

A clustering framework for conditional extremes models

Environmental and Ecological Statistics Conference 2025 Talk,
1st June, 2025,

Paddy O'Toole, University of Bath

Supervised by Christian Rohrbeck and Jordan Richards (University of Edinburgh)

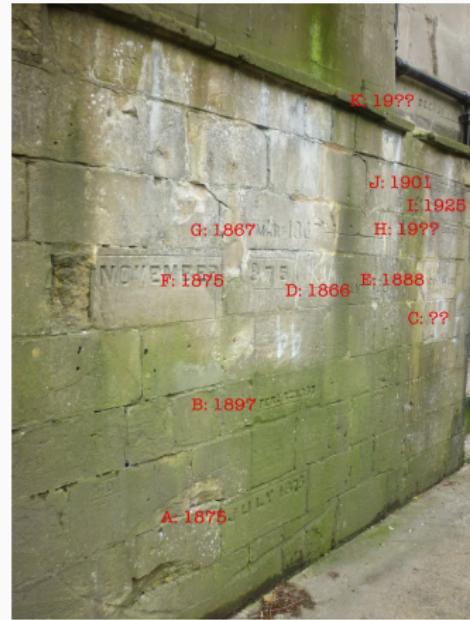


UNIVERSITY OF
BATH



**UK Research
and Innovation**

Introduction



Problem

- Often want to estimate

$$\mathbb{P}(X > x, Y > y) = \mathbb{P}(Y > y | X > x)\mathbb{P}(X > x)$$

for large x, y

- “concomitant”/concurrent extreme events for random vector \mathbf{X} often particularly devastating
- Goal: identify trends by clustering sites with similar **tail dependence**

Dependence Modelling

Asymptotic dependence

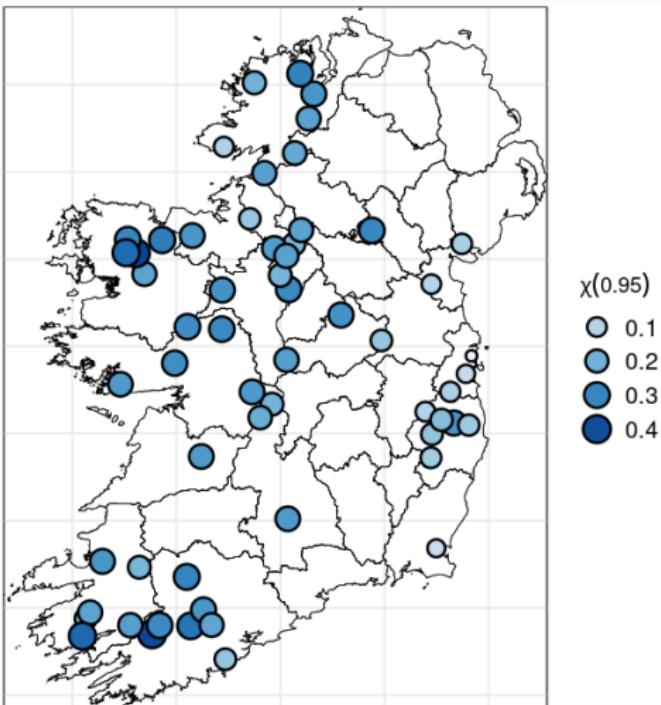
- Coefficient of extremal dependence $\chi \in [0, 1]$,

$$\chi = \lim_{u \rightarrow 1} \mathbb{P}[F_1(X_1) > u \mid F_2(X_2) > u]$$

- (Increasingly strong) asymptotic dependence for $\chi > 0$.
- However, χ only gives summary; inference requires **dependence model**.

Ireland

- Precipitation¹ & wind speed² data for 59 sites across Ireland, Winter months (Oct-Mar) 1990-2020

