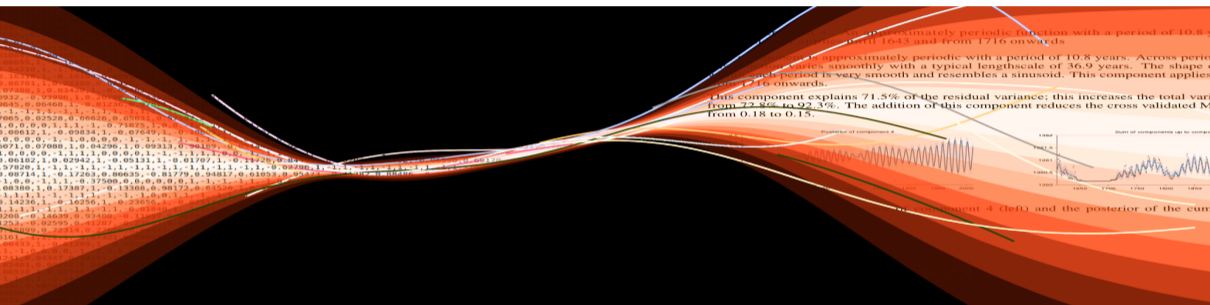


# Causal inference in geosciences with multidimensional kernel deviance measures



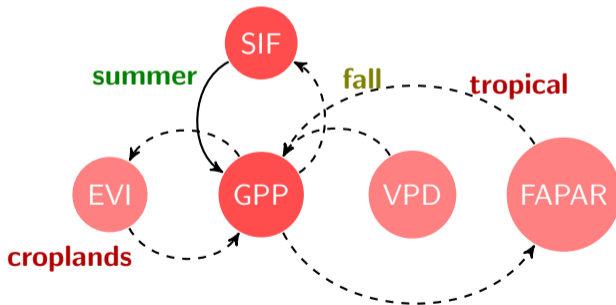
E. Díaz ,D. Sejdinovic, J. Ton  
A. Pérez-Suay, V. Laparra and G.  
Camps-Valls

Image Processing Laboratory (IPL)  
Universitat de València, Spain



## Outlook

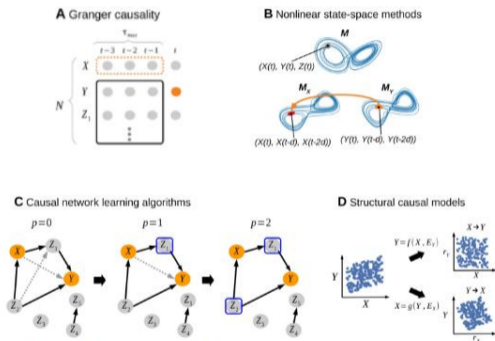
- Infer causal relations between r.v. is challenging
- Even more from pure observational data: no models involved, no ground truth!
- In GRS, causality is key to understand the Earth's system



## Outline

- ① Causal inference for instantaneous observations (CIIO)
- ② Kernel Conditional Deviance for Causal Inference (KCDC)
- ③ Our work:
  - Experiments
  - Developments to kernel deviance measures

Overview of causal inference methods



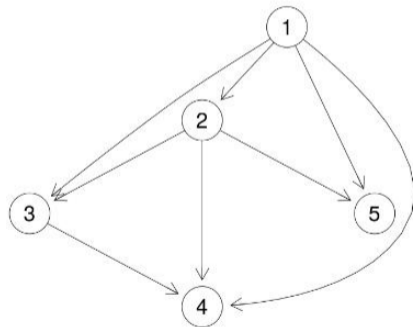
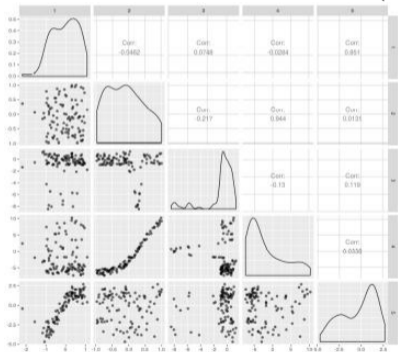
Acks: Runge et al 2019

KCDC, the method studied here and ANMs, the method to which we will compare its performance belong to group D.

- A. Multivariate granger causality tests
- B. Nonlinear state-space method CCM
- C. Causal network learning algorithms (conditional independence testing)
- D. Structural Causal models

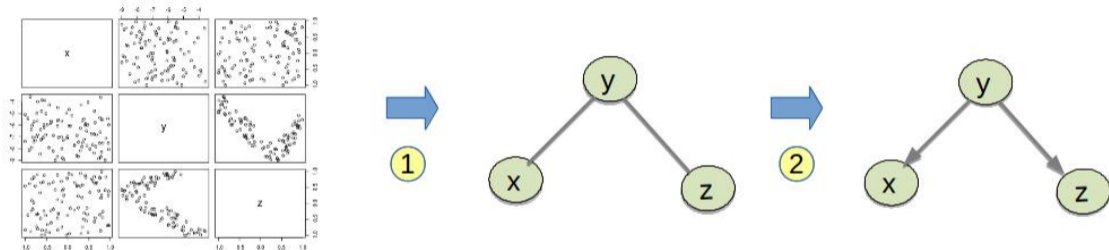
## Goal of Causal Inference for instantaneous observations

Given a system of  $p$  variables, with  $n$  observations available for each, learn underlying causal Directed Acyclic Graph (DAG)



