

# > 1. Density Estimator

Divide the space in squared bins?

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

**Need un-binned  
multivariate approach**

Find PDFs *estimators*:  $\hat{p}_B(\mathbf{x}), \hat{p}_T(\mathbf{x})$   
e.g. based on densities of points:

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

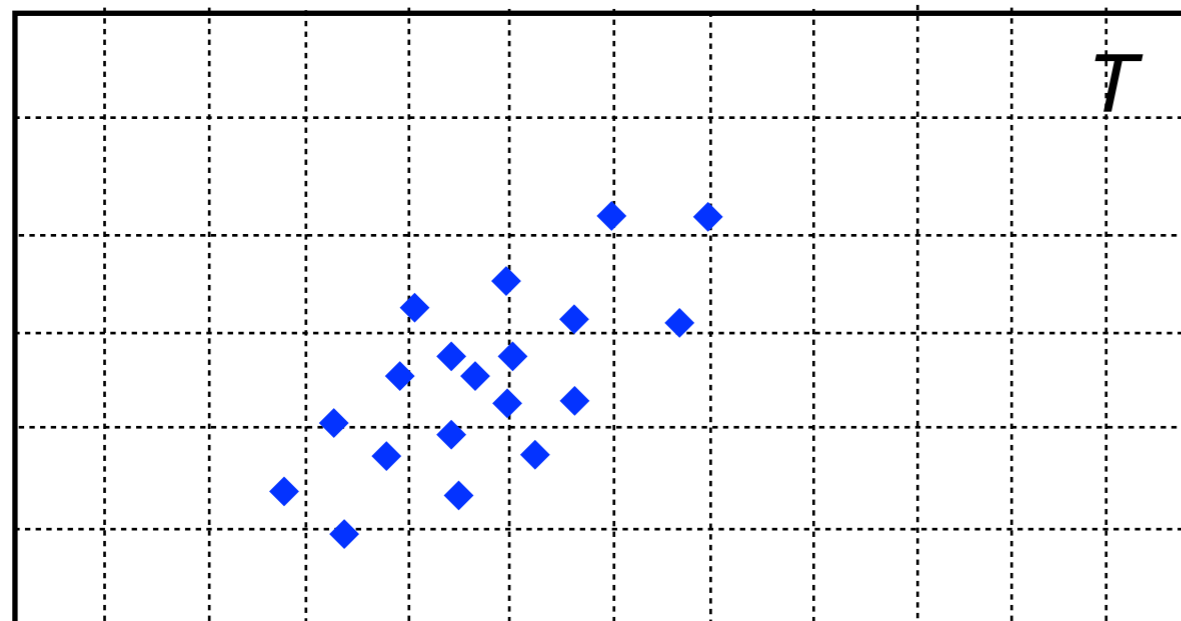
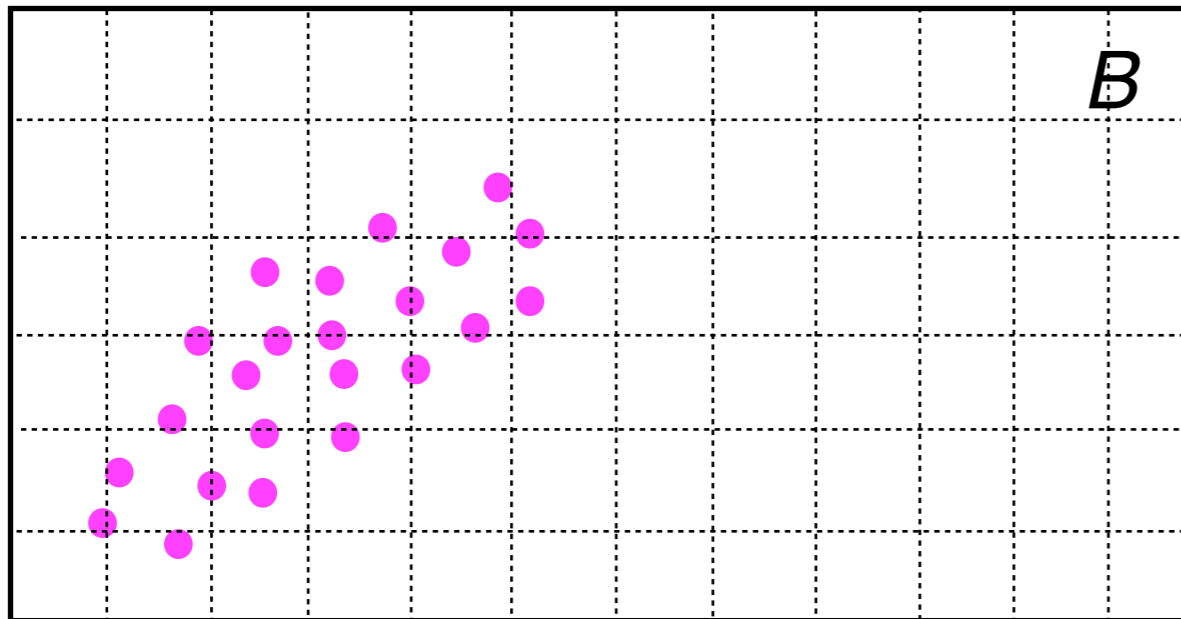
**Nearest Neighbors!**

[Schilling - 1986] [Henze - 1988]

[Wang et al. - 2005, 2006]

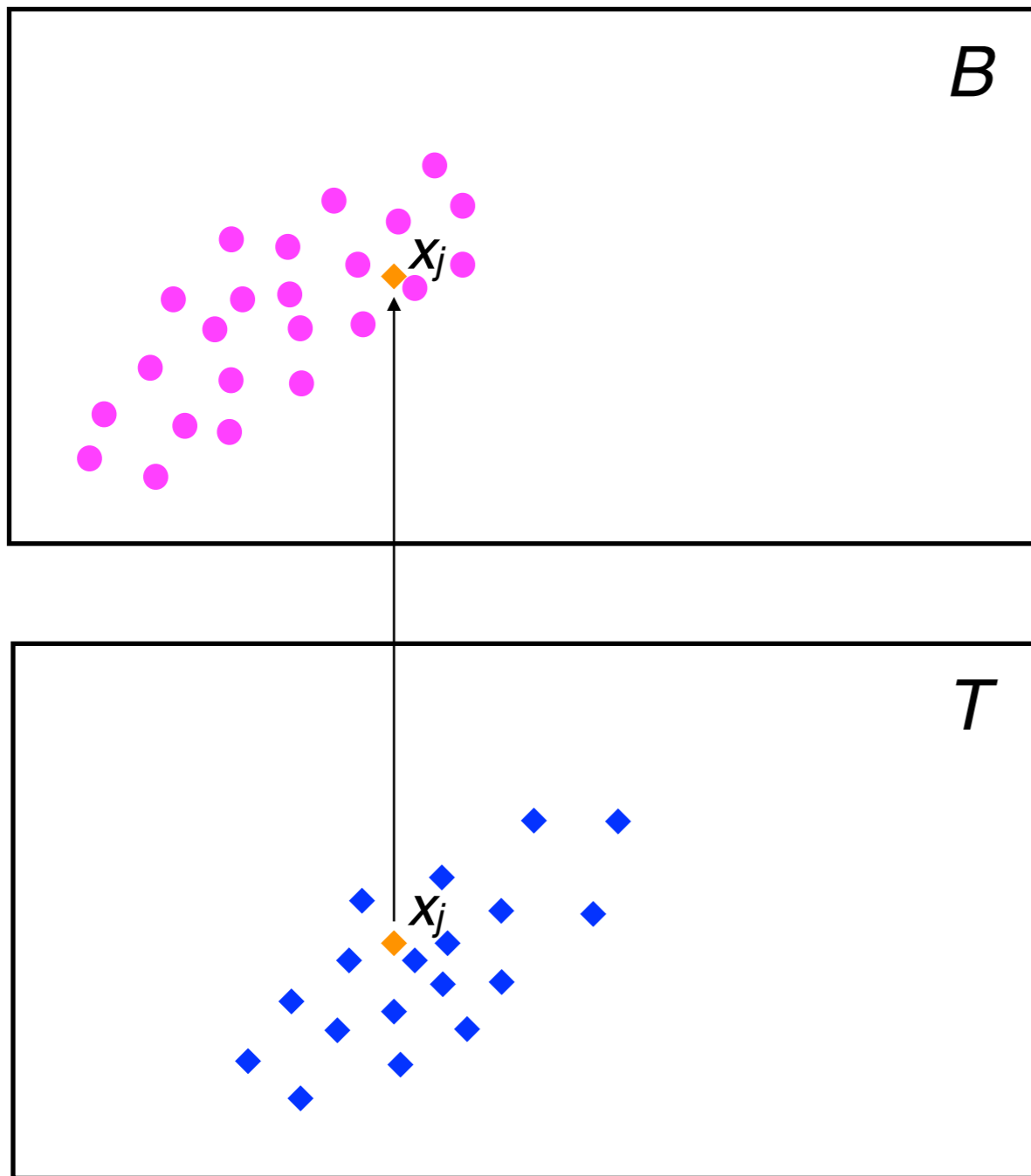
[Dasu et al. - 2006] [Perez-Cruz - 2008]

[Sugiyama et al. - 2011] [Kremer et al., 2015]

