

Discrete Multivariate Generalized Pareto Distribution with application to dry spells

Samira Aka^{1,2}, Marie Kratz², Philippe Naveau¹

¹Laboratoire des Sciences du Climat et de l'Environnement (LSCE)
CEA-CNRS-UVSQ

²Center for Research in Econo-Finance and Actuarial Sciences on Risk (CREAR)
ESSEC Business School

Data science pour les Risques Hydro-Climatiques et Côtiers 2025,
Roscoff, March 31, 2025

Defining a dry spell

Dry spell=number of consecutive dry days (below a precipitation amount threshold) ([Raymond et al., 2016, Lana et al., 2006])



Example of dry spell

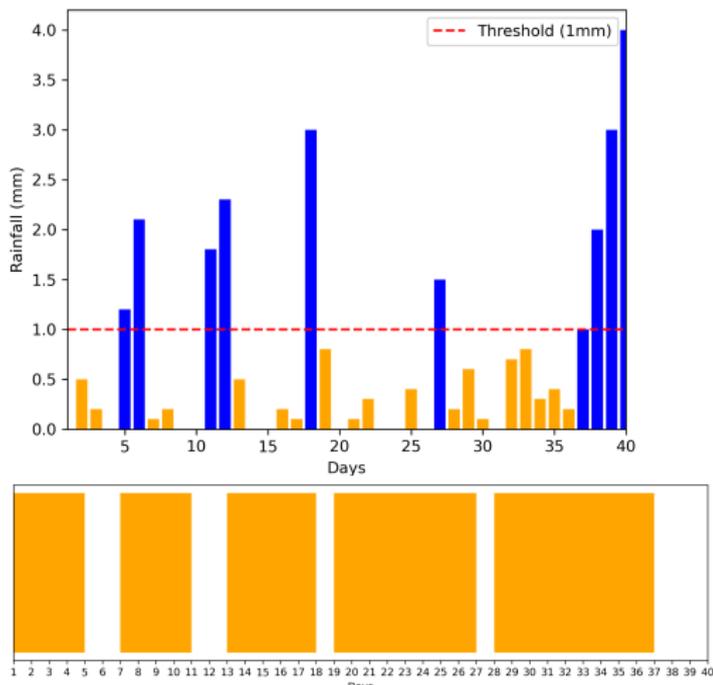
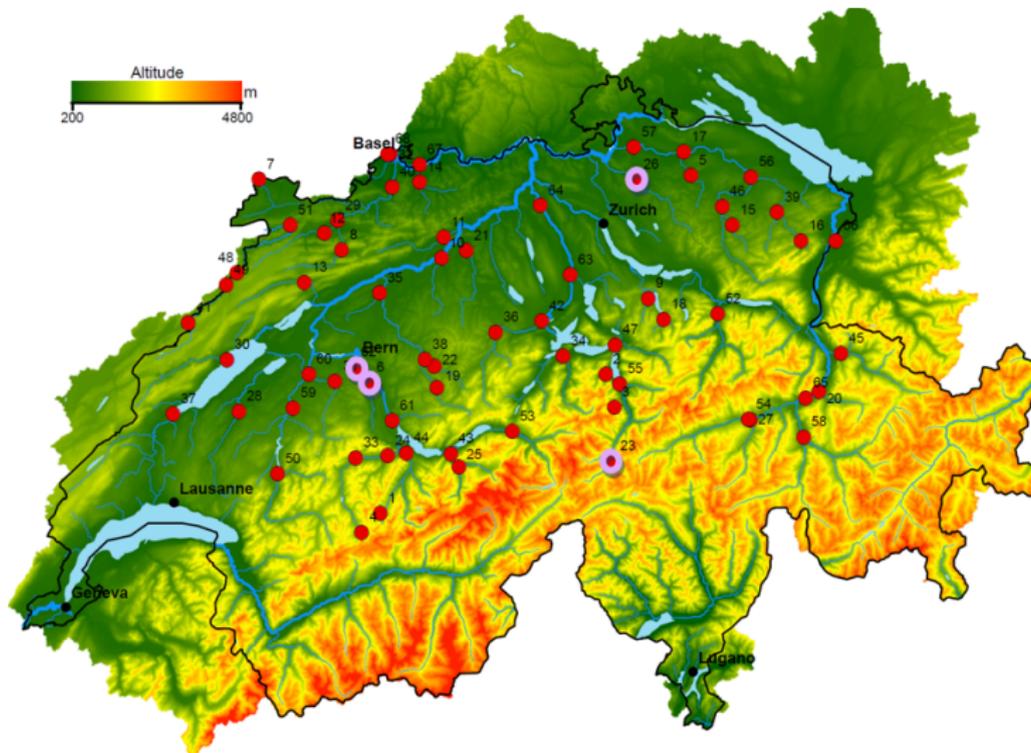
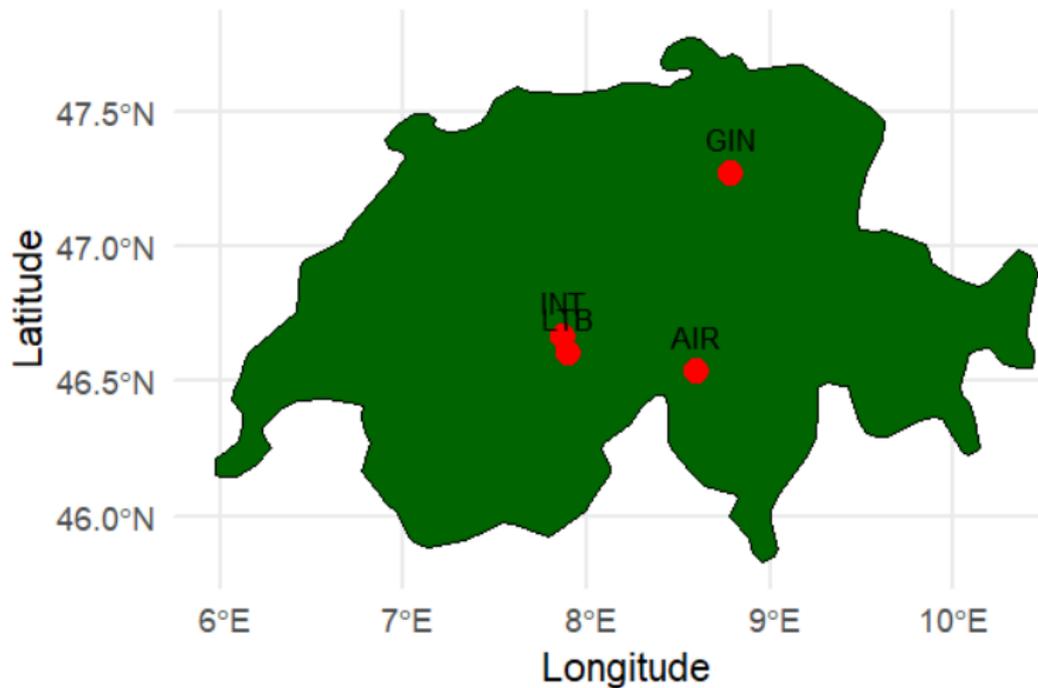


