# Identifying drivers for heat waves using ~~wavelets and~~ machine learning approaches
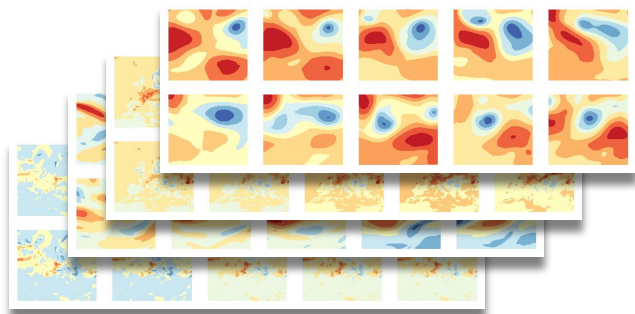
Sebastian Buschow[1], Jan Keller[2,3] and Sabrina Wahl[1,2]

[1] University of Bonn, Institute of Geosciences, Meteorology department
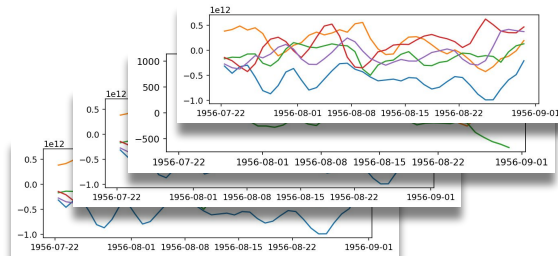[2] Hans-Ertel-Zentrum für Wetterforschung
[3] Deutscher Wetterdienst

Bundesministerium
für Bildung
und Forschung

FONA
Research for sustainability

CLIM XTREME
Climate Change and Extreme Events

Fields of geopotential, moisture, …

Principal components

PCA

predictors

Binary heatwaves: **3 or more *hot days*\* in a row**

2010-08-01

1972-07-01

2003-08-10

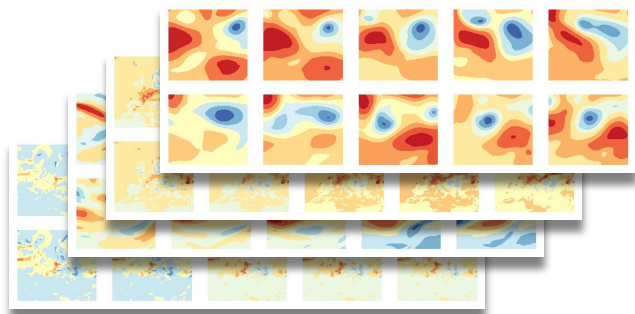logistic
PCA

Principal HW components

predictand

Statistical learner

- Simple neural net
- Dropout layer
- 7 fold CV in blocks of years
- lm for comparison

---

***hot day**: daily T$_{max}$ > 90% quantile for the current calendar day

**Data**: ERA5 JJA, 1950-2020, EURO-CORDEX domain

Fields of geopotential, moisture, …



Principal components



PCA

predictors

**Binary heatwaves: 3 or more *hot days\** in a row**



2010-08-01

1972-07-01

2003-08-10

logistic
PCA

*How does
this work?*

Principal HW components



predictand

*What does the model do?*

**Statistical learner**

- Simple neural net
- Dropout layer
- 7 fold CV in blocks of years
- lm for comparison