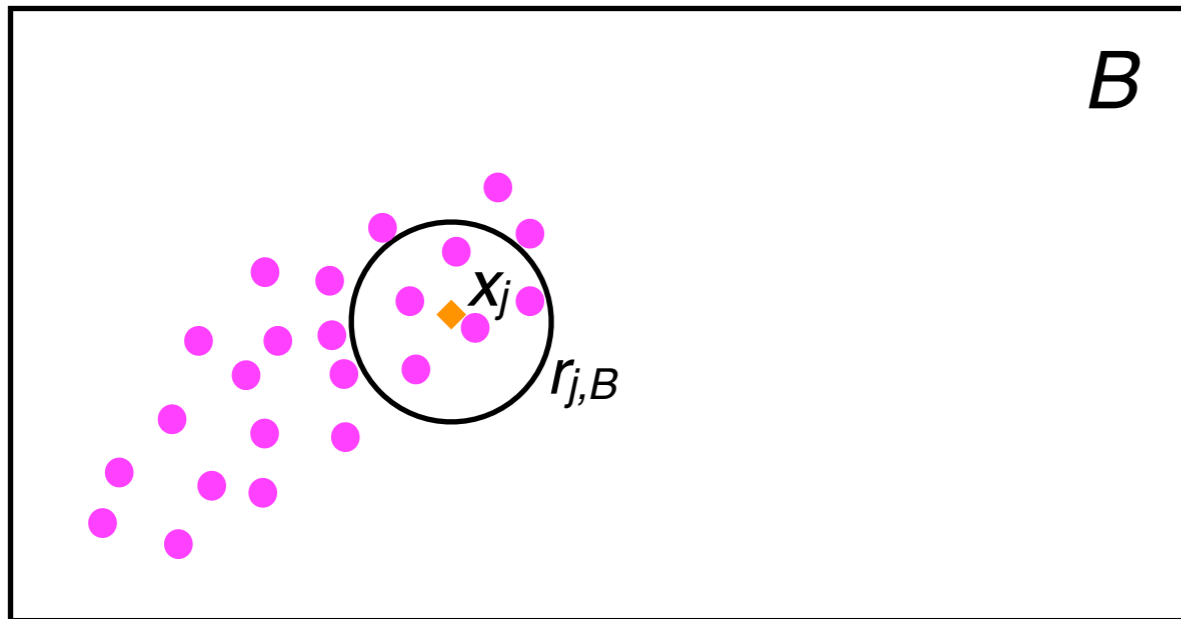


# > 1. Density Estimator

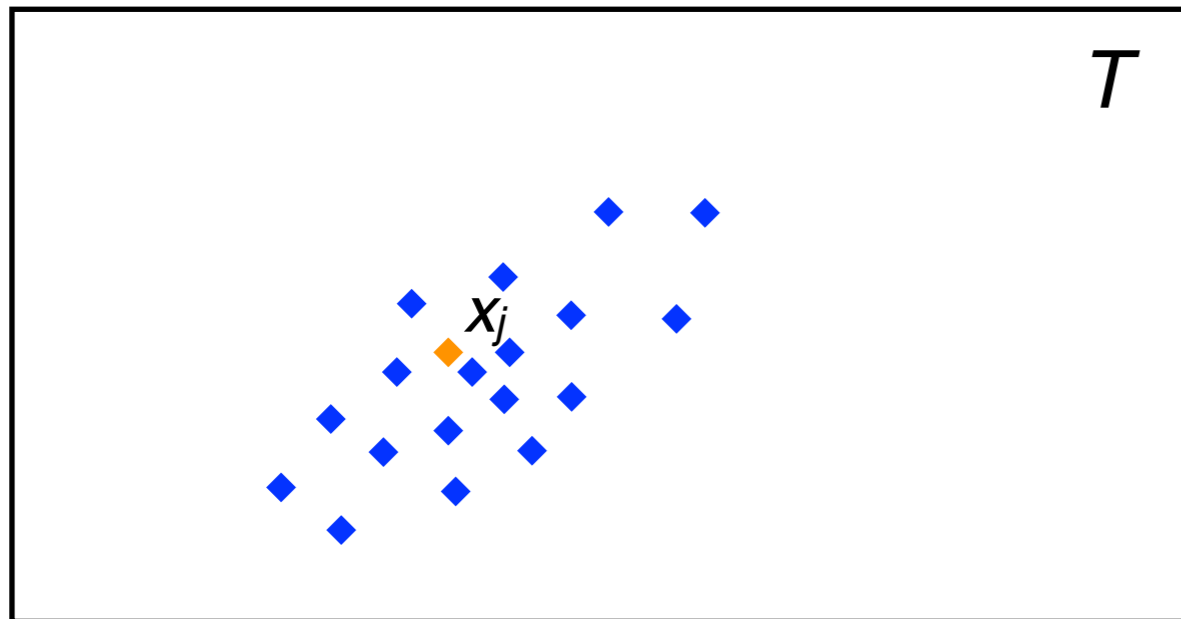
- Fix integer  $K$ .
- Choose query point  $x_j$  in  $T$  and draw it in  $B$ .



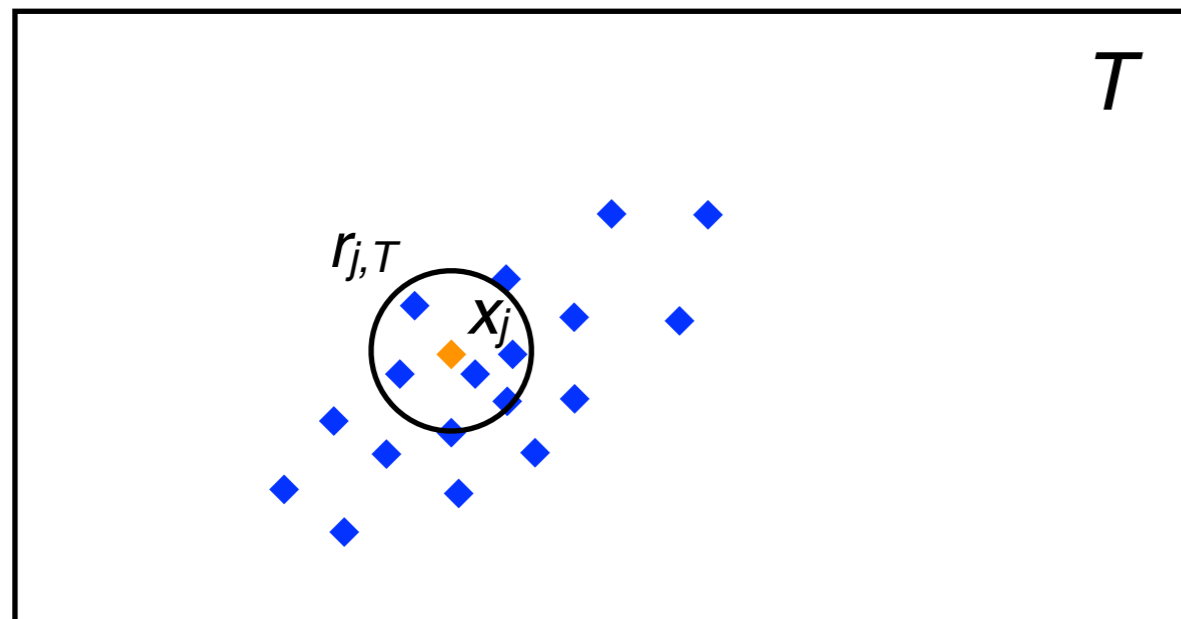
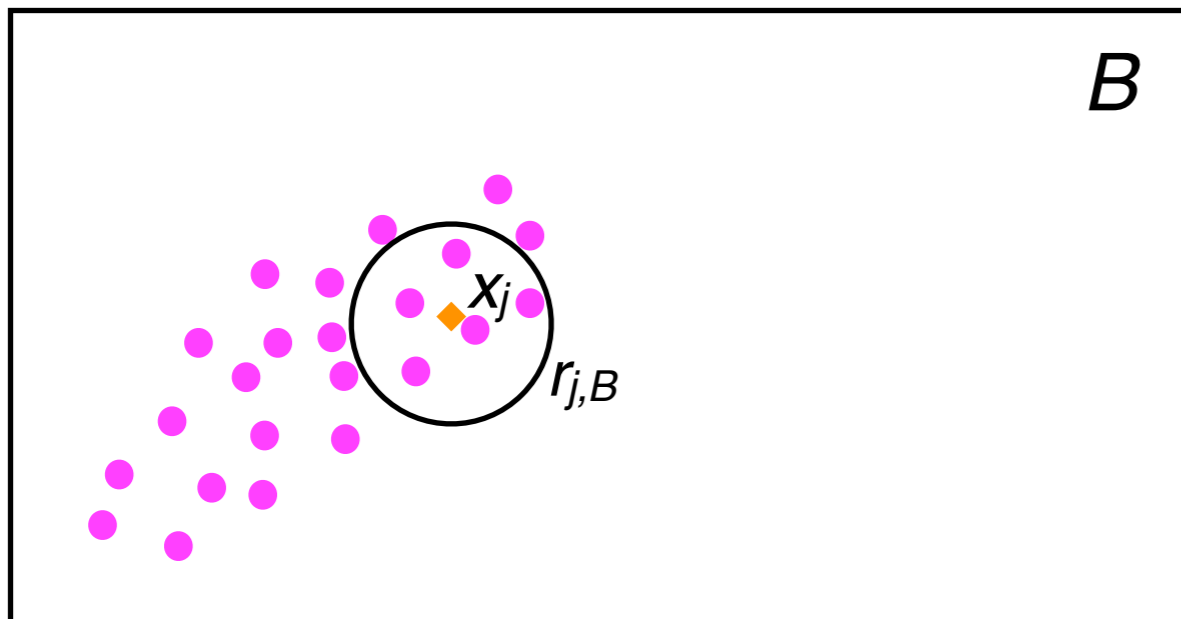
# > 1. Density Estimator