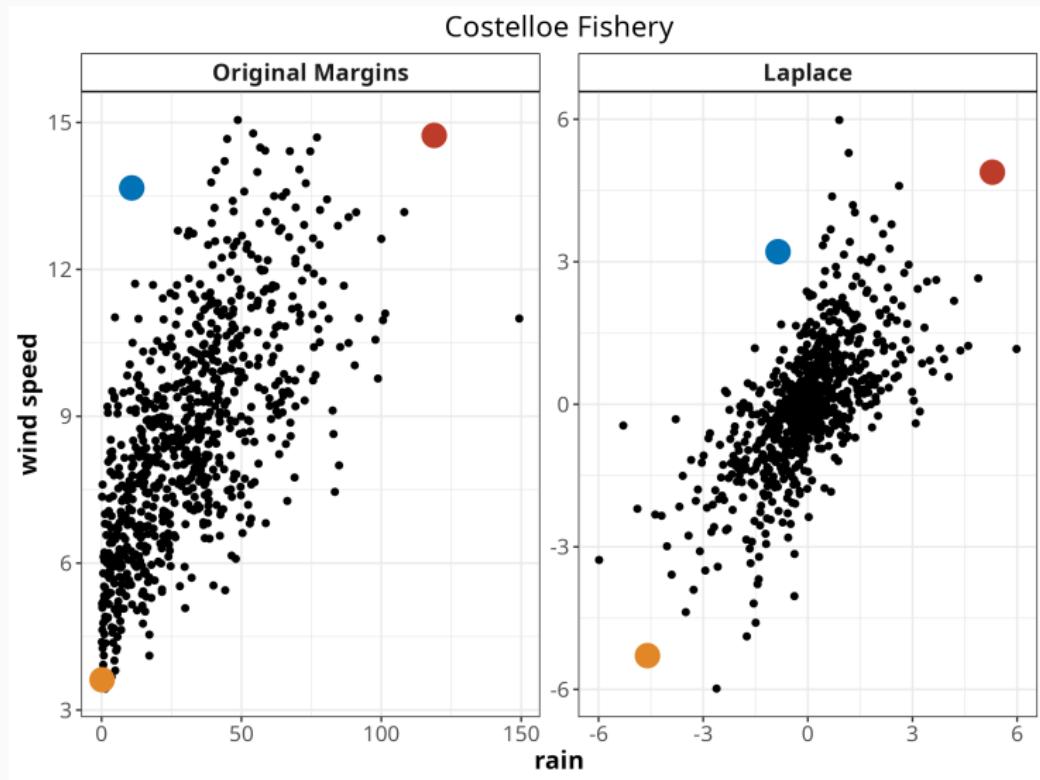


¹ Met Éireann weekly aggregate

² ERA5 reanalysis weekly mean of daily maxima

Conditional extremes

Marginal transformation



Conditional extremes

- Heteroskedastic regression dependence model:

$$(Y | X = x) = \alpha x + x^\beta Z, \text{ for } x > u$$

- slope parameter α for Y given large Y ,
- “spread” parameter $\beta \in (-\infty, 1]$ controls stochasticity of relationship between Y and large Y .

Conditional extremes

- Heteroskedastic regression dependence model:

$$(Y_{-i} \mid Y_i = y_i) = \alpha_{j|i} y_i + y_i^{\beta_{j|i}} Z_{|i}, \text{ for } y_i > u_{Y_i}$$

- slope parameter $\alpha_{j|i} \in [-1, 1]$ for Y_j given large Y_i ,
- “spread” parameter $\beta_{j|i} \in (-\infty, 1]$ controls stochasticity of relationship between Y_j and large Y_i .

Conditional extremes

- Heteroskedastic regression dependence model:

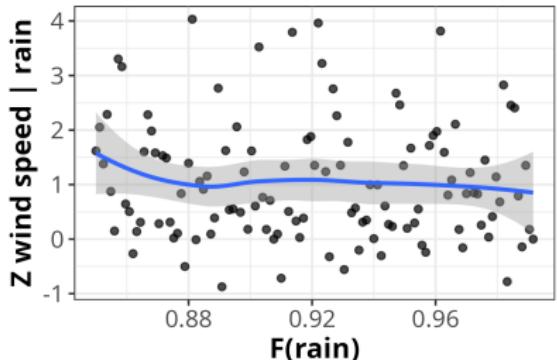
$$(Y_{-i} \mid Y_i = y_i) = \alpha_{|i} y_i + y_i^{\beta_{|i}} Z_{|i}, \text{ for } y_i > u_{Y_i}$$

- slope parameter $\alpha_{j|i} \in [-1, 1]$ for Y_j given large Y_i ,
- “spread” parameter $\beta_{j|i} \in (-\infty, 1]$ controls stochasticity of relationship between Y_j and large Y_i .
- Key assumptions:
 - Residuals $Z_{|i} \sim N(\mu_{|i}, \Sigma_{|i})$
 - $Z_{|i}, Y_i$ conditionally independent for large Y_i
- Special cases:
 - $\alpha_{|i} = 0, \beta_{|i} = 0 \implies Y_{-i}, Y_i$ independent,
 - $\alpha_{|i} = -1/1, \beta_{|i} = 0 \implies$ perfect positive/negative dependence,
 - $-1 < \alpha_{|i} < 1 \implies$ asymptotic independence

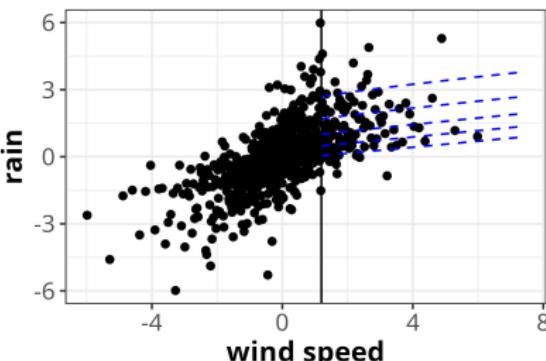
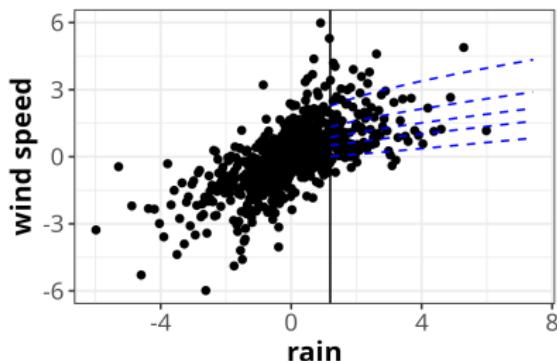
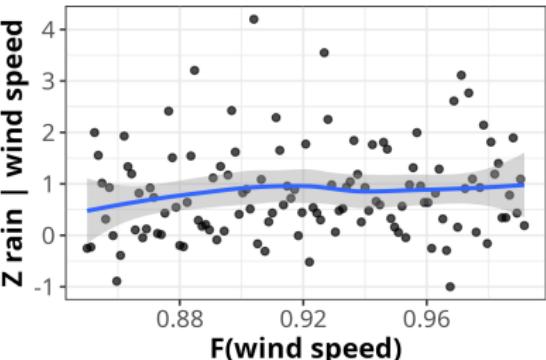
Conditional extremes

Costelloe Fishery

$$a = 0.137, b = 0.052$$



$$a = 0.147, b = 0.243$$



Inference

Inference assumes conditional distribution follows diagonal multivariate Normal (MVN) distribution:

$$(Y_{-i} \mid Y_i = y_i) \sim N \left(\boldsymbol{\alpha}_{|i} y_i + y_i^{\beta_i} \boldsymbol{\mu}_{|i}, y_i^{\beta_i} \boldsymbol{\Sigma}_{|i} \right), \text{ for } Y_i > u_{Y_i}$$