## Two step learning process

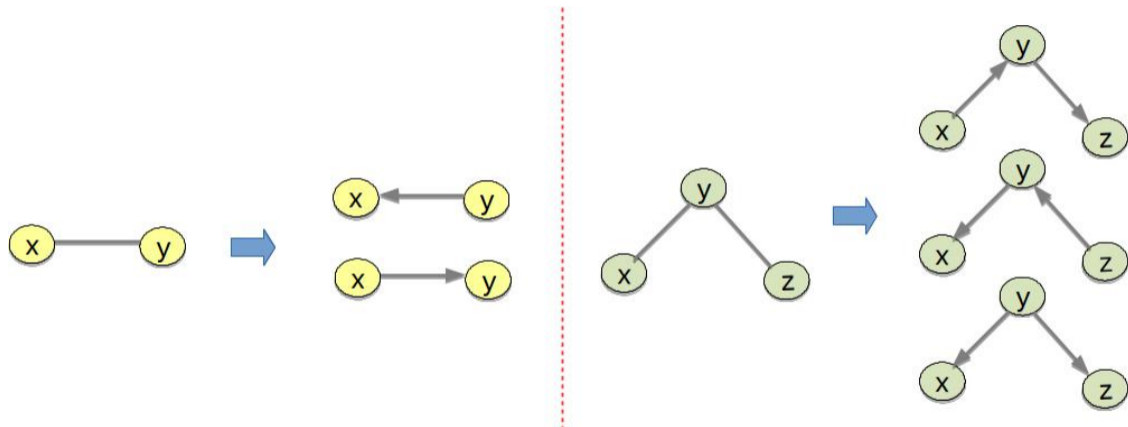
Learning a DAG can be separated into two steps:



- 1 Learn conditional independencies (learn dag skeleton and colliders)
- 2 Learn directions (learn undetermined causal relations)

Work presented here focuses on second part of learning process.

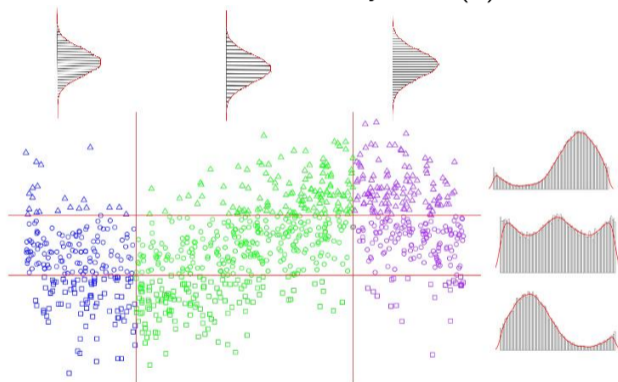
## Our task for two and three variable examples



- Given that we know  $x$  and  $y$  dependent ( $x \not\perp\!\!\!\perp y$ ): choose between  $x \rightarrow y$  or  $y \rightarrow x$
- Given that we know  $x$  and  $z$  conditionally independent given  $y$  ( $x \perp\!\!\!\perp z | y$ ): choose between  $x \rightarrow y \rightarrow z$ ,  $x \leftarrow y \leftarrow z$  or  $x \leftarrow y \rightarrow z$ .

## Idea behind KCDC

Following figure shows observations from model  $y = \sin(x) + n$  where  $n \sim N(0, 1)$



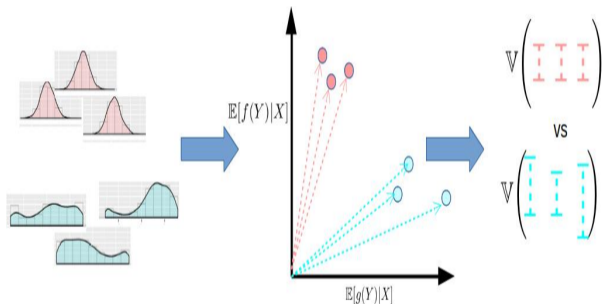
In causal direction ( $x \rightarrow y$ ) complexity of  $p(y|x)$  does not depend on  $x$  whereas in anticausal direction ( $y \rightarrow x$ ) complexity of  $p(x|y)$  varies more.



## How do we measure complexity?

Use the *spread* of the norm of vector of expected features as a proxy for the complexity of a set of distributions.

$$\mu = \mathbb{E} \phi(Y) = \mathbb{E} \begin{bmatrix} f(Y) \\ Y^2 \\ \sin(Y) \\ \vdots \\ g(Y) \end{bmatrix}$$

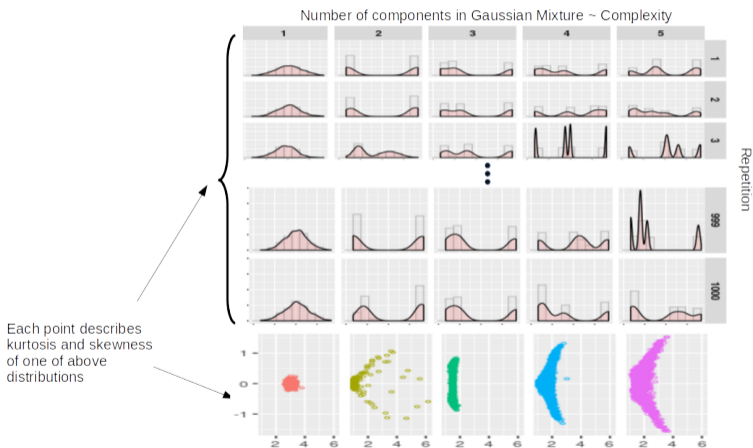


Intuition:

- 1 Expected feature vector represents distribution if adequate features chosen.
- 2 Expected feature vector of similar distributions constrained to subspace of feature space and so have similar norms.