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

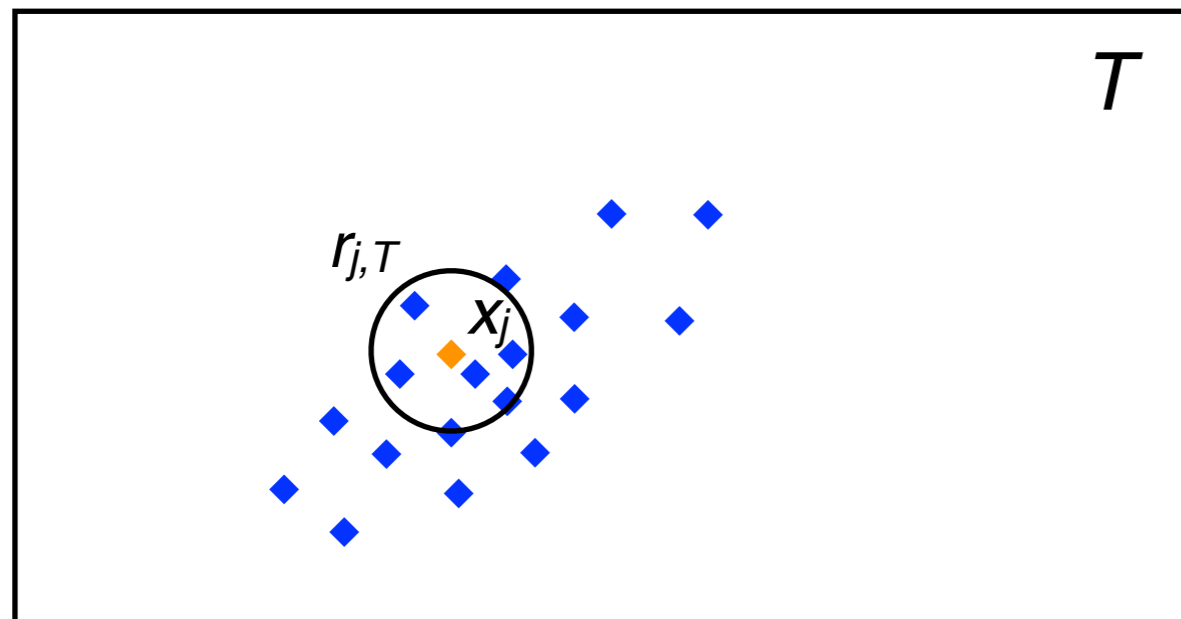
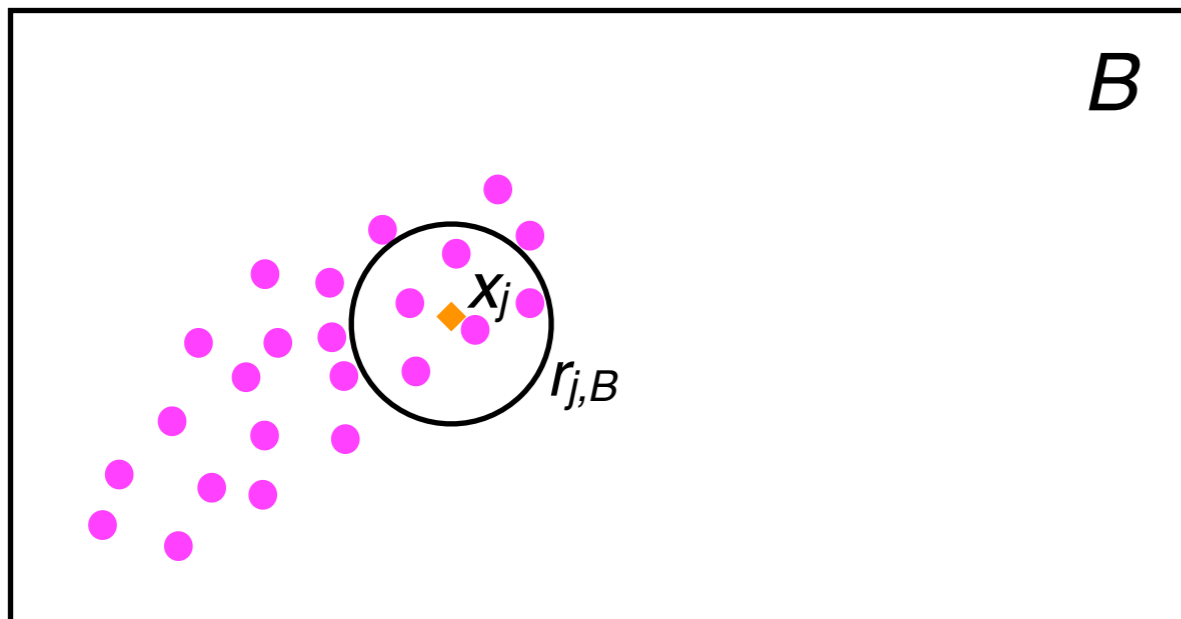


# > 1. Density Estimator



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

# > 1. Density Estimator



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

- 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 of the “distance” between 2 PDFs

- Define **Test Statistic**:  
(detect under-/over-densities)

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

- Related to Kullback-Leibler divergence as:  $\text{TS}(\mathcal{B}, \mathcal{T}) = \hat{D}_{\text{KL}}(\hat{p}_T || \hat{p}_B)$

$$D_{\text{KL}}(p||q) \equiv \int_{\mathbb{R}^D} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$

- From NN-estimated PDFs:  $\text{TS}_{\text{obs}} = \frac{D}{N_T} \sum_{j=1}^{N_T} \log \frac{r_{j,B}}{r_{j,T}} + \log \frac{N_B}{N_T - 1}$

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

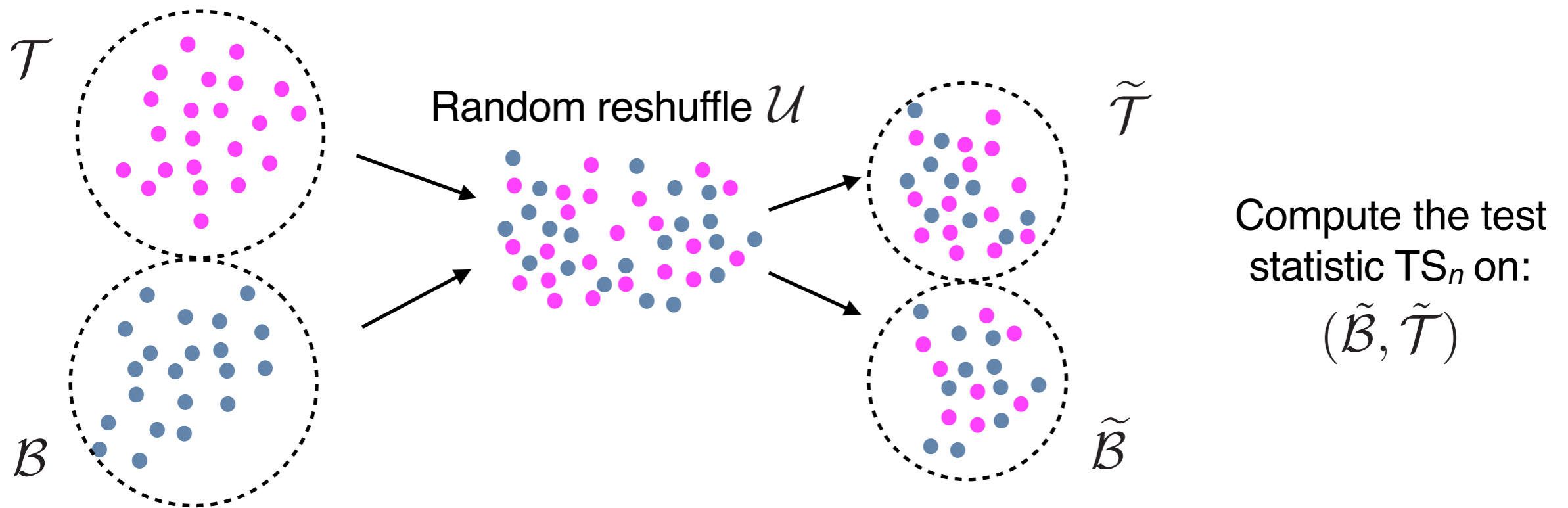
[Wang et al. - 2005, 2006]



# > 3. Test Statistic Distribution

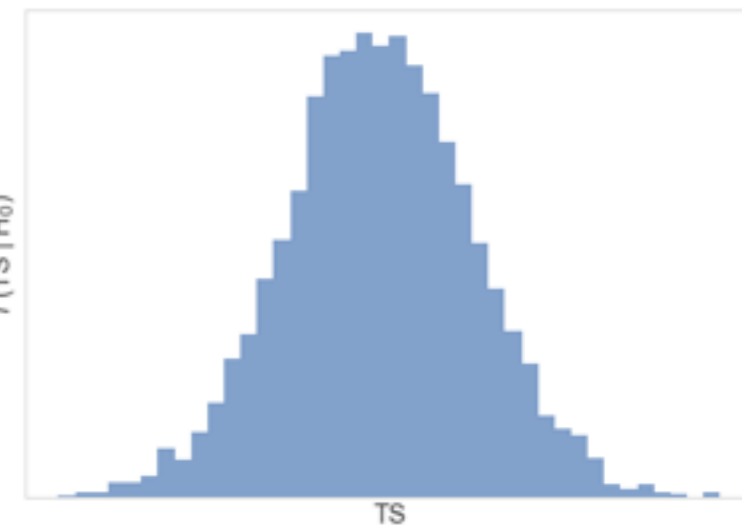
## How is TS distributed? **Permutation test!**

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



Repeat many times.

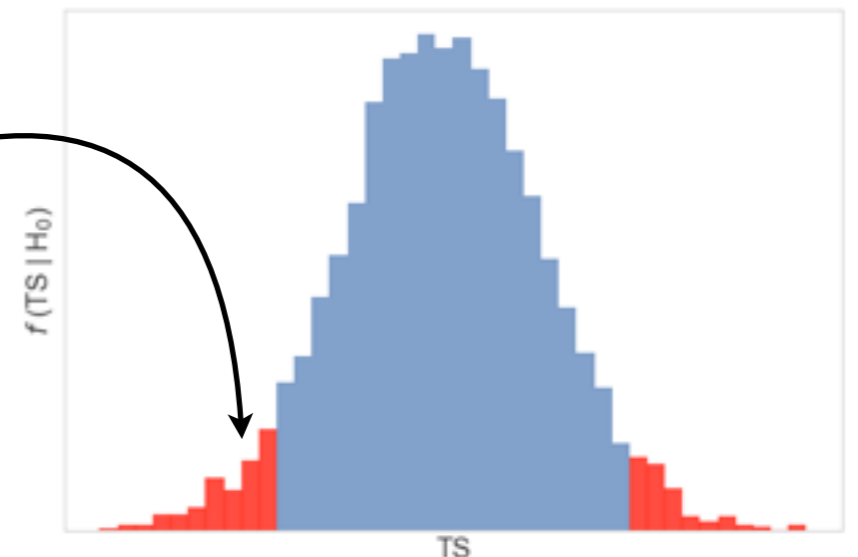
Distribution of TS under  $H_0$ :  $f(\text{TS}|H_0) \leftarrow \{TS_n\}$   
[asymptotically normal with zero mean]



## > 4. *p*-value

- $\hat{\mu}, \hat{\sigma}$  : mean, variance of TS distribution  $f(\text{TS}|H_0)$
- Standardize the TS:  $\text{TS} \rightarrow \text{TS}' \equiv \frac{\text{TS} - \hat{\mu}}{\hat{\sigma}}$
- TS' distributed according to  $f'(\text{TS}'|H_0) = \hat{\sigma} f(\hat{\mu} + \hat{\sigma}\text{TS}'|H_0)$
- Two-sided *p*-value:

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

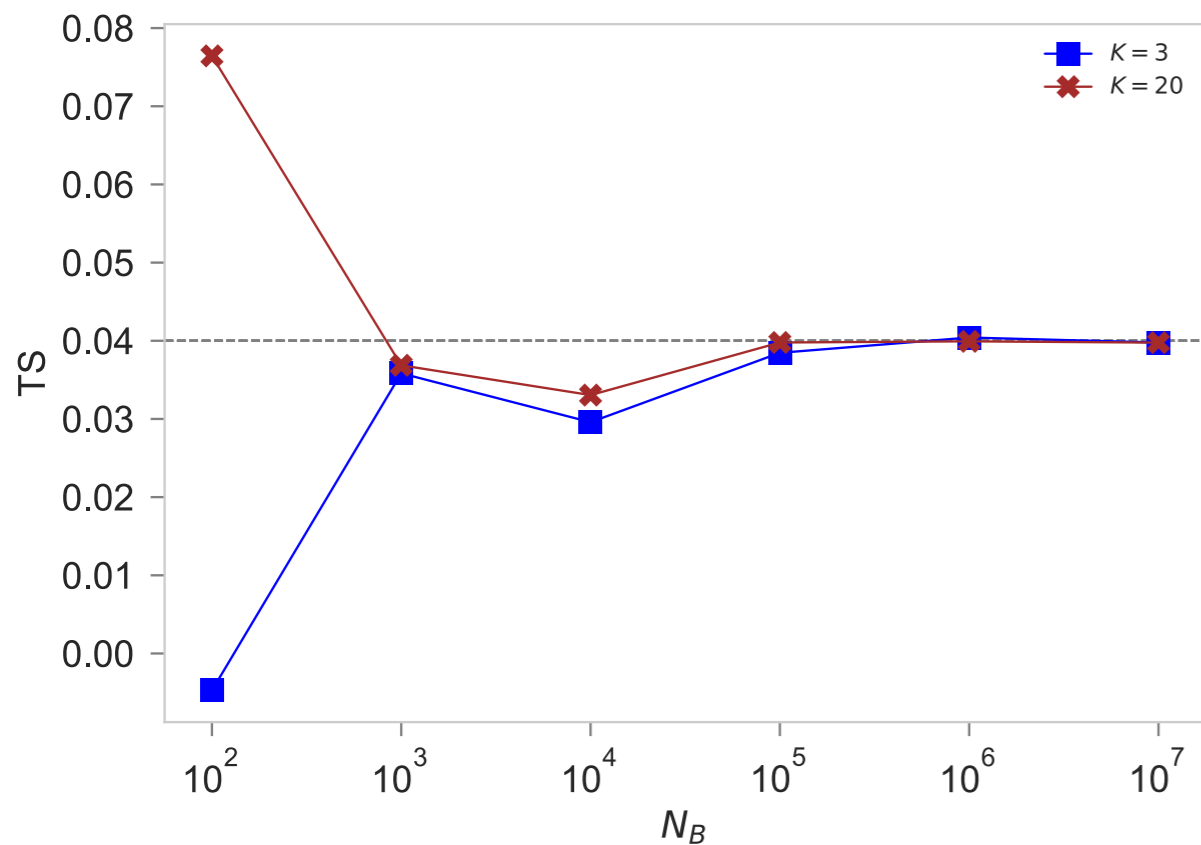


- Equivalent significance:  $Z \equiv \Phi^{-1}(1 - p/2)$

# > 2D Gaussian Example

$$p_B = \mathcal{N}(\boldsymbol{\mu}_B, \boldsymbol{\Sigma}_B) \quad p_T = \mathcal{N}(\boldsymbol{\mu}_T, \boldsymbol{\Sigma}_T) \quad \boldsymbol{\Sigma}_B = \boldsymbol{\Sigma}_T = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

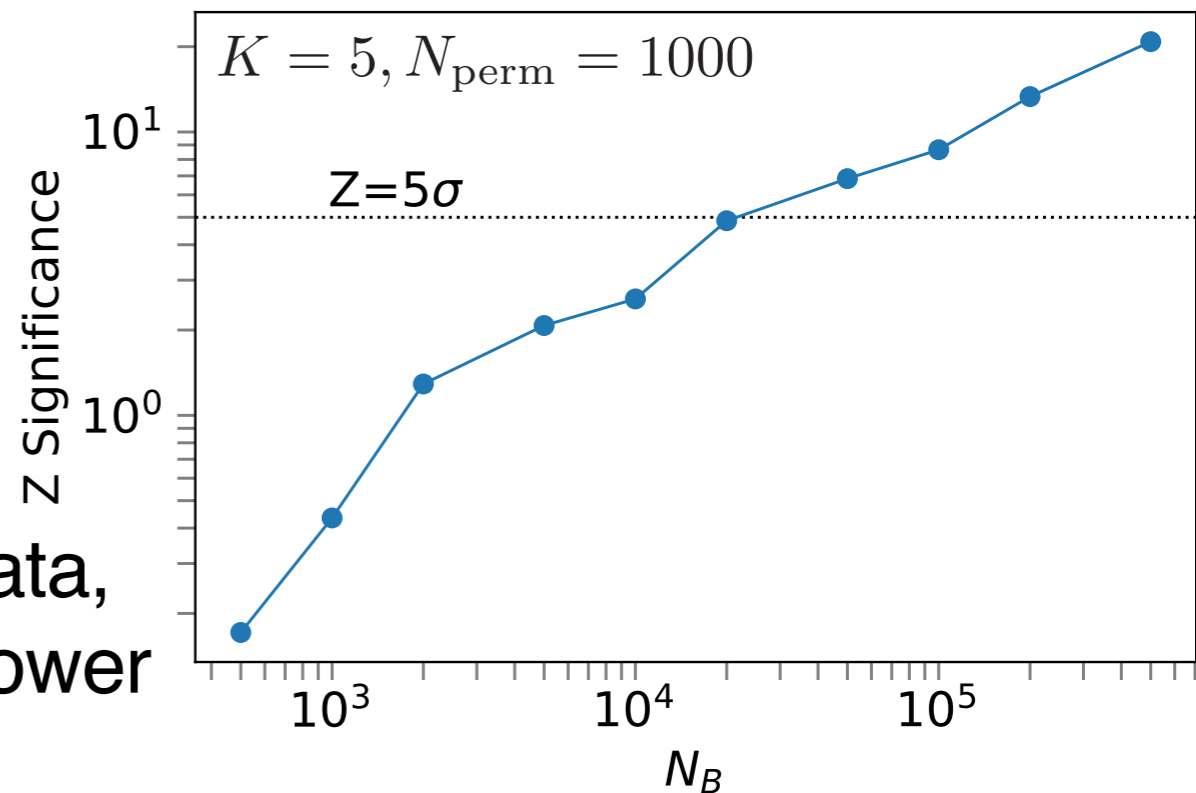
$$\boldsymbol{\mu}_B = \begin{pmatrix} 1.0 \\ 1.0 \end{pmatrix} \quad \boldsymbol{\mu}_T = \begin{pmatrix} 1.2 \\ 1.2 \end{pmatrix}$$



exact KL  
divergence

$$\boldsymbol{\mu}_B = \begin{pmatrix} 1.0 \\ 1.0 \end{pmatrix} \quad \boldsymbol{\mu}_T = \begin{pmatrix} 1.15 \\ 1.15 \end{pmatrix}$$

more data,  
more power



# > NN2ST: Summary

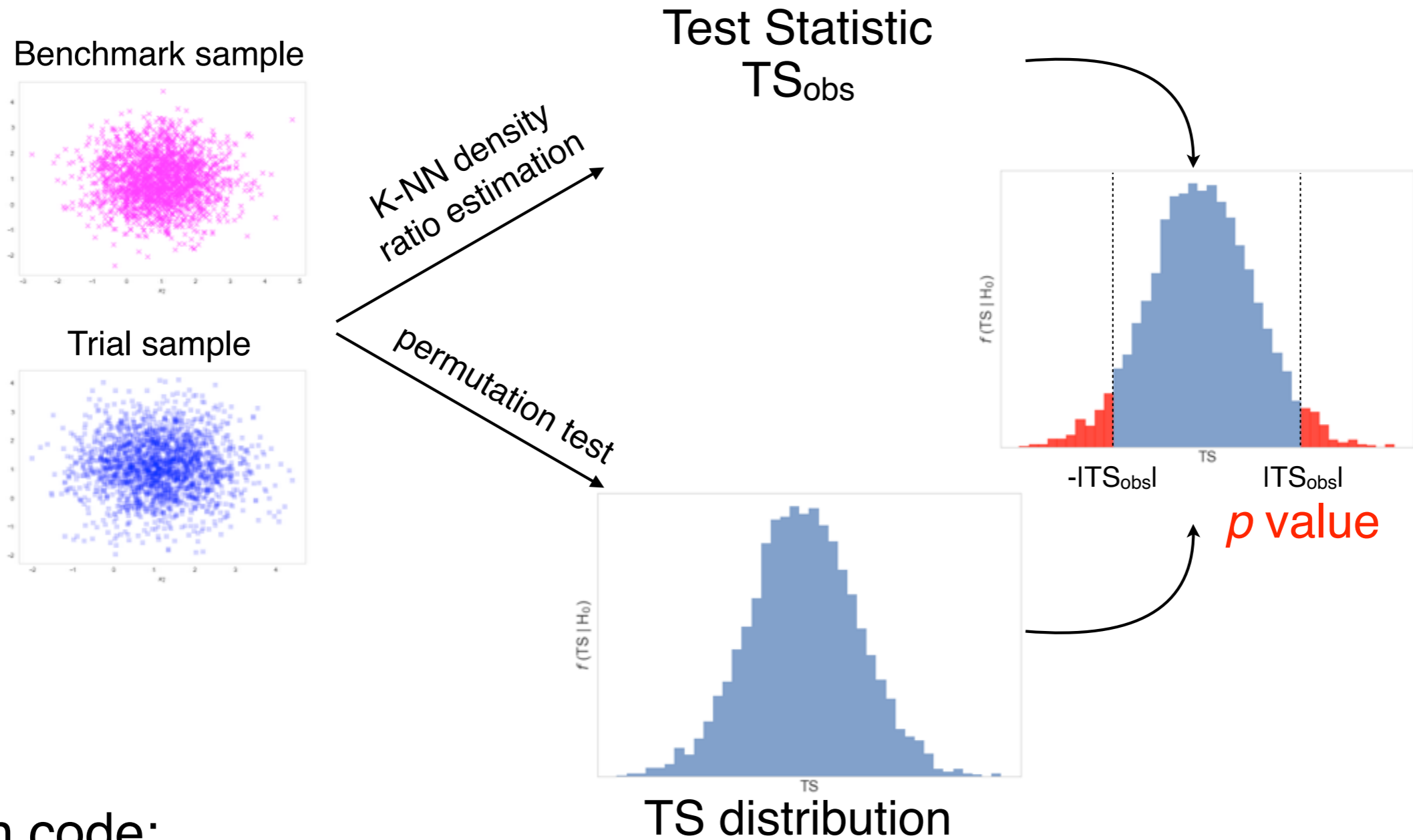
## INPUT:

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

## OUTPUT:

**$p$ -value of the null hypothesis  $H_0: p_B = p_T$**   
[check compatibility between 2 samples]

# > NN2ST: Summary



Python code:

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

### **1. Statistical test of dataset compatibility**

- Nearest-Neighbors Two-Sample Test
- **Identify Discrepancies**
- Include Uncertainties

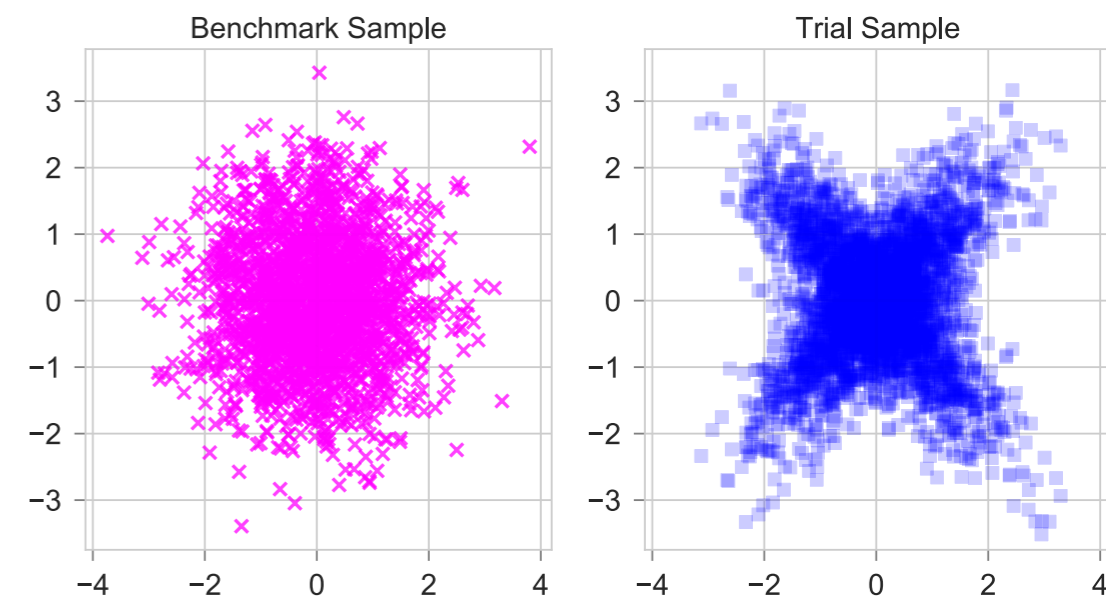
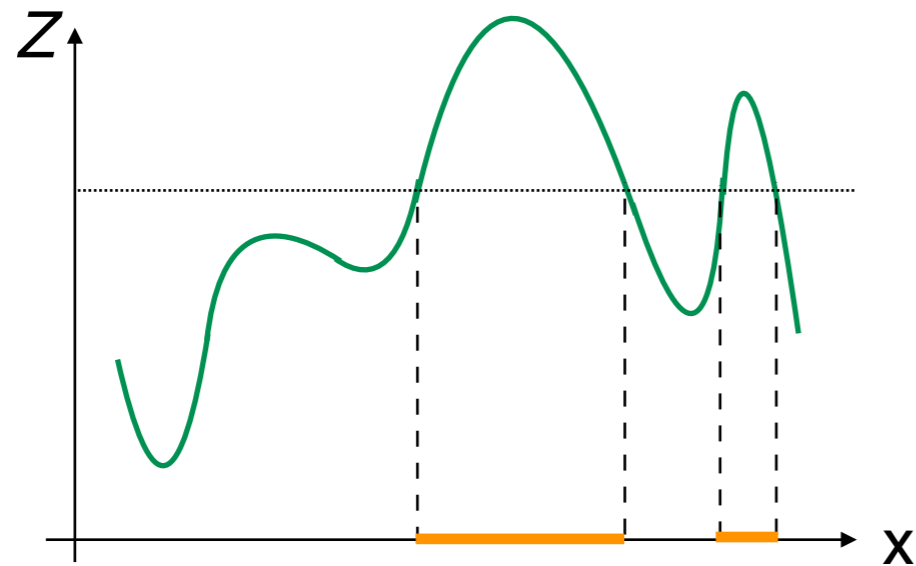
### **2. Applications to High-Energy Physics**

# > Where are the discrepancies?

## Bonus: Characterize regions with significant discrepancies

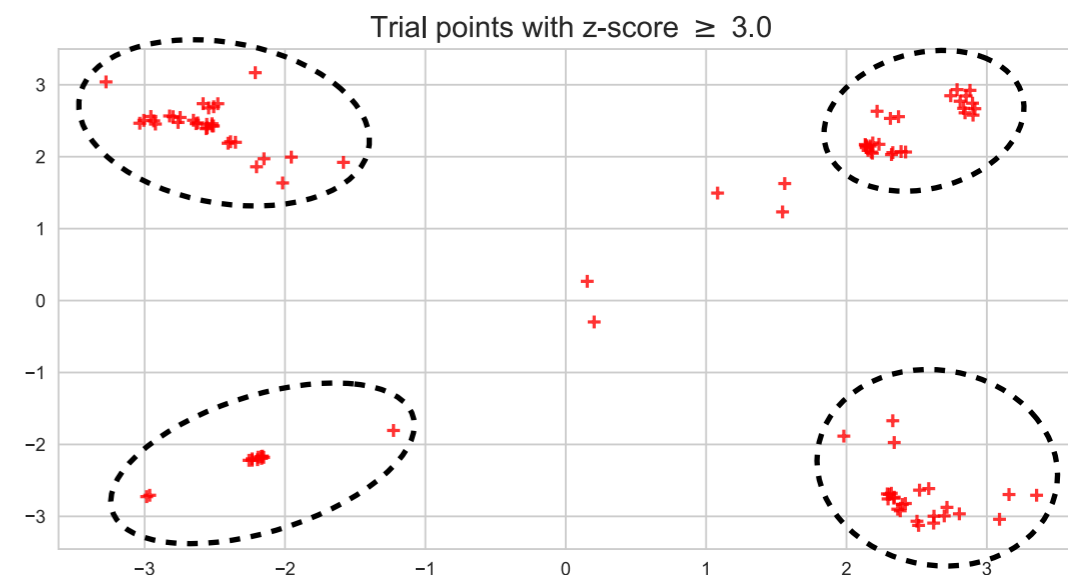
1. “Score” field over  $\mathcal{T}$ :  $Z(\mathbf{x}_j) \equiv \frac{u(\mathbf{x}_j) - \bar{u}}{\sigma_u}$

with:  $u(\mathbf{x}_j) \equiv \log \frac{r_{j,B}}{r_{j,T}}$   
 $TS_{\text{obs}} = D \bar{u} + \text{const}$



2. Identify points where  $Z(\mathbf{x}) > c$   
They contribute the most to large  $TS_{\text{obs}}$   
→ high-discrepancy (anomalous) regions

3. Apply a clustering algorithm to group them



## 1. Statistical test of dataset compatibility

- Nearest-Neighbors Two-Sample Test
- Identify Discrepancies
- Include Uncertainties

## 2. Applications to High-Energy Physics



# > Sample Uncertainties

## How to include sample uncertainties?

1. Model feature uncertainties

$$F_{\mathcal{B}}(\mathbf{x}), F_{\mathcal{T}}(\mathbf{x})$$

[e.g. zero-mean gaussians]

2. New samples by adding random noise sampled from  $F_{\mathcal{B},\mathcal{T}}$ :

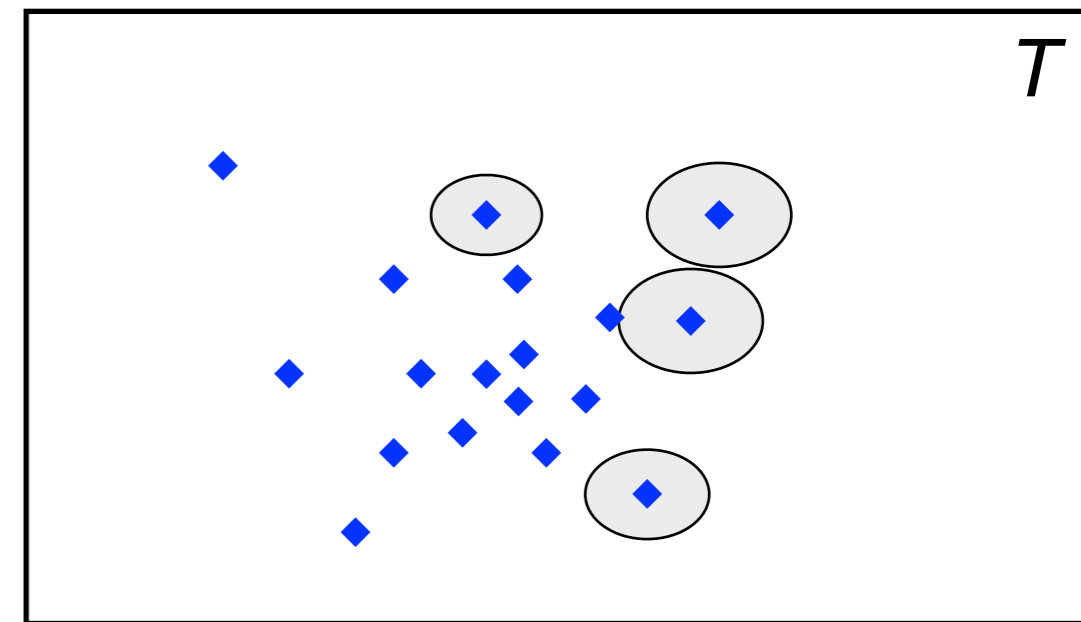
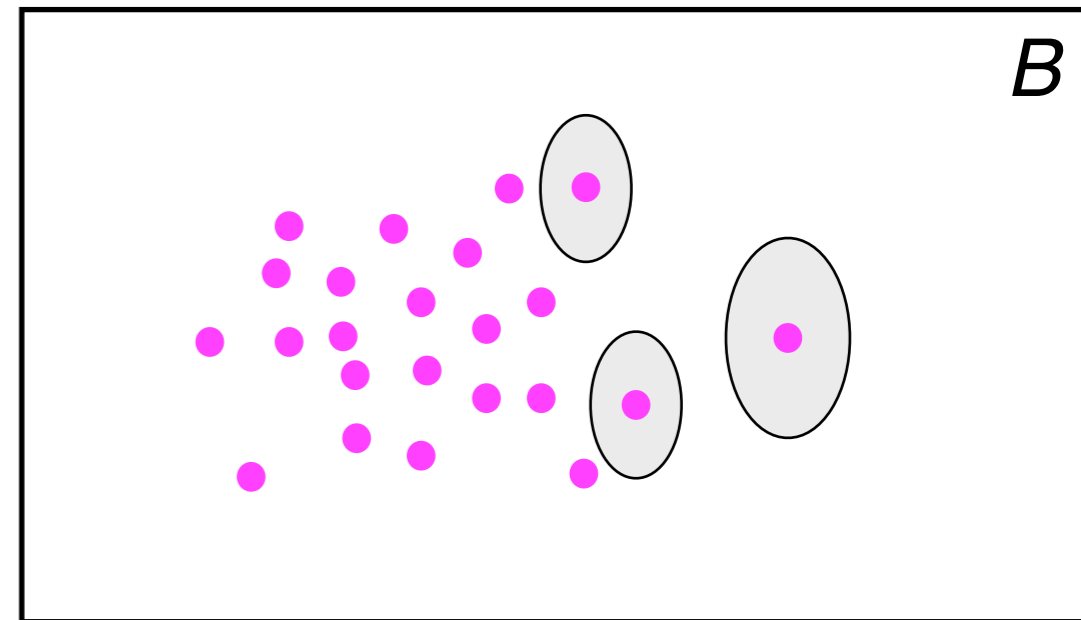
$$\mathcal{T}_u = \{\mathbf{x}_i + \Delta\mathbf{x}_i\}_{i=1}^{N_T}$$