⇒ dependence structures at different sites can be compared using their MVN distributions

Clustering

skew-geometric Jensen-Shannon divergence

- Kullback-Leibler divergence $KL(X \parallel Y)$ measures (**asymmetric**) distance from X to Y
- $KL(X \parallel Y)$ has closed form for two MVNs
- $KL(X \parallel Y) + KL(Y \parallel X)$ symmetric but does not satisfy triangle inequality \implies **not a true metric.**

skew-geometric Jensen-Shannon divergence

- Kullback-Leibler divergence $KL(X \parallel Y)$ measures (**asymmetric**) distance from X to Y
- $KL(X \parallel Y)$ has closed form for two MVNs
- $KL(X \parallel Y) + KL(Y \parallel X)$ symmetric but does not satisfy triangle inequality \implies **not a true metric.**
- Weighted geometric mean $G_\alpha(X, Y) = X^\alpha Y^{1-\alpha}, \alpha \in [0, 1]$

skew-geometric Jensen-Shannon divergence

- Kullback-Leibler divergence $KL(X \parallel Y)$ measures (**asymmetric**) distance from X to Y
- $KL(X \parallel Y)$ has closed form for two MVNs
- $KL(X \parallel Y) + KL(Y \parallel X)$ symmetric but does not satisfy triangle inequality \implies **not a true metric.**
- Weighted geometric mean $G_\alpha(X, Y) = X^\alpha Y^{1-\alpha}, \alpha \in [0, 1]$
- weighted product G_α of two exponential family members is in **same family**.

skew-geometric Jensen-Shannon divergence

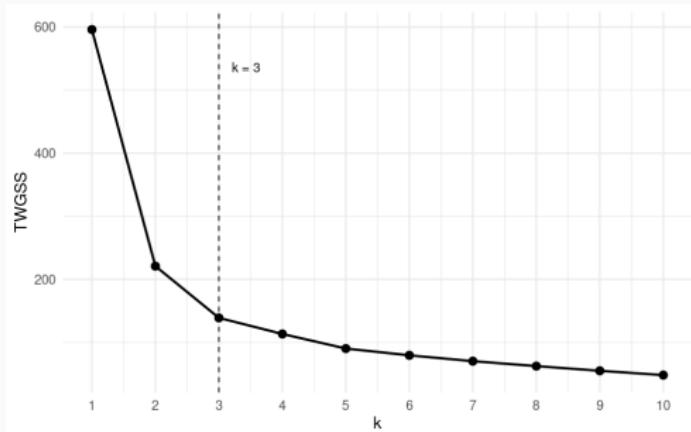
- Kullback-Leibler divergence $KL(X \parallel Y)$ measures (**asymmetric**) distance from X to Y
- $KL(X \parallel Y)$ has closed form for two MVNs
- $KL(X \parallel Y) + KL(Y \parallel X)$ symmetric but does not satisfy triangle inequality \implies **not a true metric.**
- Weighted geometric mean $G_\alpha(X, Y) = X^\alpha Y^{1-\alpha}$, $\alpha \in [0, 1]$
- weighted product G_α of two exponential family members is in **same family**.
- Skew-geometric Jensen-Shannon divergence:

$$JS^{G_\alpha}(X \parallel Y) = \frac{1}{2} \{KL(X \parallel G_\alpha(X, Y)) + KL(Y \parallel G_\alpha(X, Y))\}$$

Clustering

- Uses **k-medoids** to cluster over JS_{G_α} dissimilarity matrix between sites
- **Elbow plot:** k chosen using **Total Within Group Sum of Squares** between sites x within clusters C_i with medoids m_i :