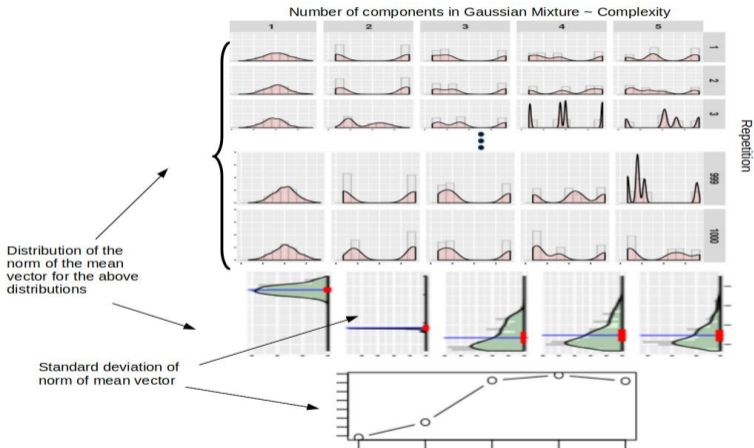
## Illustration: gaussian mixtures

The higher the number of components in a gaussian mixture the more complex it is.



Norm of mean vector can help us distinguish between distributions.

## Illustration: gaussian mixtures



In this case we replace 2 *hand-crafted* features with 1000 random fourier features to construct mean vector.

## Kernel Conditional Deviance for Causal Inference (KCDC)

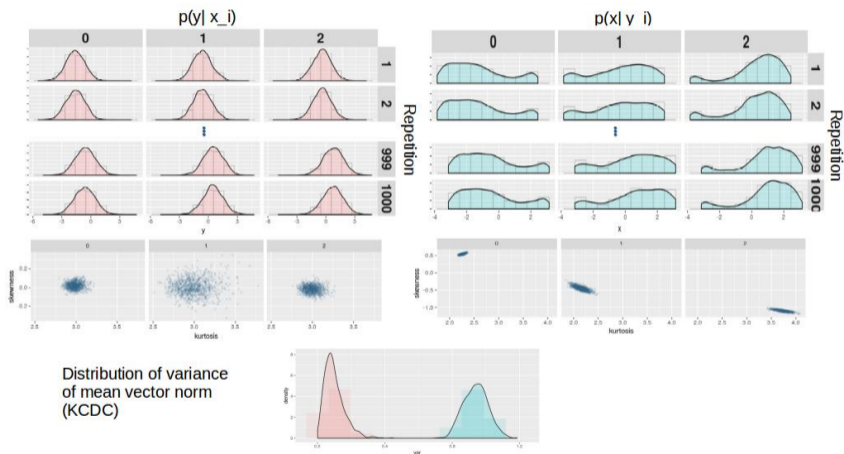
Based on this idea [Mitrovic et al, 2018] introduced KCDC to infer direction of causality for pairs of variables.

## KCDC

$$S_{x \rightarrow y} = \frac{1}{|B|} \sum_{i=1}^{|B|} \left( \|\mu_{y|x \in b_i}\|_2 - \frac{1}{|B|} \sum_{j=1}^{|B|} \|\mu_{y|x \in b_j}\|_2 \right)^2$$

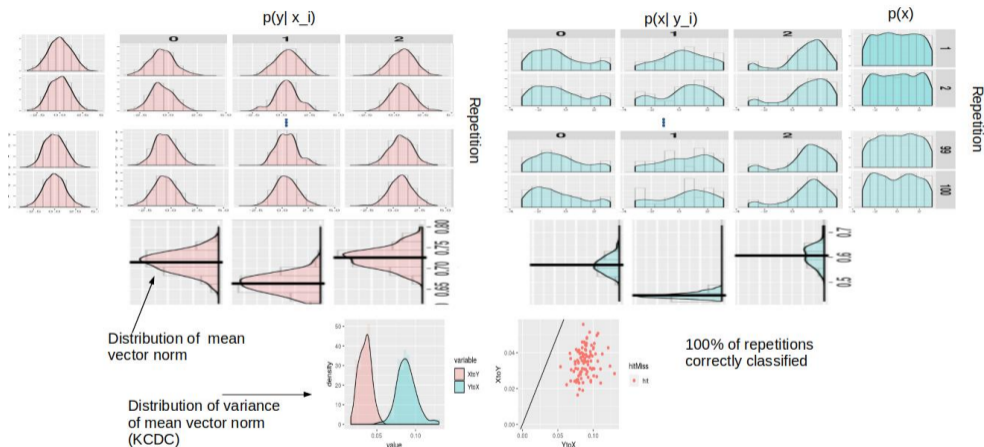
- $B = b_1, \dots, b_m$  are the bins that  $x$  is split into,
- KCDC is the variance of mean feature norms, corresponding to different bins
- Measure in both directions, direction of minimum variance is causal direction.

Back to  $\sin(x) + n$  example...



KCDC distinguishes causal direction for all 1000 repetitions.

Back to  $\sin(x) + n$  example...



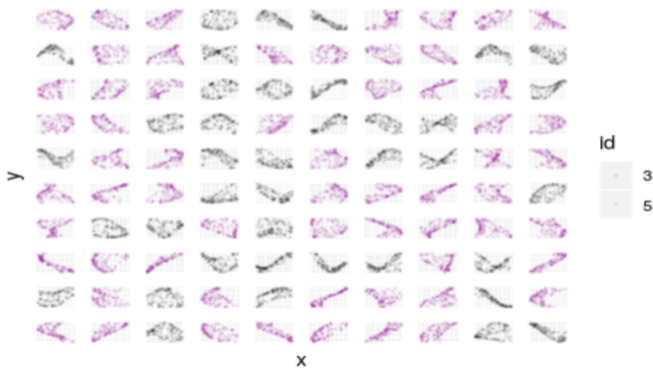
Again, we replace 2 *hand-crafted* features with 1000 random fourier features to construct mean vector.