Figure: Precipitation amount with a threshold (blue for precipitation over the threshold and orange otherwise) and corresponding dry spells (orange)

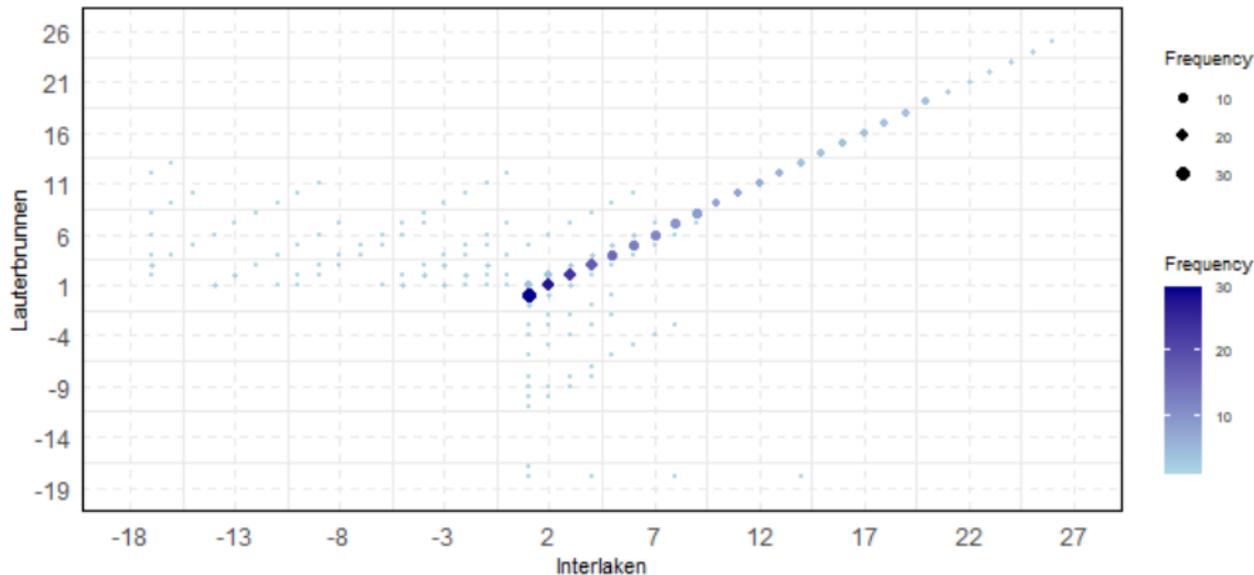
Precipitation stations in Switzerland



Selected Stations



Dry spells exceedances at two locations in Switzerland

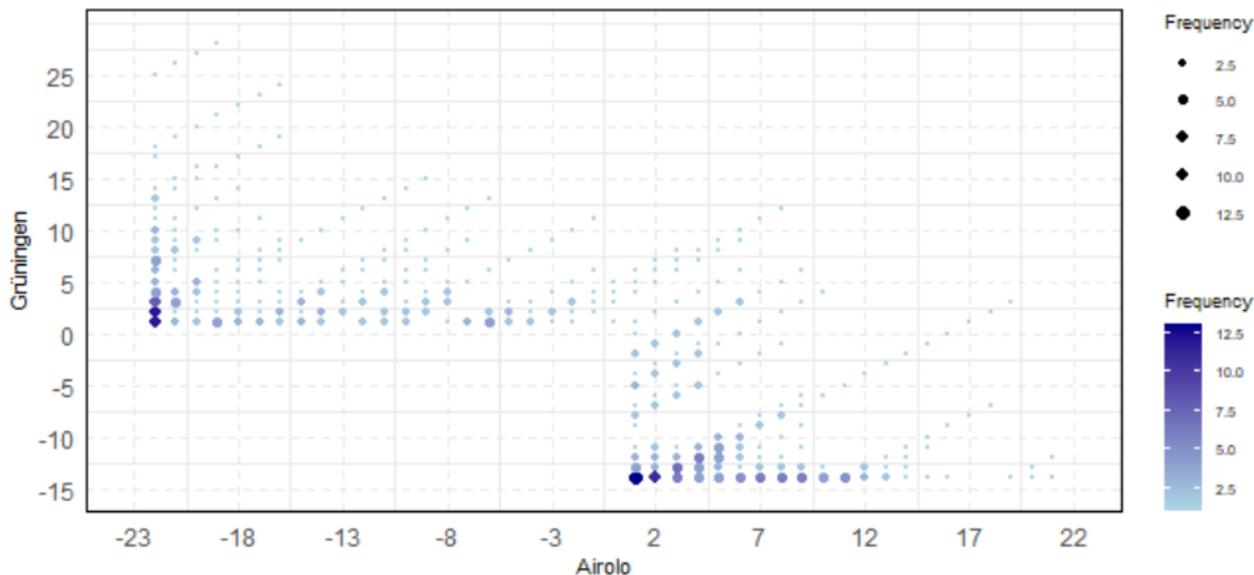


Dry spells exceeding quantile 99% : 17 days for Interlaken and 18 days for Lauterbrunnen (7 km between the stations)

What is the joint distribution ?