$$\text{TWGSS}(k) = \sum_{i=1}^k \sum_{x \in C_i} JS_{G_\alpha}(x || m_i)$$

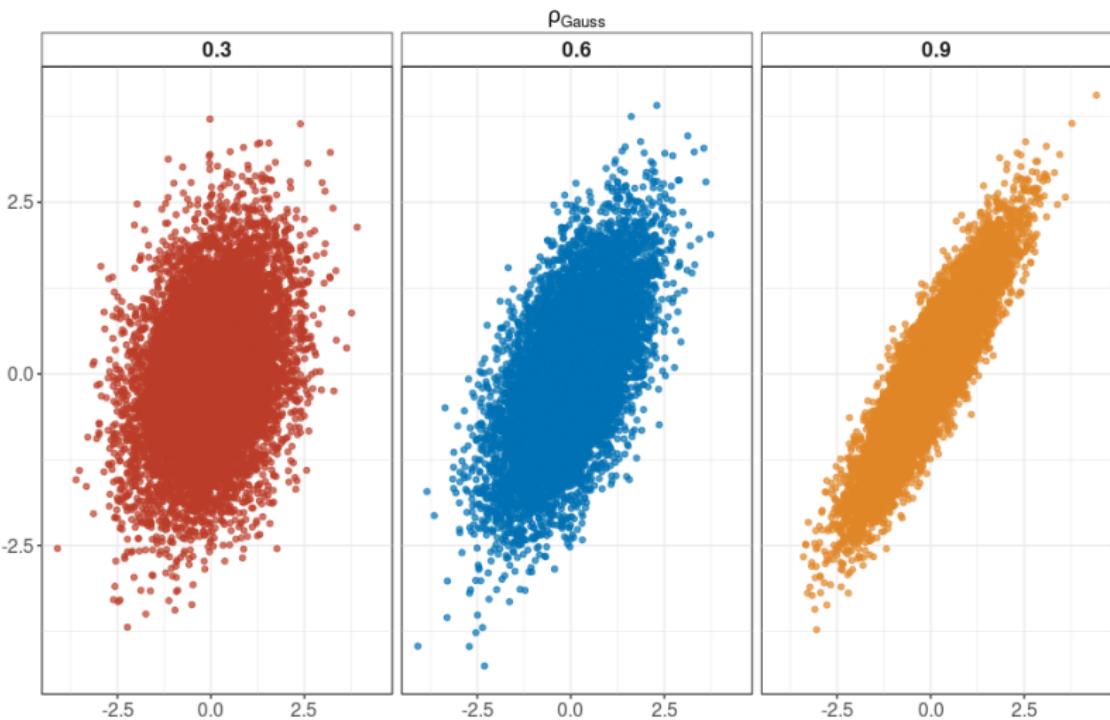


Simulations

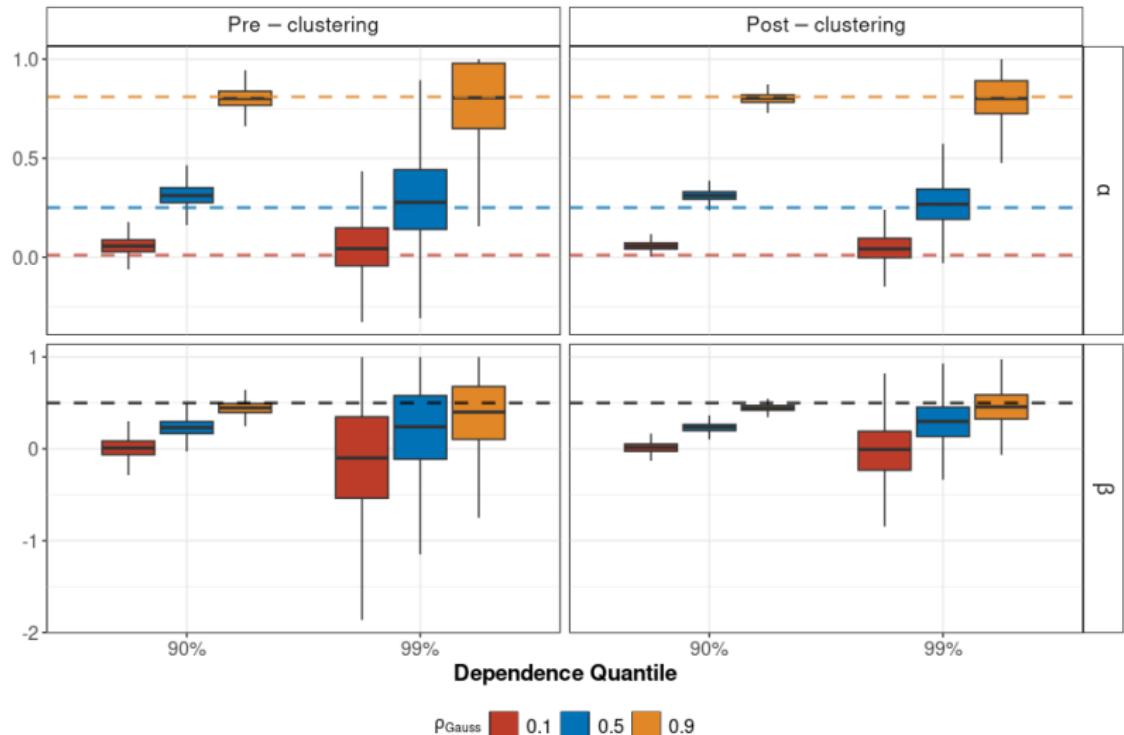
Gaussian copula

- For single “site”, generate data from bivariate Gaussian copula with correlation ρ_{Gauss} .
- (asymptotic) CE parameters are $\alpha = \rho_{\text{Gauss}}^2$, $\beta = 1/2$
- Example: 12 “sites” with 10,000 observations of 2 variables
- 3 clusters of 4 locations defined by respective ρ_{Gauss} values of 0.3, 0.6, 0.9.

Gaussian copula



Gaussian copula - results



Mixture simulation

- Extend design to mixture of Gaussian and t-copulas

Mixture simulation

- Extend design to mixture of Gaussian and t-copulas
- Idea:
 - Gaussian copula generates observations exhibiting extremal independence,
 - t-copula induces extremal dependence.

Mixture simulation

- Extend design to mixture of Gaussian and t-copulas
- Idea:
 - Gaussian copula generates observations exhibiting extremal independence,
 - t-copula induces extremal dependence.
- Design: 12 sites with 10,000 bivariate observations
- 2 clusters of 6 sites

Mixture simulation

- Extend design to mixture of Gaussian and t-copulas
- Idea:
 - Gaussian copula generates observations exhibiting extremal independence,
 - t-copula induces extremal dependence.
- Design: 12 sites with 10,000 bivariate observations
- 2 clusters of 6 sites
- Grid search performed over $\rho_{\text{Gauss}}, \rho_{t_1}, \rho_{t_2} \in [0, 1]$.