Up until now we have explored KCDC proposed by [Mitrovic et al, 2018] to infer direction of causality for pairs of variables. Our work involves developing kernel deviance measures further, specifically:

1. Pairs of variables
  - a. Test KCDC on GRS pairs to validate its effectiveness in geosciences,
  - b. Alternate kernel deviance measures: describe complexity of  $p(y|x_i)$  based on  $\hat{\mu}_{Y|x_i}$  in different ways.
  - c. Kernel deviance measures as regularizers (exploit multioutput regression interpretation of CME)
  - d. Improve how  $\hat{\mu}_{Y|x_i}$  represents  $p(y|x_i)$ : new schemes for kernel hyperparameter tuning and/or feature learning.
2. Higher Dimension dags:
  - a. Extend KCDC to multivariate systems of variables,
  - b. Test multivariate KCDC on multivariate simulated datasets.

## Experiment 1: Artificial Cause-Effect Pairs

Note: for this and following experiments kernel parameters and regularization parameter fixed (heuristically).



- 100 data sets with 100 pairs of points each
- Non-additive noise models  $y = f(x, n)$  with:
  - non-linear random function  $f$
  - $x, n \sim U(-3, 3)$

measure	ccr	auc
ANM	66.0 %	70.1%
KCDC	84.0 %	84.9 %



## Cause-Effect Pairs database

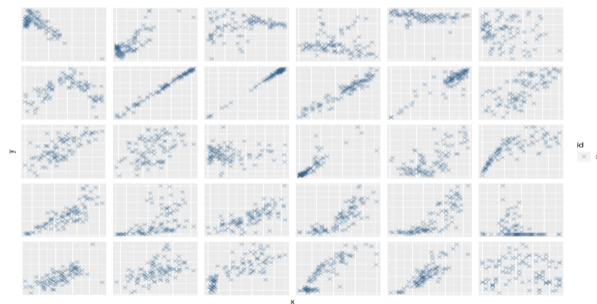
- Cause Effect Pairs (CEP) contains annotated 102 pairs<sup>1</sup>
- Unidimensional and GRS variables only (30 out of 100)

id	x	y	Cause
pair0001	Altitude	Temperature	→
pair0002	Altitude	Precipitation	→
pair0003	Longitude	Temperature	→
pair0004	Altitude	Sunshine hours	→
pair0020	Latitude	Temperature	→
pair0021	Longitude	Precipitation	→
pair0042	Day of the year	Temperature	→
pair0043	Temperature at t	Temperature at t+1	→
pair0044	Pressure at t	Pressure at t+1	→
pair0045	Sea level pressure at t	Sea level pressure at t+1	→
pair0046	Relative humidity at t	Relative humidity at t+1	→
pair0049	Ozone concentration	Temperature	←
pair0050	Ozone concentration	Temperature	←
pair0051	Ozone concentration	Temperature	←
pair0072	Sunspots	Global mean temperature	→

id	x	y	Cause
pair0073	CO2 emissions	Energy use	←
pair0077	Temperature	Solar radiation	←
pair0078	PPFD	Net Ecosystem Productivity	→
pair0079	Net Ecosystem Productivity	Diffuse PPFDdif	←
pair0080	Net Ecosystem Productivity	Diffuse PPFDdif	←
pair0081	Temperature	Local CO2 flux, BE-Bra	→
pair0082	Temperature	Local CO2 flux, DE-Har	→
pair0083	Temperature	Local CO2 flux, US-PFa	→
pair0087	Temperature	Total snow	→
pair0089	root decomposition Oct (grassl)	root decomposition Oct (grassl)	←
pair0090	root decomposition Oct (forest)	root decomposition Oct (forest)	←
pair0091	clay cont. in soil (forest)	soil moisture	→
pair0092	organic carbon in soil (forest)	clay cont. in soil (forest)	←
pair0093	precipitation	runoff	→
pair0094	hour of day	temperature	→

<sup>1</sup><https://webdav.tuebingen.mpg.de/cause-effect/>

## Experiment 2: Cause-Effect Pairs database

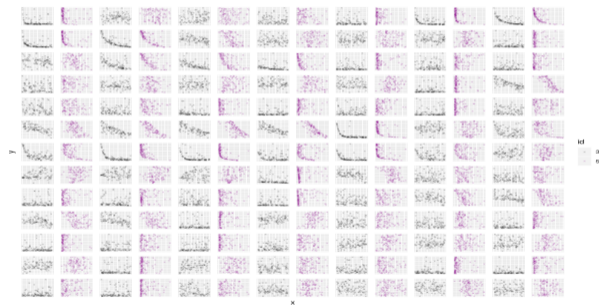


- 30 data sets with 126-10369 pairs of points each
- max 100 points used
- Non-linear, non-additive examples included

measure	ccr	auc
ANM	60.0 %	55.7 %
KCDC	66.7 %	70.2 %
SHSIC	-	70.0 %

- SHSIC result from [Pérez-Suay et al, 2019]

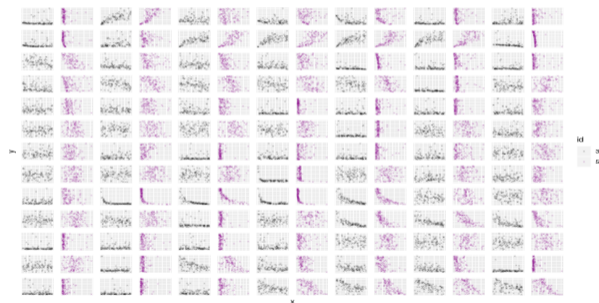
## Experiment 3: RTM Prosail Simulated Pairs