$$\mathcal{B}_u = \{\mathbf{x}'_i + \Delta\mathbf{x}'_i\}_{i=1}^{N_B}$$

3. Compute TS on new samples

$$\text{TS}_u \equiv \text{TS}(\mathcal{B}_u, \mathcal{T}_u) = \text{TS}_{\text{obs}} + U$$

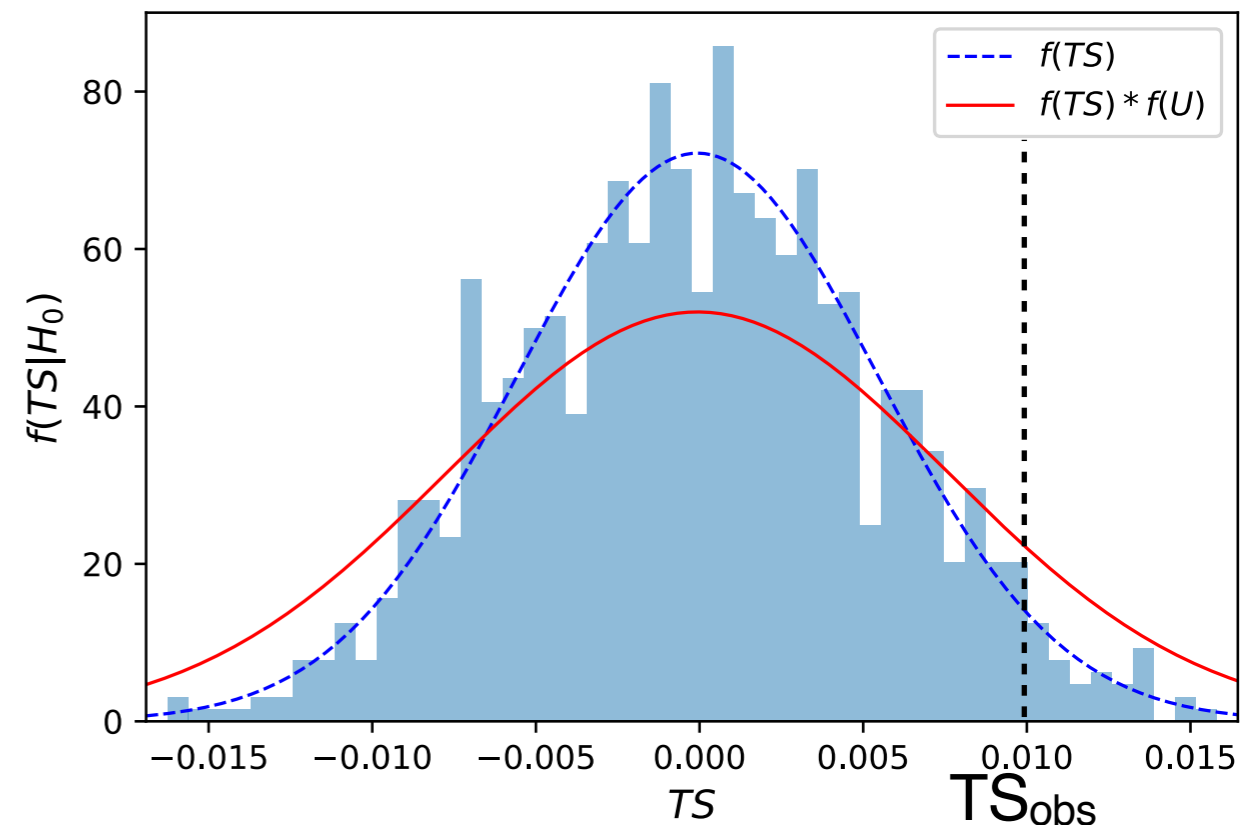
4. Repeat many times to reconstruct  $f(U)$



# > Sample Uncertainties

## How to include sample uncertainties?

- $f(TS_u)$  is a convolution:  $f(TS_u|H_0) = f(TS|H_0) * f(U)$   
 $f(TS_u)$  more spread than  $f(TS)$
- $p$ -value computed from  $f(TS_u)$
- weaker significance,  
power degradation



# > 2D Gaussian with Uncertainties

$B, T$  gaussian samples:

$$p_B = \mathcal{N}(\boldsymbol{\mu}_B, \boldsymbol{\Sigma}_B) \quad p_T = \mathcal{N}(\boldsymbol{\mu}_T, \boldsymbol{\Sigma}_T)$$

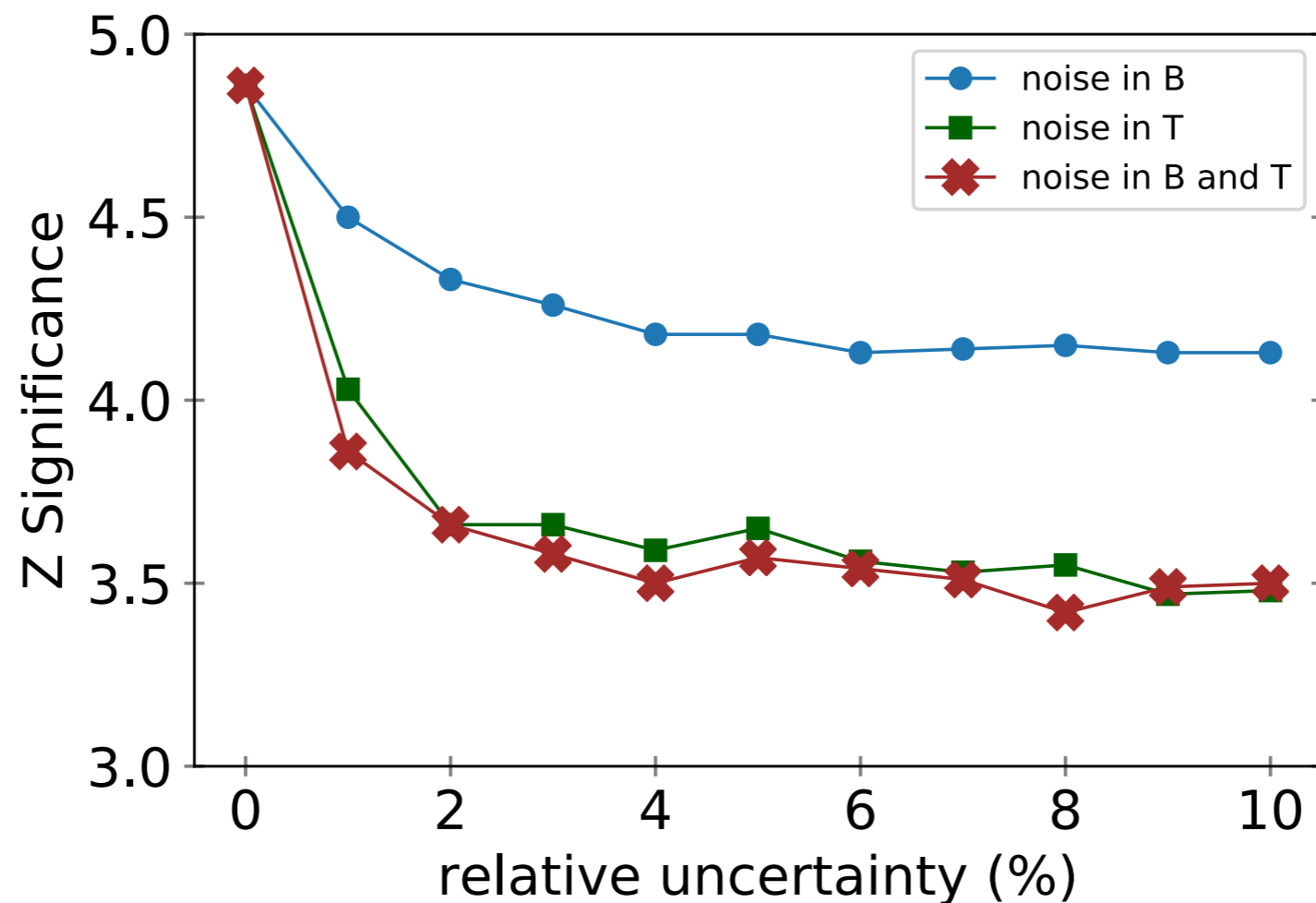
$$\boldsymbol{\mu}_B = \begin{pmatrix} 1.0 \\ 1.0 \end{pmatrix} \quad \boldsymbol{\mu}_T = \begin{pmatrix} 1.15 \\ 1.15 \end{pmatrix}$$

$$\boldsymbol{\Sigma}_B = \boldsymbol{\Sigma}_T = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

gaussian uncorrelated errors  
(diagonal covariance)  
with fixed relative uncertainty