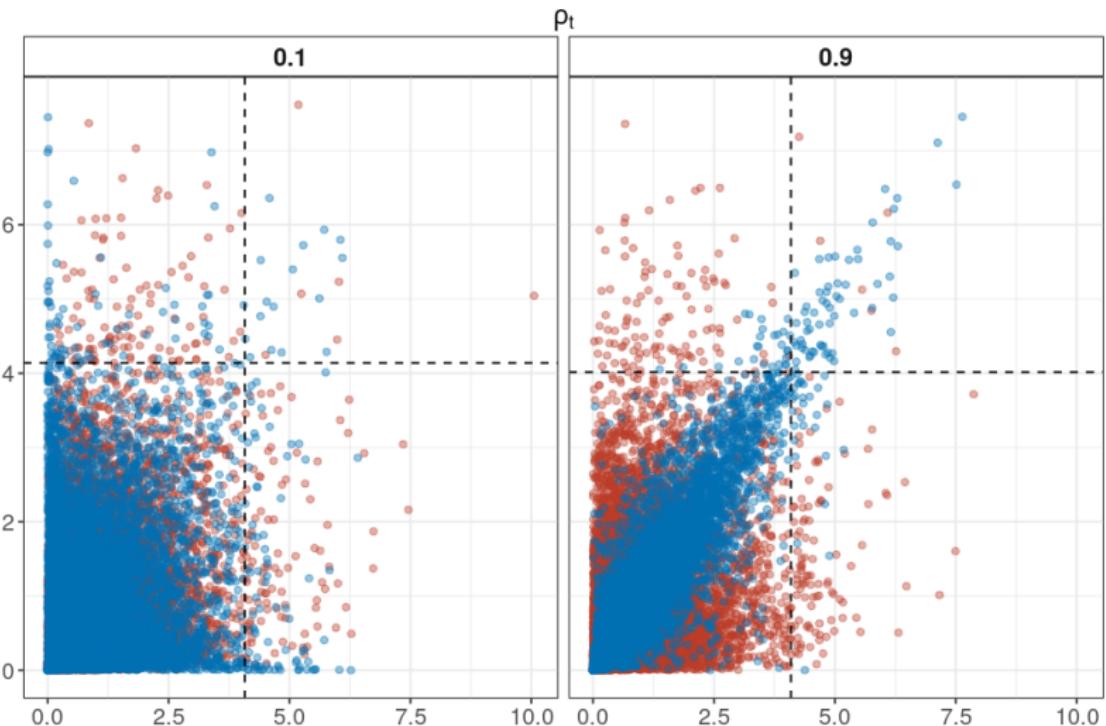
Mixture simulation

- Extend design to mixture of Gaussian and t-copulas
- Idea:
 - Gaussian copula generates observations exhibiting extremal independence,
 - t-copula induces extremal dependence.
- Design: 12 sites with 10,000 bivariate observations
- 2 clusters of 6 sites
- Grid search performed over $\rho_{\text{Gauss}}, \rho_{t_1}, \rho_{t_2} \in [0, 1]$.
- Clustering compared to competing Vignotto algorithm using Adjusted Rand Index ARI $\in [0, 1]$

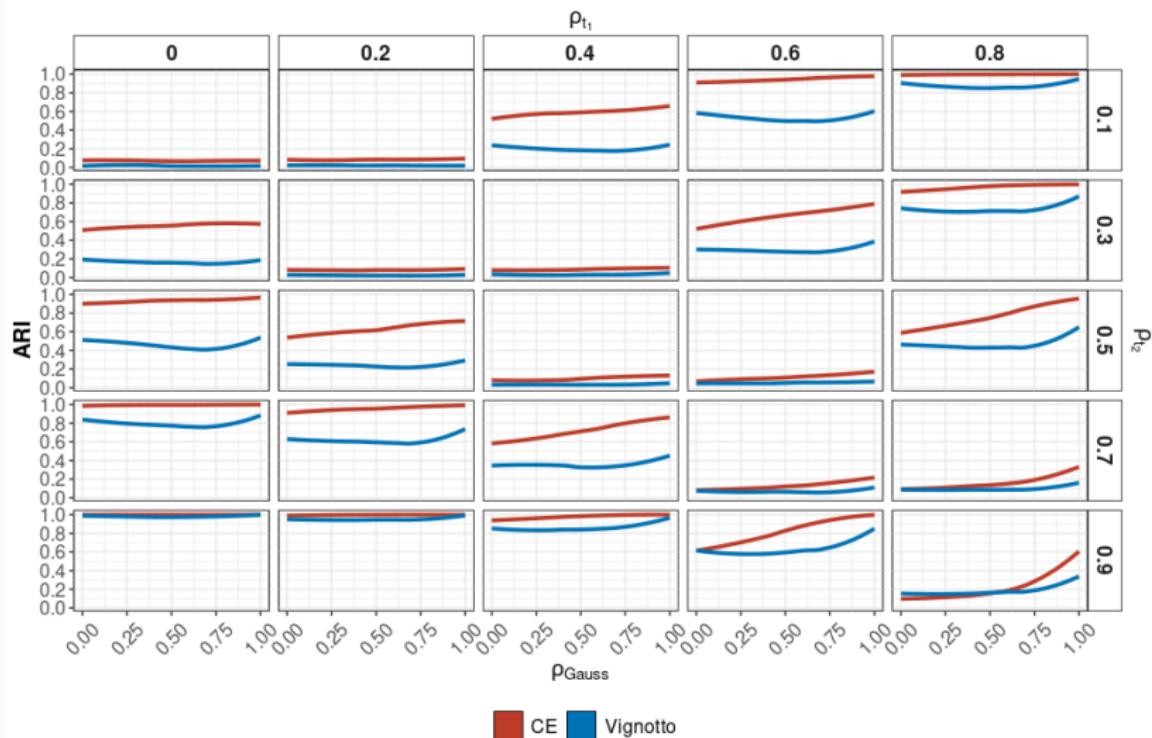
Vignotto, Engelke, and Zscheischler,

“Clustering bivariate dependencies of compound precipitation and wind extremes over Great Britain and Ireland” (2021)