- 182 data sets with 1000 pairs of points each
- max 100 points used
- causes consist of **7** biological parameters
- effects consist of reflectances for **13** different bands

measure	ccr	auc
ANM	62.6 %	60.2 %
KCDC	97.8 %	99.3 %
SHSIC	-	65.0 %

## Experiment 4: RTM Prosail Emulator Pairs

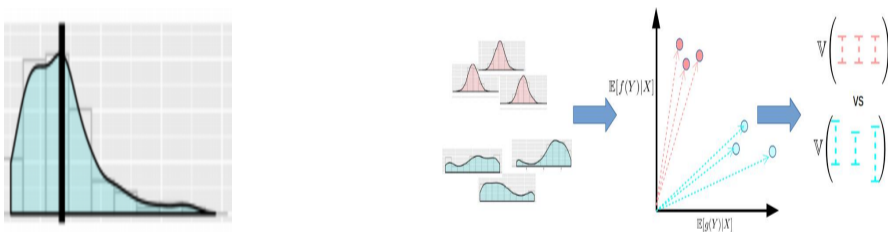


- 182 data sets with 500,000 pairs of points each
- max 100 points used
- causes consist of **7** biological parameters
- effects consist of reflectances for **13** different bands

measure	ccr	auc
ANM	58.8 %	60.4 %
KCDC	97.3 %	99.4 %
SHSIC	-	80.0 %

Alternate kernel deviance measures: recall KCDC and distribution of the norm of CMEs

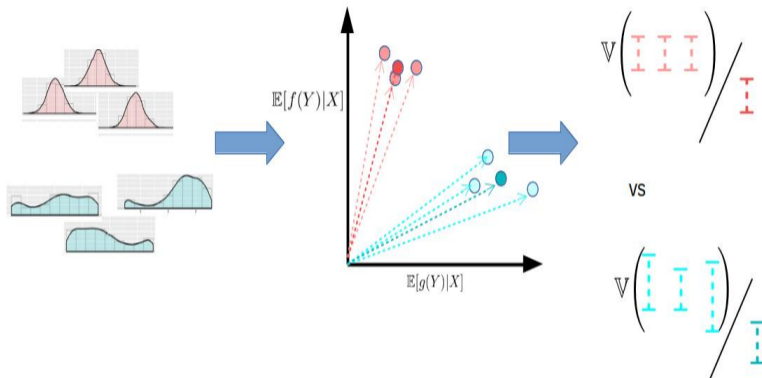
KCDC may have issues of scale for clustered data or data with outliers. Additionally we have seen that distribution of norms of CMEs is often far from gaussian, so variance may not be a good descriptor.



### KCDC

$$S_{x \rightarrow y}^{KCDC} = \hat{V}_X[\|\mu_{Y|X=x}(y)\|_{\mathcal{H}_k}] = \frac{1}{n} \sum_{i=1}^n \left( \|\hat{\mu}_{Y|X=x_i}(y)\|_{\mathcal{H}_k} - \frac{1}{n} \sum_{j=1}^n \|\hat{\mu}_{Y|X=x_j}(y)\|_{\mathcal{H}_k} \right)^2$$

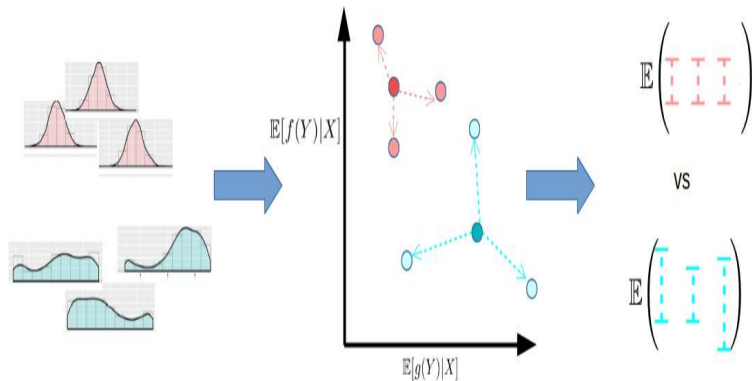
## Alternate kernel deviance measures: KCRDC



## KCRDC

$$S_{x \rightarrow y}^{KCRDC} = \frac{\sqrt{\hat{\mathbb{V}}_X[\|\mu_{Y|X=x}(y)\|_{\mathcal{H}_k}]}}{\hat{\mathbb{E}}_X[\|\mu_{Y|X=x}(y)\|_{\mathcal{H}_k}]} = \frac{\sqrt{S_{x \rightarrow y}^{KCRDC}}}{\hat{\mathbb{E}}_X[\|\mu_{Y|X=x}(y)\|_{\mathcal{H}_k}]}$$

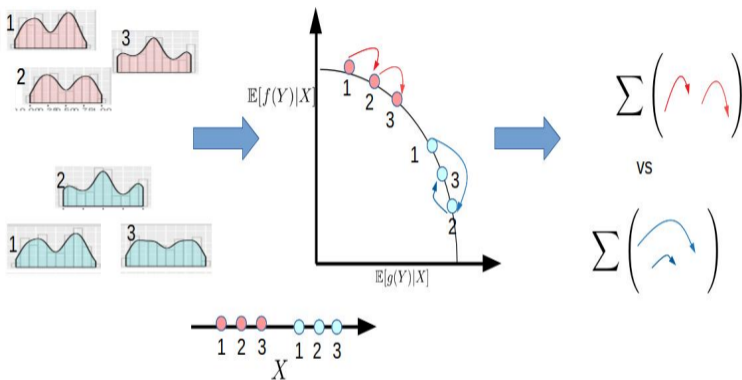
## Alternate kernel deviance measures: KCMC



## KCMC

$$S_{x \rightarrow y}^{KCMC} = \hat{E}_X \left[ \left\| \mu_{Y|X=x}(y) - \mathbb{E}_X [\mu_{Y|X=x}(y)] \right\|_{\mathcal{H}_k}^2 \right] = \frac{1}{n} \sum_{i=1}^n \left\| \hat{\mu}_{Y|X=x_i}(y) - \frac{1}{n} \sum_{j=1}^n \hat{\mu}_{Y|X=x_j}(y) \right\|_{\mathcal{H}_k}^2$$