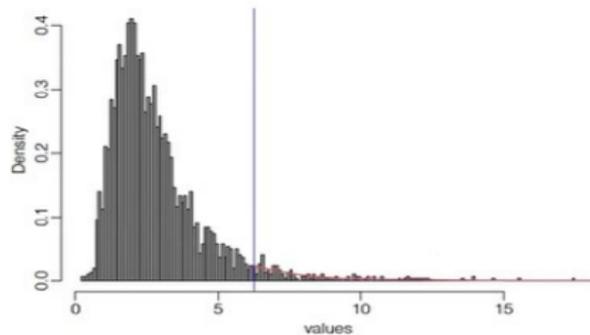
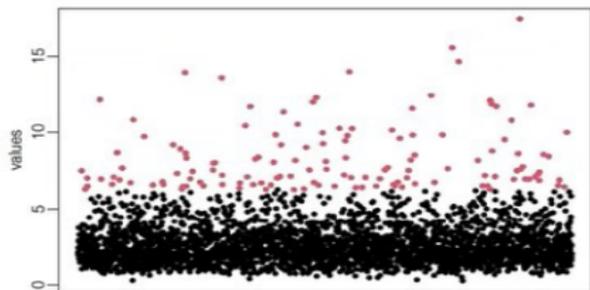
- **Main interest:** Modeling dry spells—successions of days with little or no precipitation—across various stations.
- **Broader question:** How can we model extremal dependence in multivariate discrete vectors?
- **Real-life applications:** Insurance claims, wildfire occurrences, and more.

Dry spells exceedances at two locations in Switzerland



Dry spells exceeding quantile 99% : 22 days for Airolo and 14 days for Grünigen (80 km between the stations)

Exceedances in the continuous case



Generalized Pareto Distribution (GPD)

The distribution of exceedances above a large threshold u can be approximated as (Pickands-Balkema-de Haan) :

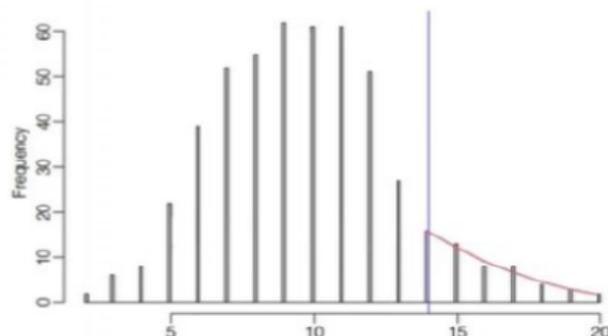
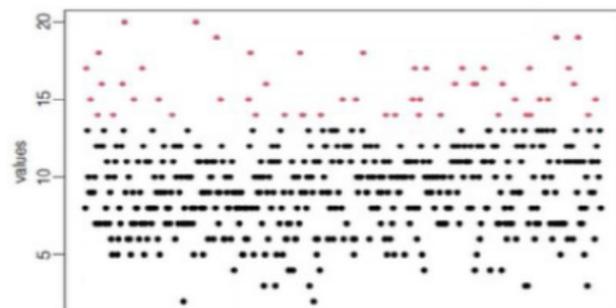
$$\mathbb{P}(Y - u > y \mid Y \geq u) \approx \overline{\text{GPD}}(y; \sigma_u, \xi)$$

where σ_u depends on u , and

$$\overline{\text{GPD}}(y; \sigma_u, \xi) = \left(1 + \xi \frac{y}{\sigma_u}\right)_+^{-\frac{1}{\xi}}, \text{ with } \sigma_u > 0,$$

is the survival function of the GPD.

Discrete GPD



Discrete generalized Pareto distribution (D-GPD) (see [Hitz et al., 2024], [Ahmad et al., 2022], [Daouia et al., 2023]).

Discrete GPD definition

From [Hitz et al., 2024], the probability mass function p_{DGPD} of the discrete GPD is defined as, for $k \in \mathbb{N}$,

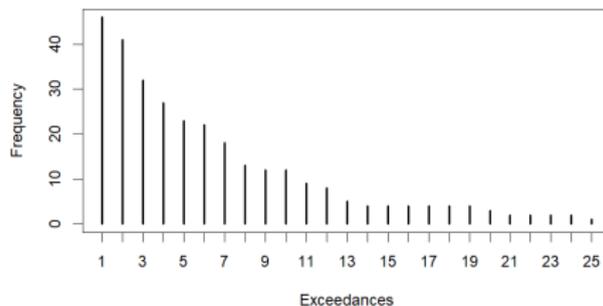
$$p_{\text{DGPD}}(k; \sigma, \xi) = \overline{\text{GPD}}(k; \sigma, \xi) - \overline{\text{GPD}}(k + 1; \sigma, \xi).$$

The fit of the DGPD can be done using the code attached to [Hitz et al., 2024].

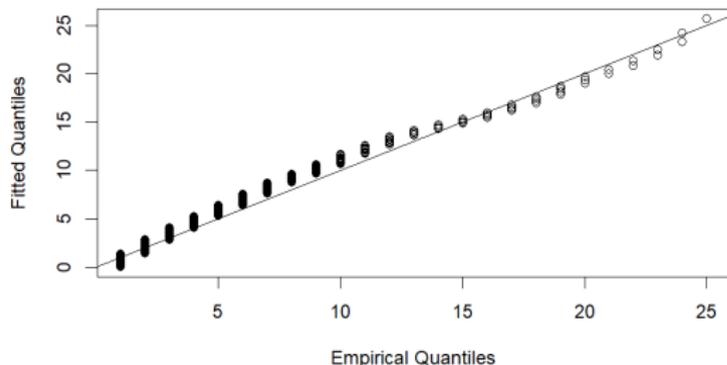
Fitting the marginal exceedances on a DGPD

$$\xi = -0.27 \text{ and } \sigma = 8.78.$$

Empirical Distribution for Station LTB



QQ Plot for Station LTB



Definition ([Rootzén et al., 2018], Theorem 7)

$\mathbf{Z} \in \mathbb{R}^d$ follows a $MGPD(\mathbf{1}, \mathbf{0}, \mathbf{S})$ if:

- $\max(\mathbf{Z})$ unit exponential distribution,
- $\mathbf{S} = \mathbf{Z} - \max(\mathbf{Z})$ with \mathbf{S} independent of $\max(\mathbf{Z})$.