$$\sigma_i = \epsilon x_i$$

for each feature component  $i$



## > NN2ST: Summary

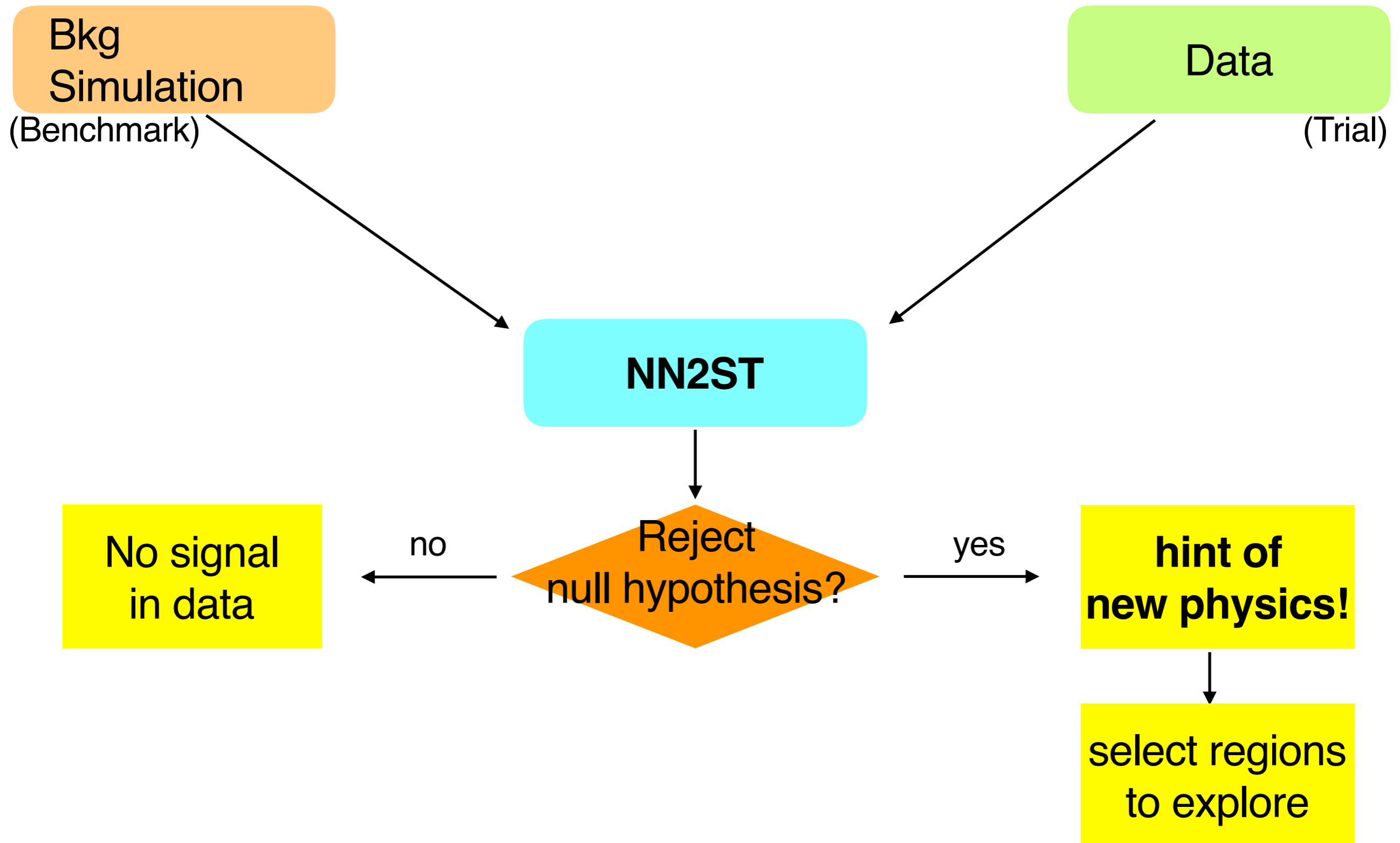
- ✓ general, model-independent
- ✓ fast, no optimization
  - [  $N_{B,T}=20k$ ,  $K=5$ ,  $N_{perm}=1k$ ,  $D=2$ :  $t \sim 2$  mins
  - $N_{B,T}=20k$ ,  $K=5$ ,  $N_{perm}=1k$ ,  $D=8$ :  $t \sim 50$  mins ]
- ✓ sensitive to **unspecified** signals
- ✓ useful when no variable can separate sig/bkg
- ✓ helps finding signal regions, optimal cuts, ...
- ✓ flexible to incorporate uncertainties
  
- ✗ need to run for each sample pair
- ✗ permutation test is bottleneck

## 1. Statistical test of dataset compatibility

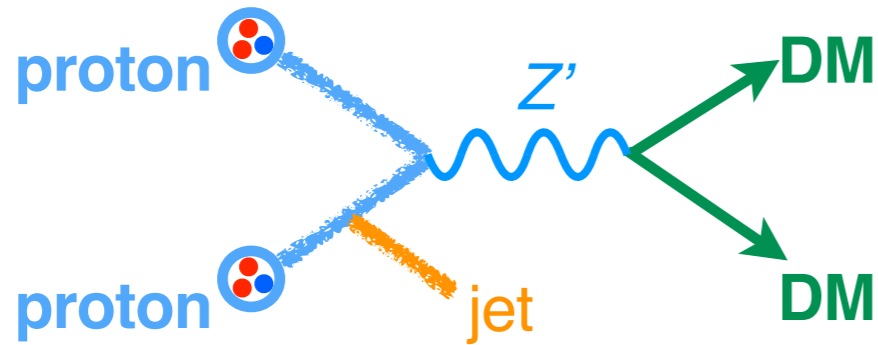
- Nearest-Neighbors Two-Sample Test
- Identify Discrepancies
- Include Uncertainties

## 2. Applications to High-Energy Physics

# > Our Method



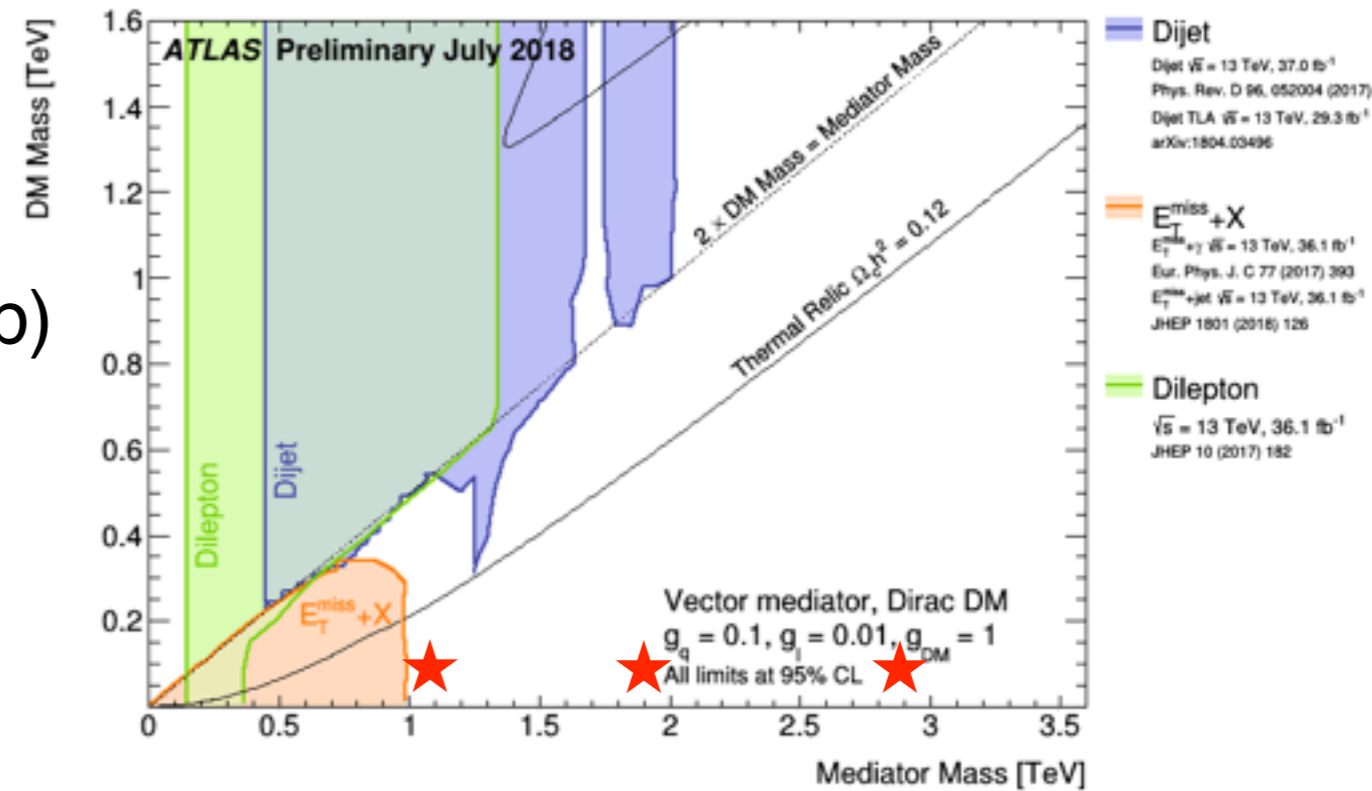
# > DM search @ LHC



DM +  $Z'$   
vector mediator

$m_{\text{DM}} = 100 \text{ GeV}$   
 $m_{Z'} = 1.2, 2, 3 \text{ TeV}$   
 $g_{\text{DM}} = 1, g_q = 0.1$   
 $\sqrt{s} = 13 \text{ TeV}$

- “proof-of-principle” study
- bkg:  $Z \rightarrow \nu\bar{\nu} + (1, 2) j$  ( $\sigma_{\text{bkg}}=202.6 \text{ pb}$ )  
sub-leading bkg's not included
- no full detector effects  
(generic Delphes profile)



**Benchmark: BKG<sub>1</sub>**  
**Trials: BKG<sub>2</sub> + SIG**  
 $K = 5$   
 $N_{\text{perm}} = 3000$

## 8 features:

- n. of jets
- $p_T, \eta$  of 2 leading jets
- $E_T^{\text{miss}}, H_T$
- $\Delta\phi_{E_T^{\text{miss}}, j_1}$