## Alternate kernel deviance measures: KCSC

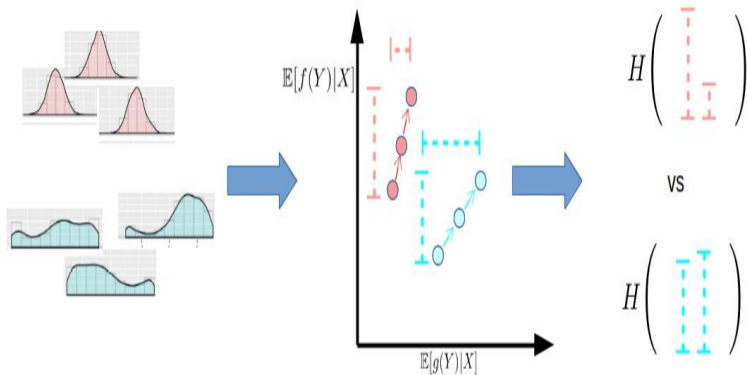


## KCSC

$$S_{x \rightarrow y}^{KCSC} = \mathbb{E}_Y \mathbb{E}_X \left[ \left\| \nabla_x \mu_{Y|X=x}(y) \right\|_{\mathcal{H}_k}^2 \right] \quad (1)$$



## Alternate kernel deviance measures: KCCC



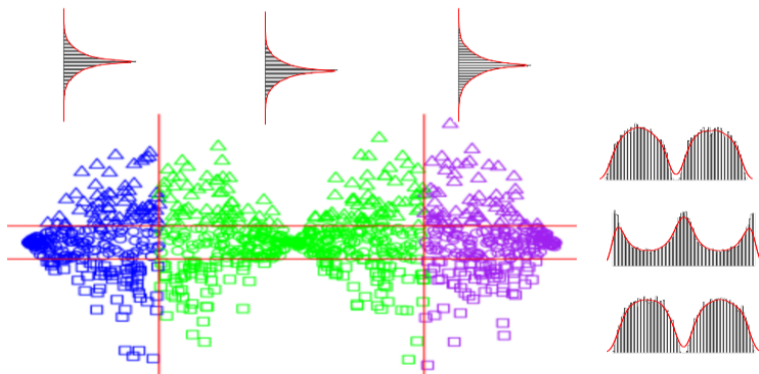
## KCCC

$$R = \int \nabla_x \mu_{Y|X=x} \nabla_x \mu_{Y|X=x}^T dP_x$$

$$S = \text{diag}(R \circ R) \in \mathbb{R}^m$$

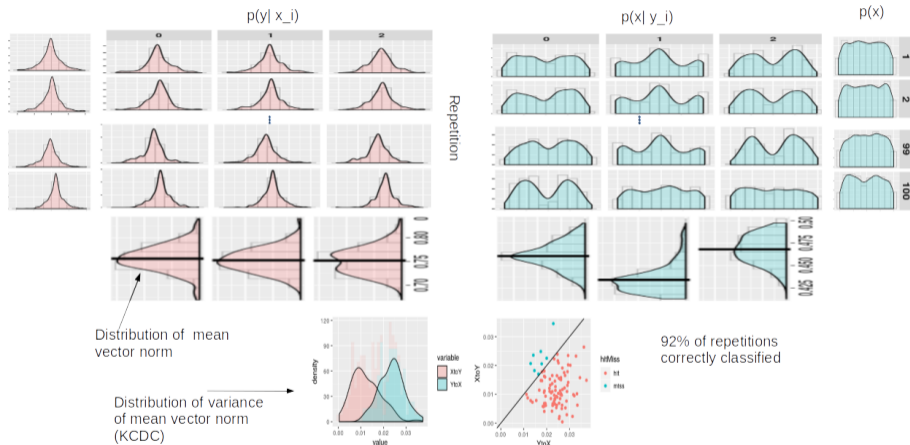
$$S_{x \rightarrow y}^{KCCC} = H\left(\underbrace{\tilde{S}}_{\in \mathbb{R}^m}\right)$$

## Alternate kernel deviance measures: KCMC and KCRDC a toy example



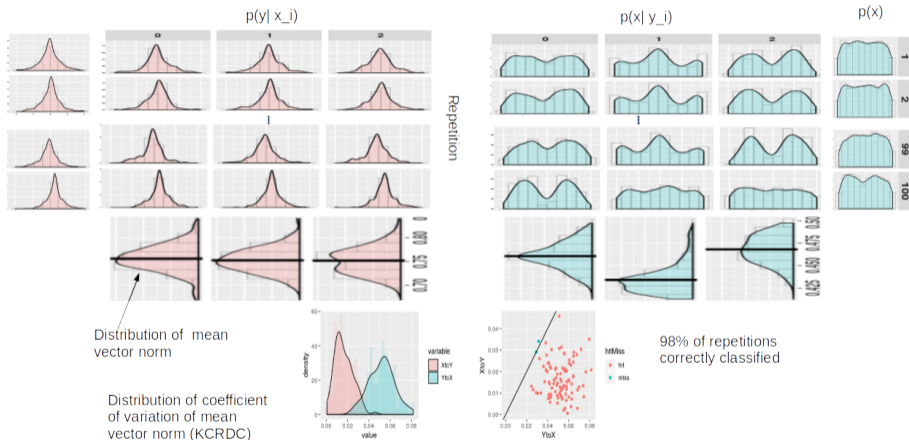
Data from  $y = \sin(x) * n$  with  $n \sim N(0, 1)$