Definition ([Rootzén et al., 2018])

$\mathbf{Z} \in \mathbb{R}^d$ follows a $MGPD(\mathbf{1}, \mathbf{0}, \mathbf{S})$ then

$$\mathbf{X} = \sigma \frac{e^{\gamma \mathbf{Z}} - \mathbf{1}}{\gamma},$$

follows a $MGPD(\sigma, \gamma, \mathbf{S})$.

What is well-known about the GPD:

- Univariate GPD for exceedances in continuous univariate data,

What is well-known about the GPD:

- Univariate GPD for exceedances in continuous univariate data,
- Univariate discrete GPD for exceedances in discrete univariate data,

What is well-known about the GPD:

- Univariate GPD for exceedances in continuous univariate data,
- Univariate discrete GPD for exceedances in discrete univariate data,
- Multivariate GPD for exceedances in continuous multivariate data.

What is well-known about the GPD:

- Univariate GPD for exceedances in continuous univariate data,
- Univariate discrete GPD for exceedances in discrete univariate data,
- Multivariate GPD for exceedances in continuous multivariate data.

→ **Aim: Construction of a GPD distribution for discrete multivariate data.**

Definition ([Aka et al., 2024])

$\mathbf{N} \in \mathbb{Z}^d$ follows a multivariate discrete Generalized Pareto Distribution $MDGPD(\mathbf{1}, \mathbf{0}, \mathbf{S})$ if :

- $\max(\mathbf{N})$ geometric distribution with parameter $1 - e^{-1}$

$$\mathbb{P}(\max(\mathbf{N}) \leq k) = 1 - e^{-k}, k \in \mathbb{N}^*$$

- $\mathbf{S} = \mathbf{N} - \max(\mathbf{N})$ with \mathbf{S} independent of $\max(\mathbf{N})$.

Proposition ([Aka et al., 2024])

$\mathbf{N} \sim \text{MDGPD}(\mathbf{1}, \mathbf{0}, \mathbf{S})$ and $\mathbf{m} \in \mathbb{N}^n$:

$$\mathcal{L}(\mathbf{N} - \mathbf{m} | \mathbf{N} \not\leq \mathbf{m}) = \text{MDGPD}(\mathbf{1}, \mathbf{0}, \mathbf{S}_d)$$

Proposition ([Aka et al., 2024])

$\mathbf{N} \sim \text{MDGPD}(\mathbf{1}, \mathbf{0}, \mathbf{S})$, $\mathbf{A} = (a_{ij})$ a matrix $\in \mathbb{N}^{n \times d}$ such that $\mathbb{P}(\sum_{j=1}^d a_{ij} N_j > 0) > 0, \forall i = 1, \dots, n$, and $\mathbf{m} \in \mathbb{N}^n$, then

$$\mathcal{L}(\mathbf{AN} - \mathbf{m} | \mathbf{AN} \not\leq \mathbf{m}) = \text{MDGPD}(\mathbf{A}\mathbf{1}, \mathbf{0}, \mathbf{S}_n).$$

Non-standard *MDGPD*

$\mathbf{N} \in \mathbb{Z}^d$ follows a *MDGPD*($\mathbf{1}, \mathbf{0}, \mathbf{S}$) then

$$\sigma \frac{e^{\gamma \mathbf{N}} - 1}{\gamma},$$

follows a *MDGPD*($\sigma, \gamma, \mathbf{S}$).

Non-standard MDGPD

$\mathbf{N} \in \mathbb{Z}^d$ follows a $MDGPD(\mathbf{1}, \mathbf{0}, \mathbf{S})$ then

$$\sigma \frac{e^{\gamma \mathbf{N}} - 1}{\gamma},$$

follows a $MDGPD(\sigma, \gamma, \mathbf{S})$.

NO!

Non-standard MDGPD

$\mathbf{N} \in \mathbb{Z}^d$ follows a $MDGPD(\mathbf{1}, \mathbf{0}, \mathbf{S})$ then

$$\sigma \frac{e^{\gamma \mathbf{N}} - 1}{\gamma},$$

follows a $MDGPD(\sigma, \gamma, \mathbf{S})$.

NO!

Definition ([Aka et al., 2024])

$\mathbf{K} \in \mathbb{Z}^d$ follows a non-standard multivariate discrete Generalized Pareto Distribution MDGPD $(\boldsymbol{\sigma} = \frac{\boldsymbol{\beta}}{\alpha}, \gamma = \frac{1}{\alpha}, \mathbf{S})$ if

$$\mathbb{P}(\mathbf{K} \leq \mathbf{k}) = 1 - \mathbb{E} \left(1 \wedge e^{\max(\mathbf{S} - \alpha \log(\frac{\mathbf{k}+1}{\boldsymbol{\beta}} + 1))} \right).$$

Generalization of the *MDGPD* by ratio on a Gamma

Proposition ([Aka et al., 2024])

$\mathbf{Z} \text{ MGPD}(\mathbf{1}, \mathbf{0}, \mathbf{S})$ and Λ a $\text{Gamma}(\alpha, \beta)$ random variable independent from \mathbf{Z} . Then,

$$\left[\frac{\mathbf{Z}}{\Lambda} \right] \text{ follows a MDGPD} \left(\frac{\beta}{\alpha}, \frac{\mathbf{1}}{\alpha}, \mathbf{S} \right).$$

Proposition ([Aka et al., 2024])

Let $\mathbf{M} = \frac{\mathbf{Z}}{\Lambda}$, where \mathbf{Z} and Λ are defined as in the previous proposition, $\mathbf{A} = (a_{ij})$ be a matrix $\in \mathbb{N}^{n \times d}$ such that $\mathbb{P}(\sum_{j=1}^d a_{ij} N_j > 0) > 0, \forall i = 1, \dots, n$, and $\mathbf{m} \in \mathbb{N}^n$, then $\mathcal{L}(\lfloor \mathbf{AM} \rfloor - \mathbf{m} \mid \lfloor \mathbf{AM} \rfloor \not\leq \mathbf{m}) = \text{MDGPD}(\frac{\beta}{\alpha} \mathbf{A} \mathbf{1}, \frac{1}{\alpha}, \mathbf{S}_n)$.