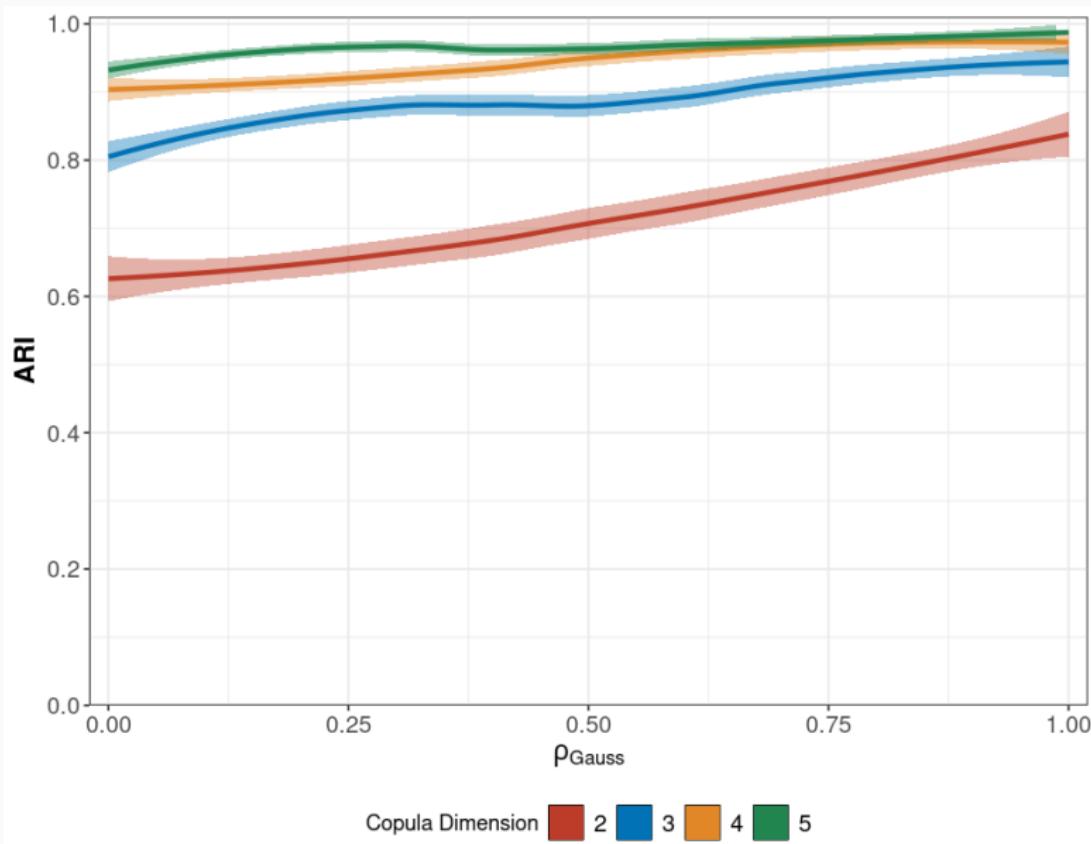
Mixture simulation



Mixture simulation

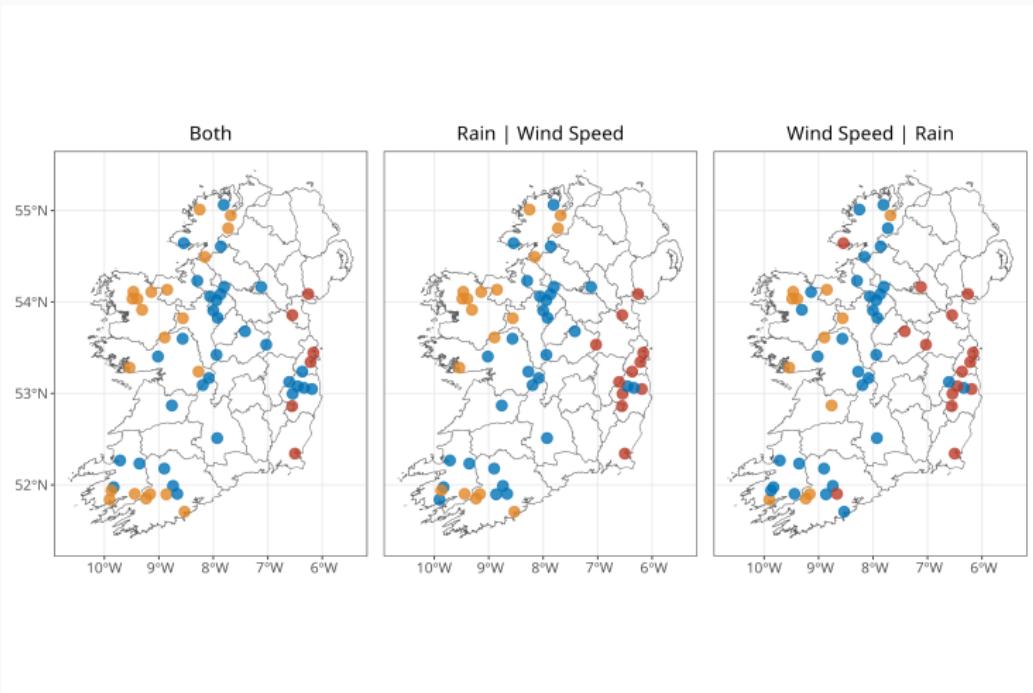


Mixture simulation



Ireland meteorological data

Ireland clustering



Discussion

Discussion

- Conclusion:
 - Principled & effective clustering framework for CE models
 - Simulations & application show clustering can group sites with similar tail dependence, which may aid interpretation

Discussion

- Conclusion:
 - Principled & effective clustering framework for CE models
 - Simulations & application show clustering can group sites with similar tail dependence, which may aid interpretation
- Limitations:
 - CE: Gaussian assumption for $Z_{|i}$
 - Clustering: Uncertainty in CE fits is ignored in clustering