## Alternate kernel deviance measures: KCMC and KCRDC a toy example

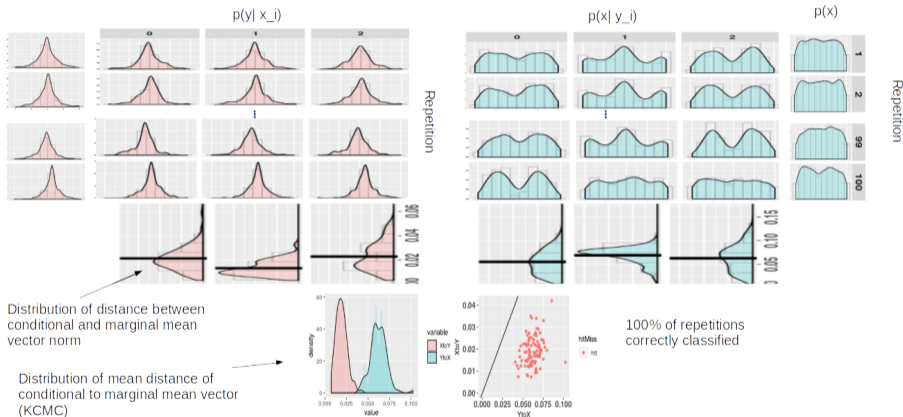


KCDC doesn't do as well as with additive noise in this case. Observe higher mean level for causal distributions.

Alternate kernel deviance measures: KCMC and KCRDC a toy example



# Alternate kernel deviance measures: KCMC and KCRDC a toy example



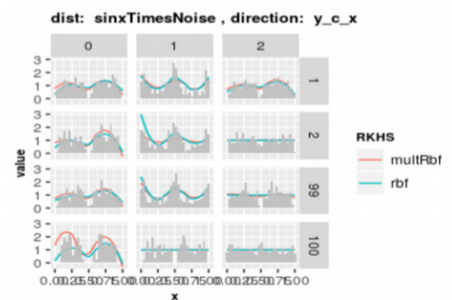
## Kernel deviance measures as regularizers: causal regularizers?

- We can express all kernel deviance measures in terms of the  $\alpha$  kernel regression parameter.
- We can simply replace L2 regularizer and optimize loss.
- In causal direction, kernel deviance restriction should not hurt accuracy too much? If so we can use regression loss as causal predictor.
- For KCMC and KCSC we obtain a convex loss function!
- No experiments yet

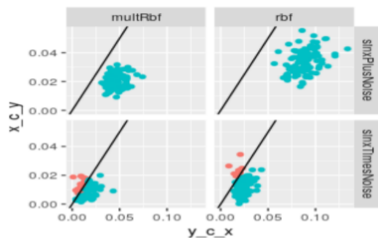
$$\begin{aligned}\mathcal{L}^{KCMC}(\mathbf{x}, \mathbf{y}, \alpha, \lambda) &= (\Phi_y \Phi_y^T) - 2(K_x \alpha \Phi_y^T) + (\alpha^T K_x K_x \alpha) \\ &\quad + \lambda n \left( \frac{1}{n} (K_x \alpha \alpha^T K_x) - \frac{1}{n^2} \mathbf{1}^T K_x \alpha \alpha^T K_x \mathbf{1} \right) \\ \nabla_{\alpha} \mathcal{L}^{KCMC}(\mathbf{x}, \mathbf{y}, \alpha, \lambda) = 0 &\Rightarrow \hat{\alpha}^{KCMC} = (K_x (I + \lambda H))^{-1} \Phi_y\end{aligned}\tag{2}$$

Choosing/learning  $\mathcal{H}_y$ 

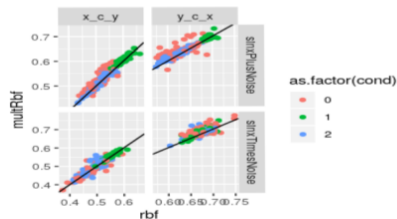
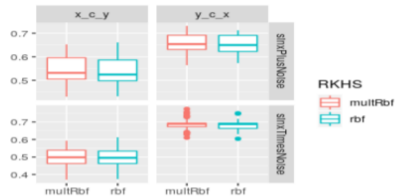
Squared error across different  $\mathcal{H}_y$  not comparable. Work so far here at Oxford has involved using Noise Contrastive Estimation as in [J. Ton et al 2019] to choose appropriate  $\mathcal{H}_y$ .



Example of Noise Contrastive Estimation for  $y = \sin(x) * n$  in anti-causal direction.

Choosing/learning  $\mathcal{H}_y$ 

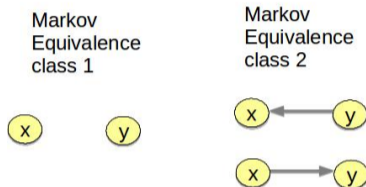
	RKHS	
dist	multRbf	rbf
sinxPlusNoise	100	100
sinxTimesNoise	73	90



For the  $y = \sin(x) * n$  example the loss using the rbf  $\mathcal{H}_y$  is lower than if we use a direct sum rbf  $\mathcal{H}_y$ . This corresponds with a better KCDC performance suggesting we can use NCE to supervise learning/tuning of  $\mathcal{H}_y$ .



## Extending KCDC to systems with more than two variables



To extend KCDC to DAGs with more than two nodes (higher dimensional systems) we note that:

- KCDC only serves to distinguish between DAGs in the same Markov Equivalence class (those graphs with same set of conditional independencies).
- The distribution of nodes with no parents is not taken into account since the causal mechanism is encoded in the conditional distributions of nodes with parents.