# > DM search @ LHC

**B:** **BKG<sub>1</sub>** (20k events)

**T1:** **BKG<sub>2</sub>** (20k events) + **SIG<sub>1</sub>** (2010 events)

**T2:** **BKG<sub>2</sub>** (20k events) + **SIG<sub>2</sub>** (375 events)

**T3:** **BKG<sub>2</sub>** (20k events) + **SIG<sub>3</sub>** (59 events)

$$N_{\text{sig}} = N_B \times \frac{\sigma_{\text{signal}}}{\sigma_{\text{bkg}}}$$

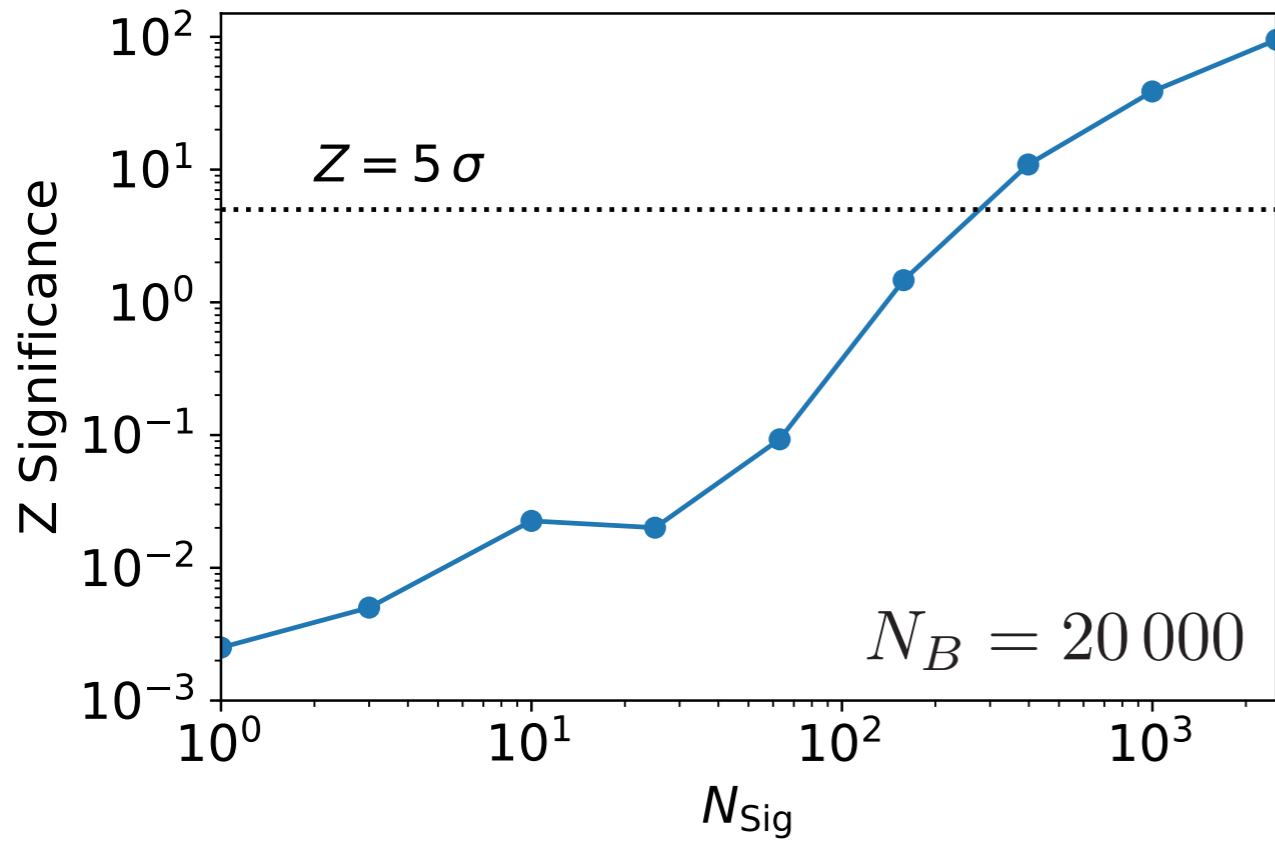
Sample	$M_{Z'}$	$\sigma_{\text{signal}}$	$Z_{\text{no uncert.}}$	$Z_{10\% \text{ rel uncert.}}$
T1	1.2 TeV	20.4 pb	40 $\sigma$	26 $\sigma$
T2	2 TeV	3.8 pb	13 $\sigma$	12 $\sigma$
T3	3 TeV	0.6 pb	2.7 $\sigma$	2.5 $\sigma$

still not  
real-world

- systematics: expect further degradation of results
- **the method has value, it is worth exploring**



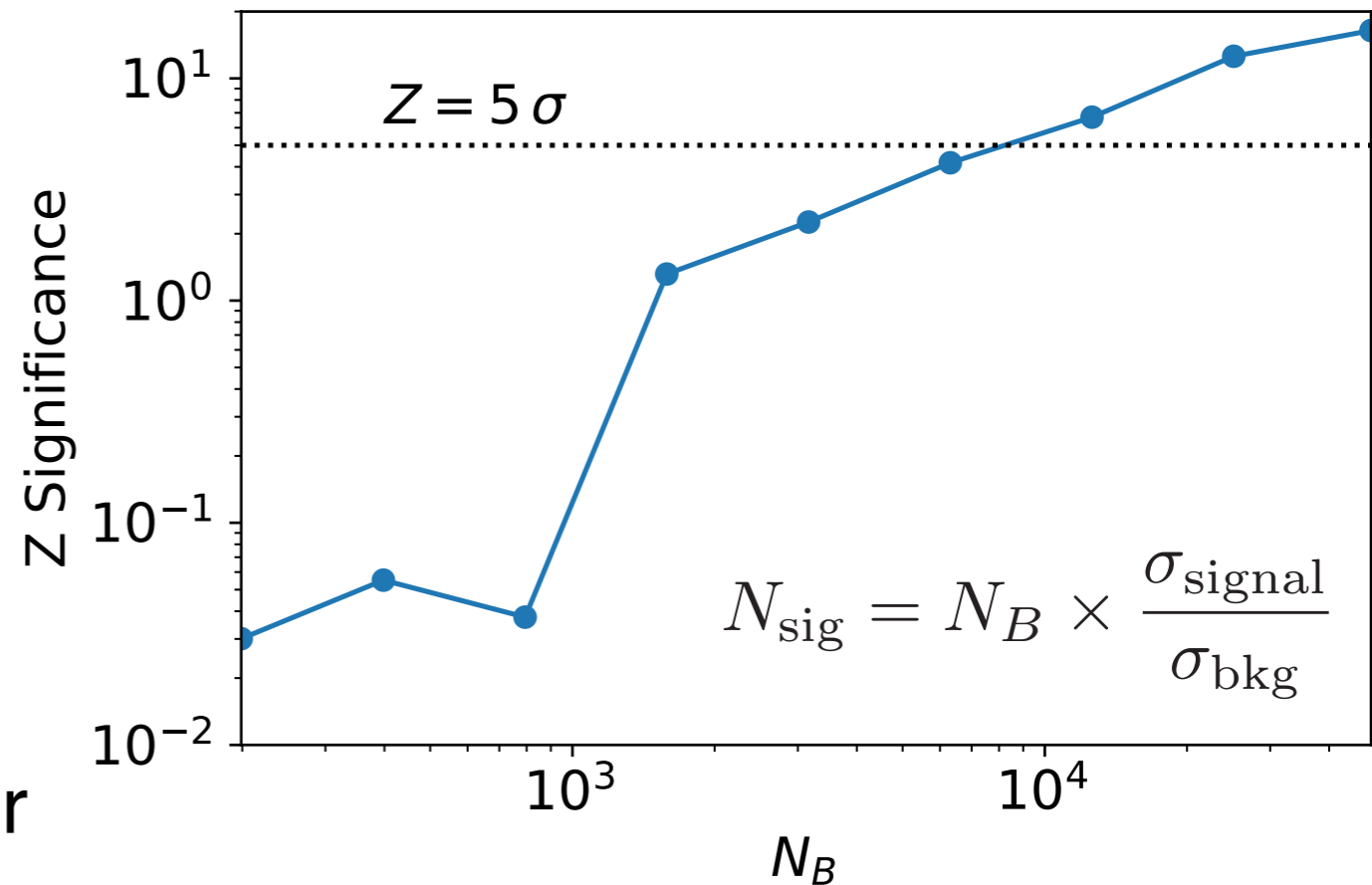
# > DM search @ LHC



stronger signal  
easier to discover

$$N_T = N_B + N_{\text{sig}}$$

more data,  
more power



## Directions for future work:

- adaptive choice of  $K$
- identifying discrepant regions in realistic situations (with  $Z$ -score method)
- validation tool for bkg:  
compatibility between MC sims. and data in control regions
- scalability
- ... your suggestions?

# > Take-Home Messages

## 1. New Statistical Test for BSM Physics

- assess degree of compatibility between 2 samples
- rooted on nearest neighbors, solid math foundations

## 2. NN2ST as a discovery tool

- powerful and model-independent
- lots of applications

## 3. NN2ST to guide searches

- identify regions of discrepancies

**BACK UP**

# > Model Selection

## how to choose $K$ ? **Model Selection!**

$$\text{True: } r(\mathbf{x}) = \frac{p_T(\mathbf{x})}{p_B(\mathbf{x})} \qquad \text{Estimated: } \hat{r}(\mathbf{x}) = \frac{\hat{p}_T(\mathbf{x})}{\hat{p}_B(\mathbf{x})}$$

Define the mean-square error:

$$\begin{aligned} L(r, \hat{r}) &= \frac{1}{2} \int [\hat{r}(\mathbf{x}') - r(\mathbf{x}')]^2 p_B(\mathbf{x}') d\mathbf{x}' \\ &= \frac{1}{2} \int \hat{r}(\mathbf{x}')^2 p_B(\mathbf{x}') d\mathbf{x}' - \int \hat{r}(\mathbf{x}) p_T(\mathbf{x}) d\mathbf{x} + \frac{1}{2} \int r(\mathbf{x}')^2 p_B(\mathbf{x}') d\mathbf{x}' \end{aligned}$$

$$\text{Estimate loss: } \hat{L}(r, \hat{r}) = \frac{1}{2N_B} \sum_{\mathbf{x}' \in \mathcal{B}} \hat{r}(\mathbf{x}')^2 - \frac{1}{N_T} \sum_{\mathbf{x} \in \mathcal{T}} \hat{r}(\mathbf{x})$$

Select optimal  $K$  minimizing the loss.

Alternatively: Point-Adaptive k-NN (PAk) [\[1802.10549\]](#)