Thank you! For slides and supplementary materials, please scan the QR code below:



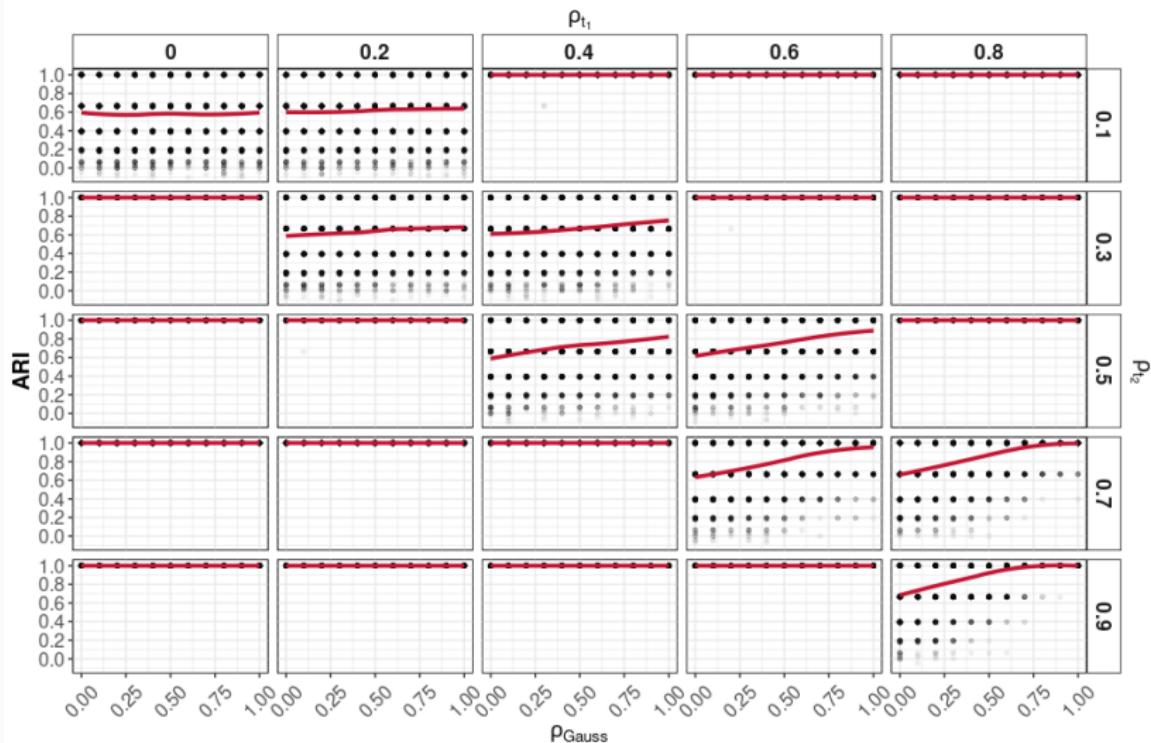
Email: *pot23@bath.ac.uk*

References

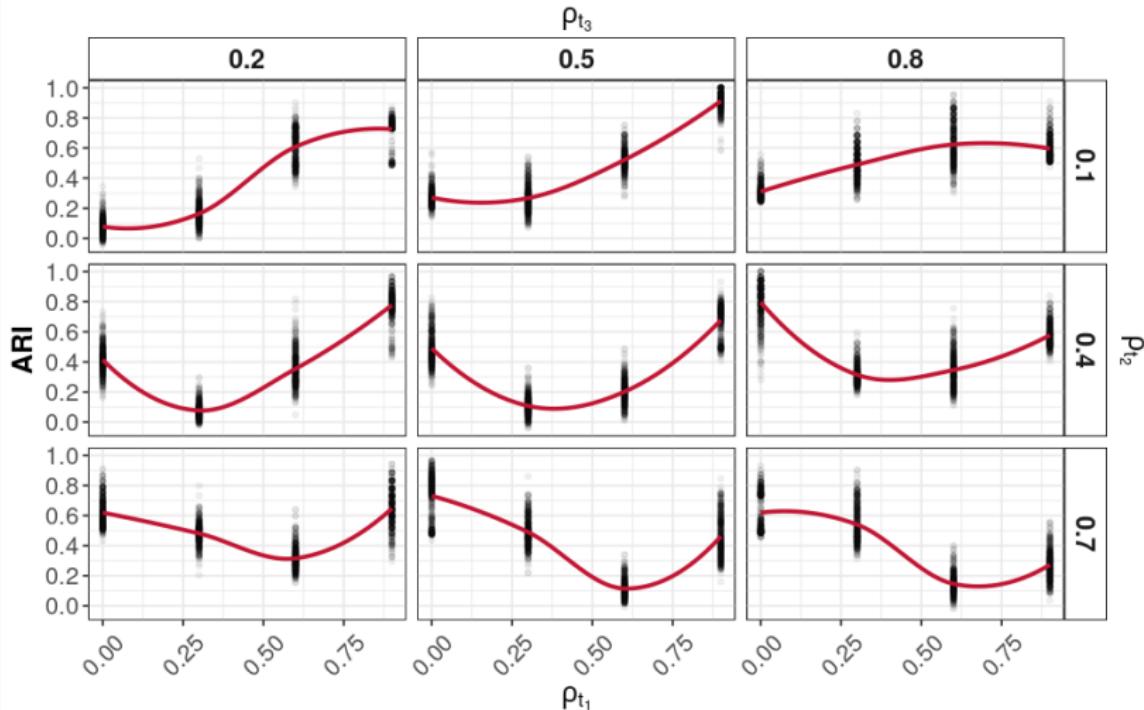
-  Coles, Stuart, Janet Heffernan, and Jonathan Tawn (1999). "Dependence measures for extreme value analyses". In: *Extremes (Boston)* 2.4, pp. 339–365.
-  Heffernan, Janet E. and Jonathan A. Tawn (July 2004). "A Conditional Approach for Multivariate Extreme Values (with Discussion)". In: *Journal of the Royal Statistical Society Series B: Statistical Methodology* 66.3, pp. 497–546. ISSN: 1467-9868. DOI: [10.1111/j.1467-9868.2004.02050.x](https://doi.org/10.1111/j.1467-9868.2004.02050.x). URL: <http://dx.doi.org/10.1111/j.1467-9868.2004.02050.x>.
-  Nielsen, Frank (May 2019). "On the Jensen–Shannon Symmetrization of Distances Relying on Abstract Means". In: *Entropy* 21.5, p. 485. ISSN: 1099-4300. DOI: [10.3390/e21050485](https://doi.org/10.3390/e21050485). URL: <http://dx.doi.org/10.3390/e21050485>.
-  Vignotto, Edoardo, Sebastian Engelke, and Jakob Zscheischler (2021). "Clustering bivariate dependencies of compound precipitation and wind extremes over Great Britain and Ireland". In: *Weather and Climate Extremes* 32, p. 100318.