## Extending KCDC to systems with more than two variables

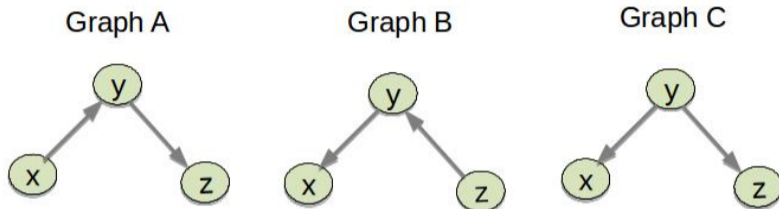
Taking this into account we write the KCDC of a general p-node DAG as:

$$KCDC(\mathcal{G}) = \sum_{i \in \mathcal{A}} KCDC\left(p(x_i | pa(x_i))\right) \quad (3)$$

where

- $\mathcal{A}$  is the set of nodes in the dag  $\mathcal{G}$  that have at least one parent, and
- $pa(x_i)$  is the set of parents of node  $x_i$ .

## An example

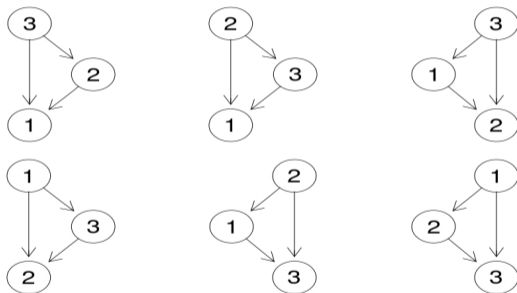


With previous definition:

- $KCDC(\mathcal{G}_A) = KCDC(p(y|x)) + KCDC(p(z|y))$
- $KCDC(\mathcal{G}_B) = KCDC(p(x|y)) + KCDC(p(y|z))$
- $KCDC(\mathcal{G}_C) = KCDC(p(x|y)) + KCDC(p(z|y))$

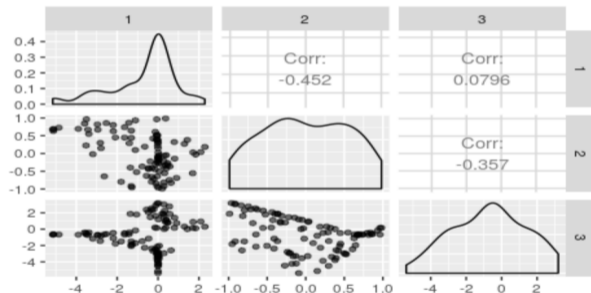
Lets see some experimental results for multi-variate KCDC.

## Experiment 5: Artificial Cause-Effect 3-tuples



- 100 datasets with 100 3-tuples each
- Non-additive noise models  $z = f(x, y, n)$  with:
  - non-linear random function  $f$
  - $x, y, n \sim U(-1, 1)$
- true causal structure one of 6 dags on the left.

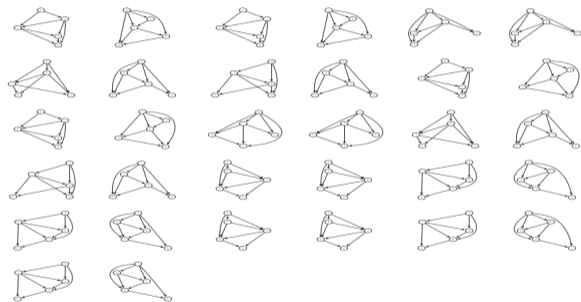
## Experiment 5: Artificial Cause-Effect 3-tuples



- Data for 1 of 100 datasets plotted on left.

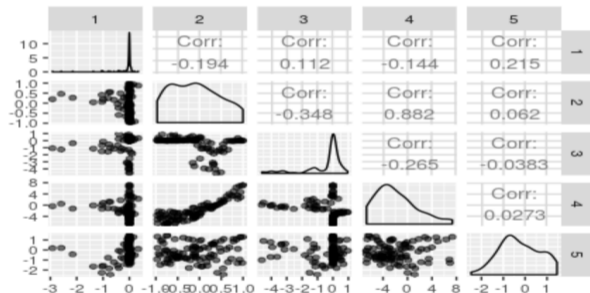
measure	ccr	edgeCCR
ANM	18.0 %	48.0 %
KCDC	38.0 %	67.0 %
Rnd	23.0 %	55.0 %

## Experiment 5: Artificial Cause-Effect 5-tuples



- 100 datasets with 100 5-tuples each
- Non-additive noise models  $e = f(a, b, c, n, )$  with:
  - non-linear random function  $f$
  - $a, b, c, d, n \sim U(-1, 1)$
- true causal structure one of 32 dags on the left.

## Experiment 5: Artificial Cause-Effect 5-tuples



- Data for 1 of 100 datasets plotted on left.

measure	ccr	edgeCCR
ANM	8.0 %	67.3 %
KCDC	23.0 %	80 %
Rnd	6.0 %	63.3 %

## References

- 📄 Mitrovic et al., 'Causal Inference via Kernel Deviance Measures,' NIPS 2018.
- 📄 Ton et al., 'Noise contrastative estimation for meta learning using kernel mean embeddings' AISTATS 2019 submitted.
- 📄 Pérez-Suay and Camps-Valls, 'Causal Inference in Geoscience and Remote Sensing From Observational Data,' IIIE-GRSL, 2019.
- 📄 Pérez-Suay and Camps-Valls, 'Sensitivity Maps of the Hilbert-Schmidt Independence Criterion,' Applied Soft Computing, 2017.
- 📄 Mooij et al., 'Distinguishing cause from effect using observational data,' JMLR, 17(1), 2016.
- 📄 Gretton et al., 'Measuring statistical dependence with Hilbert-Schmidt norms,' ALT, 2005
- 📄 Hoyer et al., 'Nonlinear causal discovery with additive noise models,' NIPS 2008.
- 📄 Camps-Valls, Mooij and Schölkopf, 'Remote sensing feature selection by kernel dependence measures,' IEEE-GRSL 2010