Building \mathbf{N} from the increments of a distribution

If $\mathbf{N} \sim \text{MDGPD}(\mathbf{1}, \mathbf{0}, \mathbf{S})$, then:

$$\mathbf{N} = \underbrace{\mathbf{T} - \max(\mathbf{T})}_{\mathbf{S} \perp G} + \underbrace{G}_{\max(\mathbf{N})},$$

with \mathbf{T} any discrete random vector called the generator.

Other definition of N in the bivariate case

In the bivariate case,

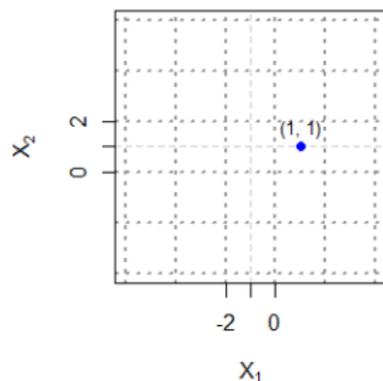
$$N_1 = G + (T_1 - T_2)\mathbb{1}_{((T_1 - T_2) < 0)},$$

$$N_2 = G - (T_1 - T_2)\mathbb{1}_{((T_1 - T_2) \geq 0)}.$$

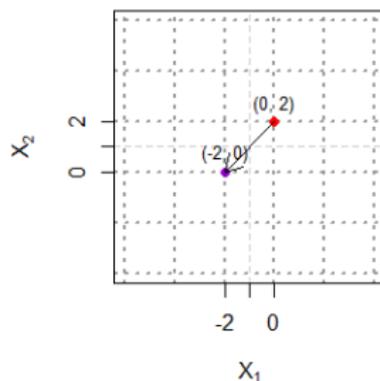
Note that $N_1 - N_2 = T_1 - T_2$.

Creating a Bivariate DGPD (1)

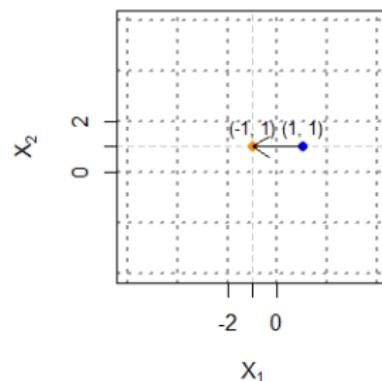
G Point (Blue)



(T_1, T_2) Point (Red)

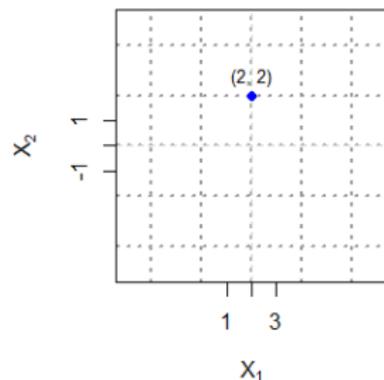


Bivariate DGPD Point (Yellow)

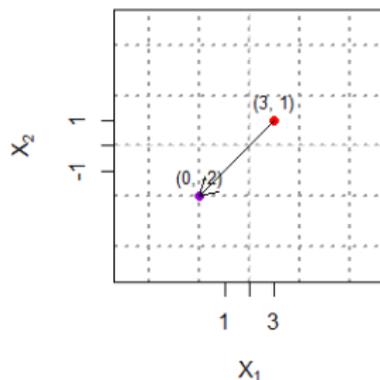


Creating a Bivariate DGPD (2)

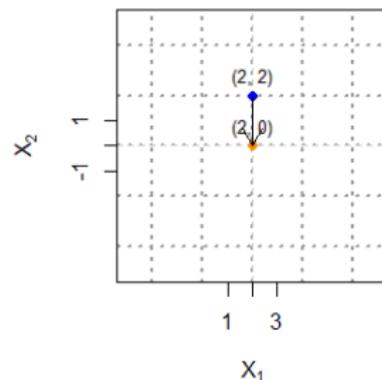
G Point (Blue)



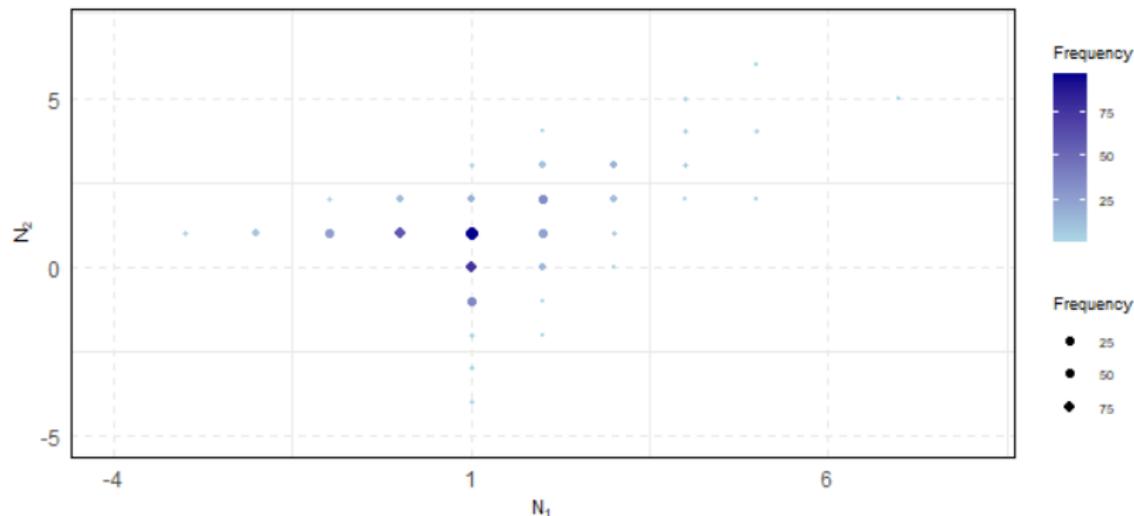
(T_1, T_2) Point (Red)