Appendix

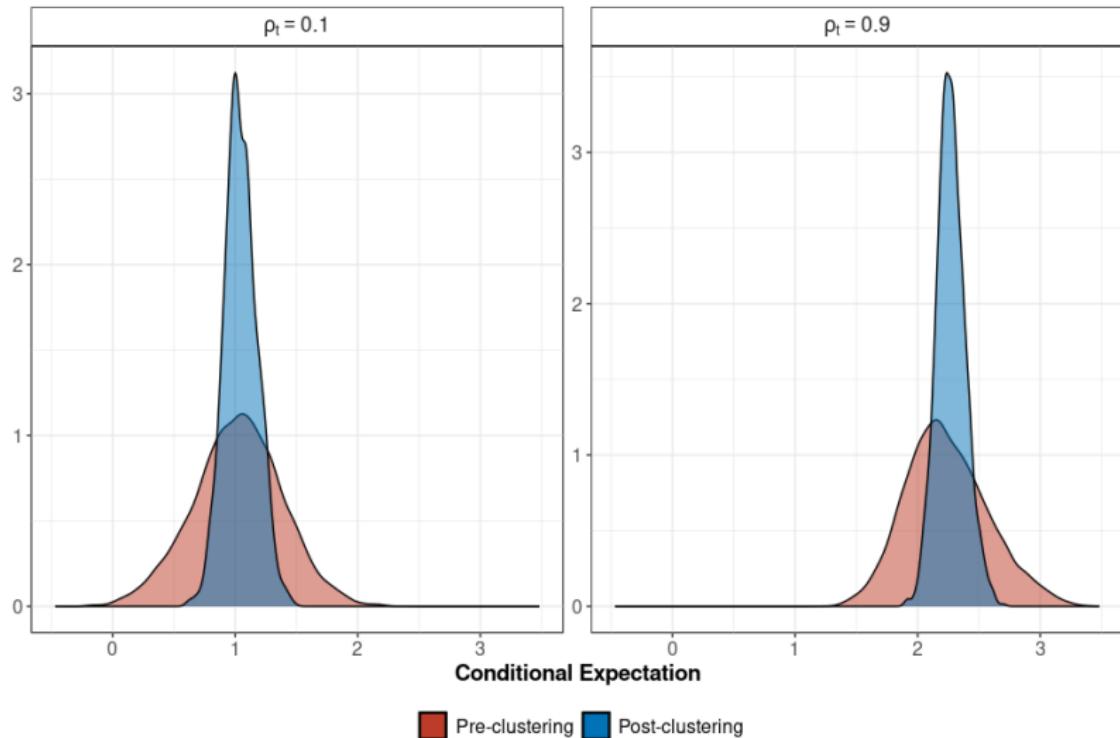
Mixture simulation (three variables)



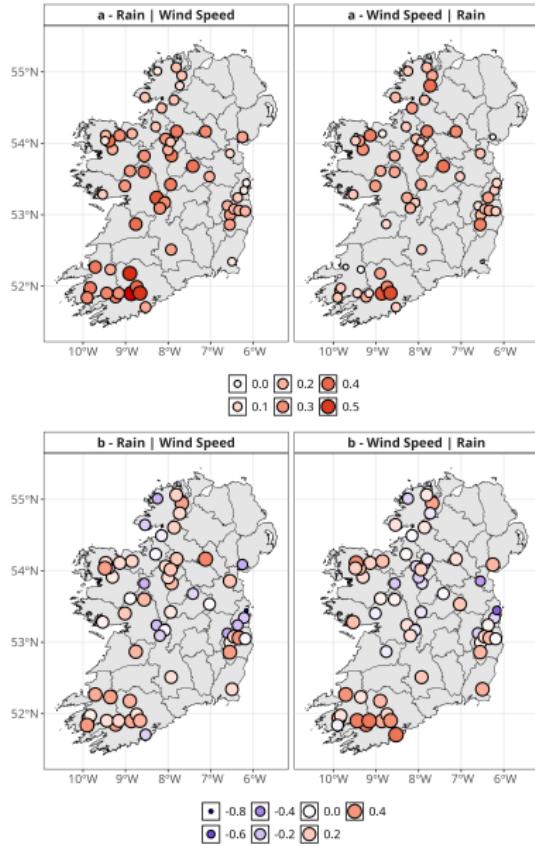
More realistic example



Uncertainty Reduction



α, β pre-clustering



α, β post-clustering