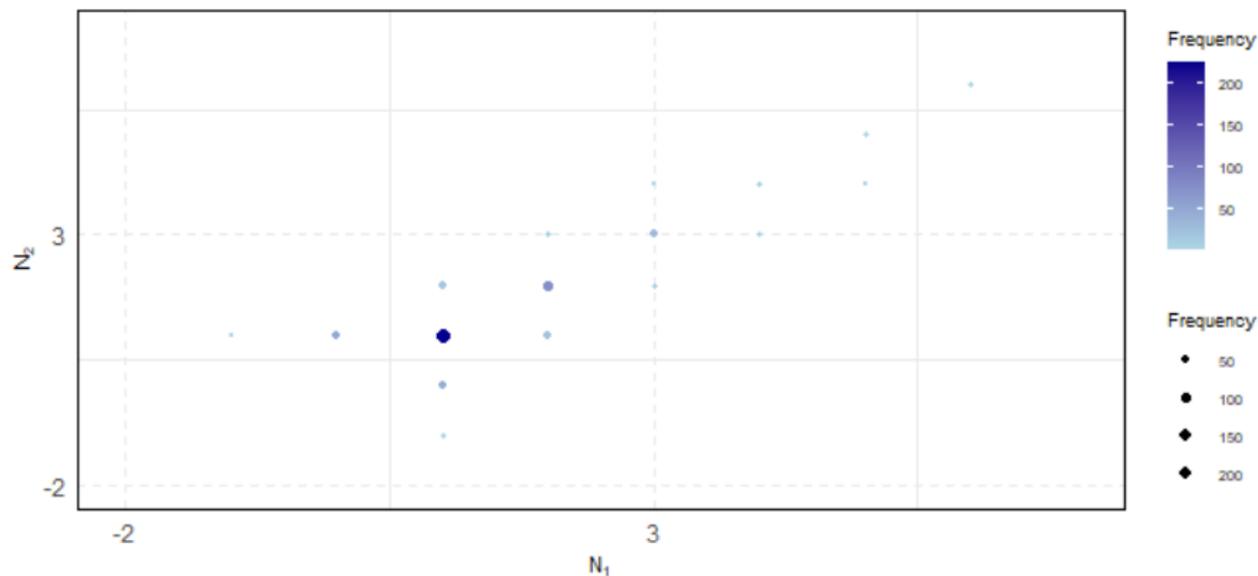
Bivariate DGPD Point (Yellow)



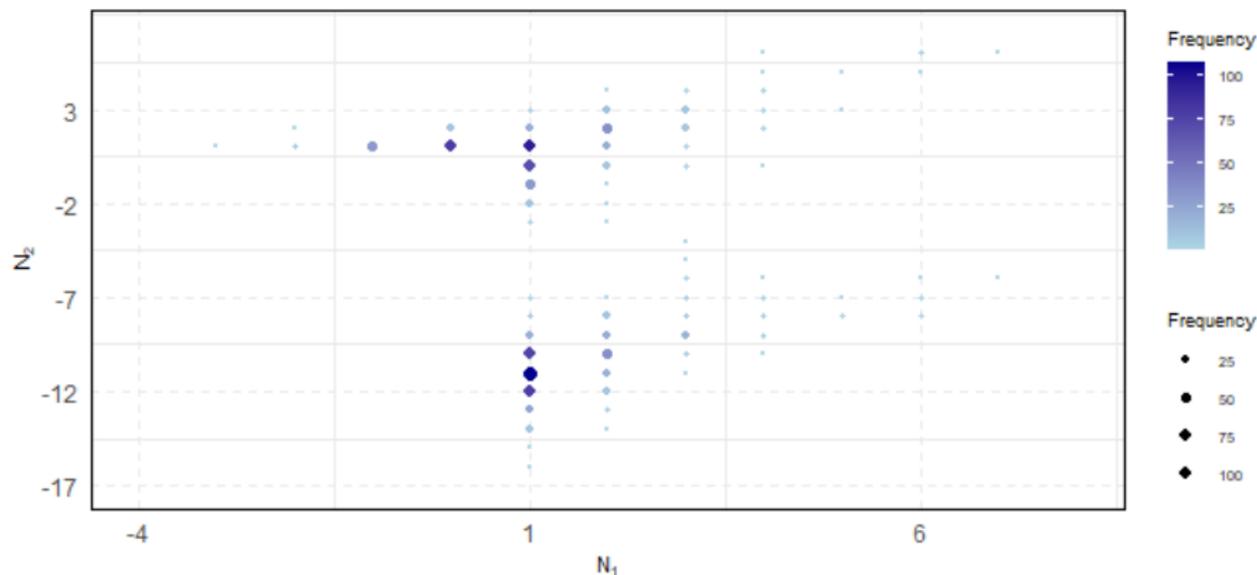
Bivariate DGPD: T_1 and T_2 Independent Poisson Variables



Bivariate DGPD: T_1 and T_2 Dependent Poisson Variables



Bivariate DGPD: T_1 and T_2 Bimodal Poisson Variables



Bootstrap algorithm to simulate a bivariate $DGPD$ from unknown (T_1, T_2)

From Bootstrap simulations algorithm in [Legrand et al., 2021].

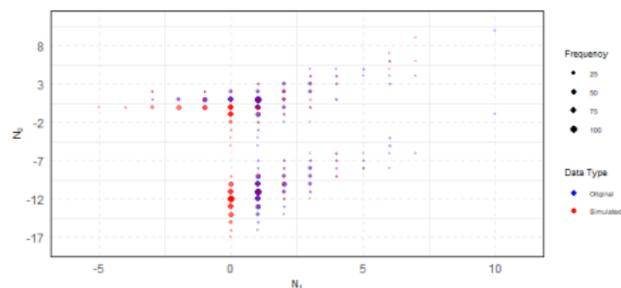
Algorithm Bootstrap $MDGPD$ simulation

Require: A sample of $(N_{1,i}, N_{2,i})_{1 \leq i \leq n} \sim MDGPD$

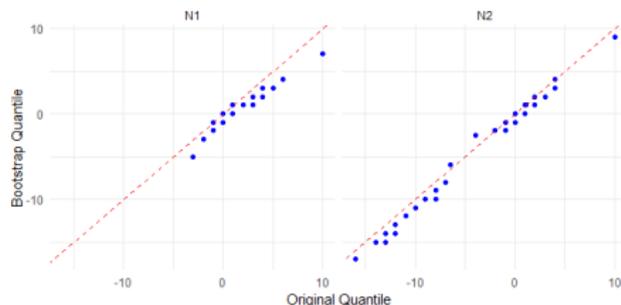
Ensure: A discrete simulated sample $(N_{1,k}^*, N_{2,k}^*)_{1 \leq k \leq m}$ to choose

- 1: Define $\Delta_i = N_{1,i} - N_{2,i}$, $1 \leq i \leq n$,
 - 2: Generate m generalizations $G_k \sim \text{Geom}(1 - e^{-1})$, $1 \leq k \leq m$, independently from Δ_i ,
 - 3: Bootstrap m realization Δ_k^* from $(\Delta_1, \dots, \Delta_n)$
 - 4: **Return** $N_{1,k} := G_k + \Delta_k^* \mathbb{1}_{(\Delta_k^* < 0)}$ and $N_{2,k} := G_k - \Delta_k^* \mathbb{1}_{(\Delta_k^* \geq 0)}$, for $1 \leq k \leq m$.
-

Bivariate DGPD $N = (N_1, N_2)$ from unknown (T_1, T_2)



(a) Scatter plot of simulated data with the parametric model of sample size $n = 500$ (blue dots) and sampled data from one simulation with sample size $m = 500$ (red dots).



(b) Q-Q plot for marginals N_1 and N_2 .
x-axis : Original sample
y-axis : Bootstrap sample

Figure: Bootstrap simulations of (N_1, N_2) using bimodal Poisson (T_1, T_2) and Geometric variable G .

How to Perform Inference on the *MDGPD*?

Discrete Data

↓
Challenges

Marginals \neq *Geometric*($1 - e^{-1}$).

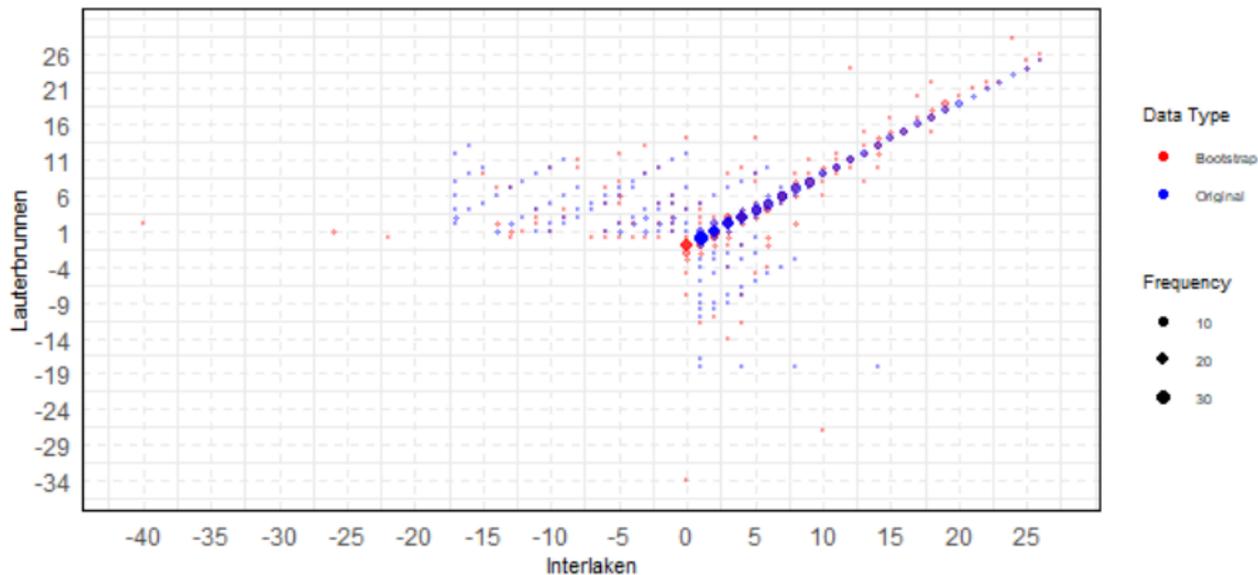
How to transform discrete distributions into another?

- **If continuous:** 😊 Easy
- **If discrete:** 😞 Not so easy

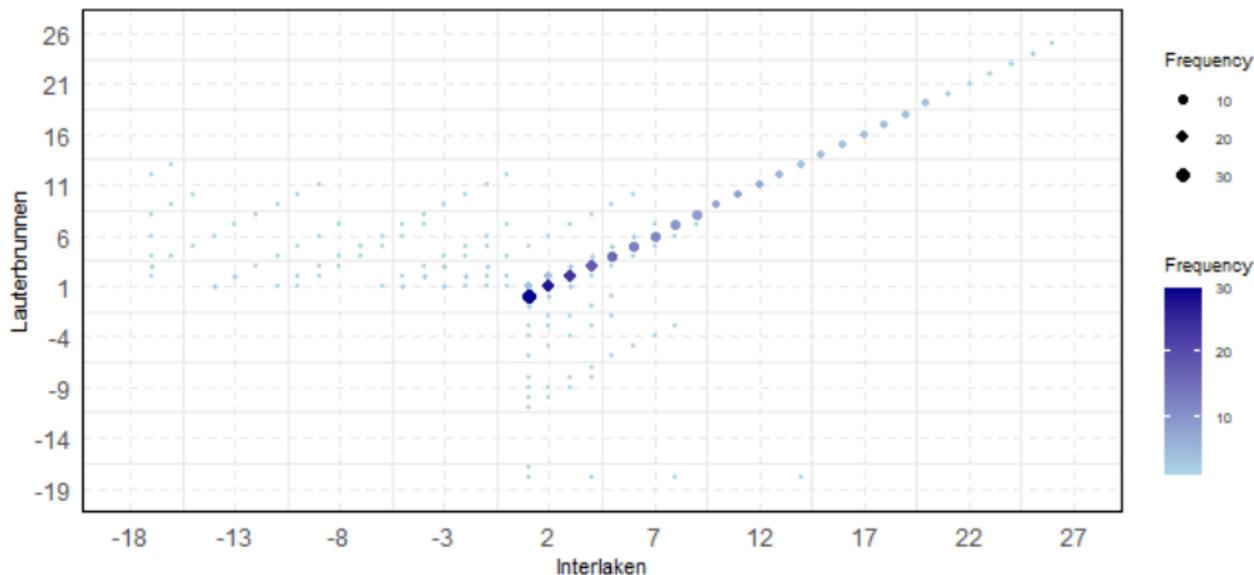
**A solution: Simulation-based
inference with Neural Networks**

- Neural networks (NN) are used for inference in intractable models, such as those in [Pacchiardi et al., 2021] and [Lenzi et al., 2023], which focus on max-stable models.
- In this work, we apply **Neural Bayes Estimators** as introduced by [Sainsbury-Dale et al., 2024].

Bootstrap simulations at Two Close Stations in Switzerland

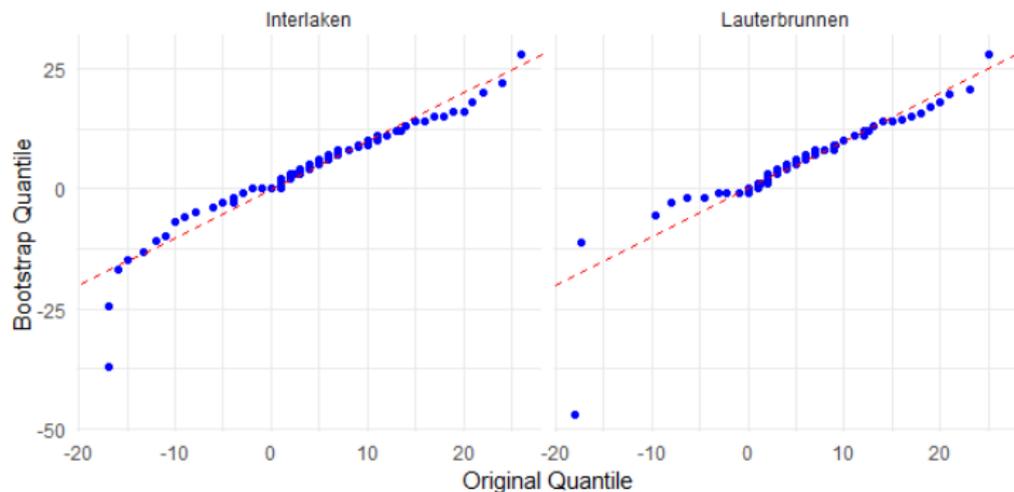


Dry spells exceedances at two locations in Switzerland



Dry spells exceeding quantile 99% : 17 days for Interlaken and 18 days for Lauterbrunnen (7 km between the stations)

Quantile-Quantile plot of marginals (close pair)



- Develop a dedicated package for the MDGPD, facilitating easier application and wider use.
- Investigate extreme value regression using covariates such as geographical features, temperature, and soil moisture, within the MDGPD framework.

 Ahmad, T., Gaetan, C., and Naveau, P. (2022).

Modelling of discrete extremes through extended versions of discrete generalized pareto distribution.

 Aka, S., Kratz, M., and Naveau, P. (2024).

Discrete multivariate generalized Pareto distribution with application to dry spells.

 Balkema, A. A. and de Haan, L. (1974).

Residual life time at great age.

The Annals of Probability, 2.



Daouia, A., Stupfler, G., and Usseglio-Carleve, A. (2023).

Extreme value modelling of sars-cov-2 community transmission using discrete generalized pareto distributions.

Royal Society Open Science, 10(3):220977.



Hitz, A. S., Davis, R. A., and Samorodnitsky, G. (2024).

Discrete extremes.

Journal of Data Science, pages 1–13.

References III



Lana, X., Martínez, M. D., Burgueño, A., Serra, C., Martín-Vide, J., and Gómez, L. (2006).

Distributions of long dry spells in the iberian peninsula, years 1951–1990.

International Journal of Climatology, 26:1999–2021.



Legrand, J., Naveau, P., and Oesting, M. (2021).

Evaluation of binary classifiers for asymptotically dependent and independent extremes.



Lenzi, A., Bessac, J., Rudi, J., and Stein, M. L. (2023).

Neural networks for parameter estimation in intractable models.

Computational Statistics Data Analysis, 185:107762.

References IV



Pacchiardi, L., Adewoyin, R., Dueben, P., and Dutta, R. (2021).

Probabilistic forecasting with generative networks via scoring rule minimization.

Journal of Machine Learning Research, 25:45:1–45:64.



Pickands, J. (1975).

Statistical inference using extreme order statistics.

The Annals of Statistics, 3.



Sainsbury-Dale, M., Zammit-Mangion, A., and Huser, R. (2024).
Likelihood-free parameter estimation with neural bayes estimators.
The American Statistician, 78:1–14.