---

***hot day**: daily $T_{max}$ > 90% quantile for the current calendar day          **Data**: ERA5 JJA, 1950-2020, EURO-CORDEX domain
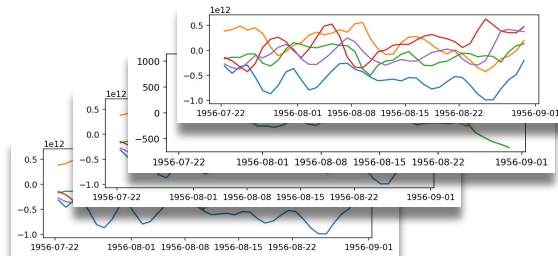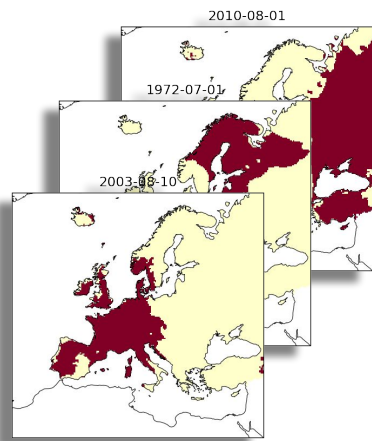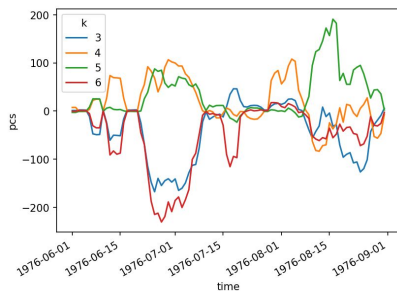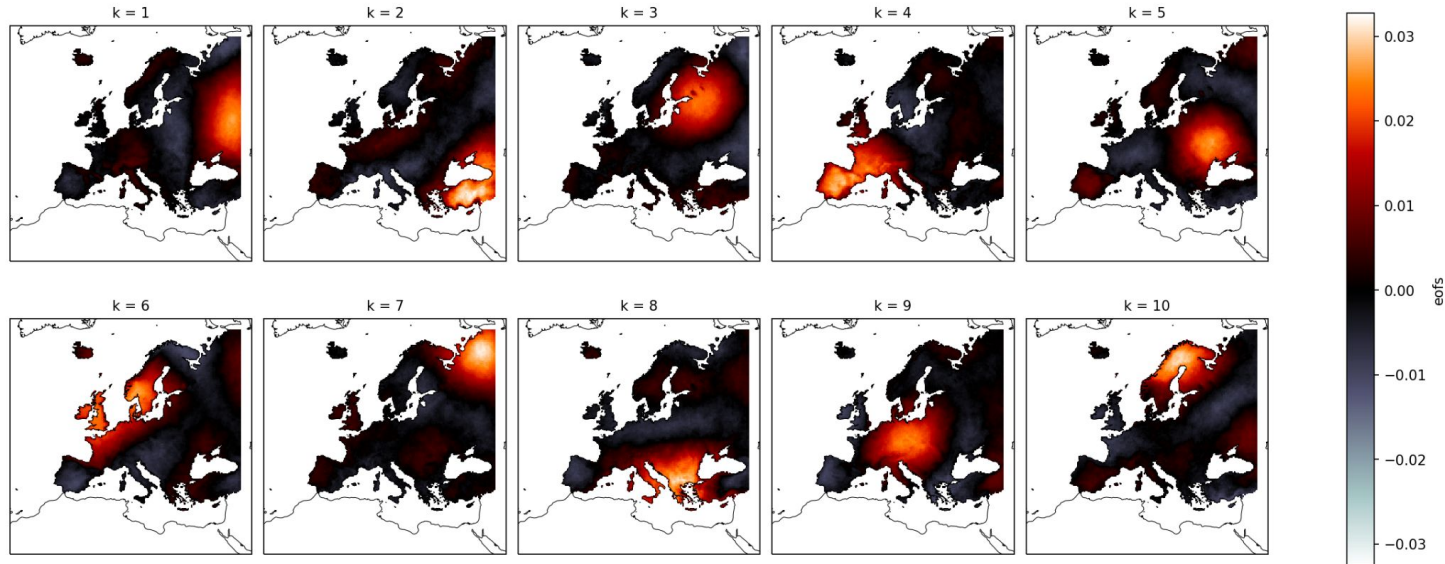
# PCA for binary fields (!)

Regular PCA ("EOFs"):   given data $\mathbf{X}$, find vectors $\mathbf{U}$ such that $|\mathbf{X} - \mathbf{U}\mathbf{U}^T\mathbf{X}|^2$ is minimal
→ minimize *Gauss deviance*, solution: eigenvectors of $\mathbf{X}^T\mathbf{X}$

Landgraf and Lee (2020): assume exponential family, compute natural parameters $\boldsymbol{\theta}$
→ **iteratively** search a projection $\boldsymbol{\theta}\mathbf{U}\mathbf{U}^T$ that minimizes the relevant deviance

Binary data:  Bernoulli distribution, $\boldsymbol{\theta} = \log(\ \mathbf{p}\ /\ (\ 1 - \mathbf{p}\ )\ )$



*10 rotated "logistic EOFs" for European heatwaves*

4

# Modelling heatwaves in the reduced space



Predictand: 10 rotated logistic PCs of heatwaves
Predictors: 20 PCs of soil moisture and geopotential at 1000, 800, 500, 300 and 100 hPa
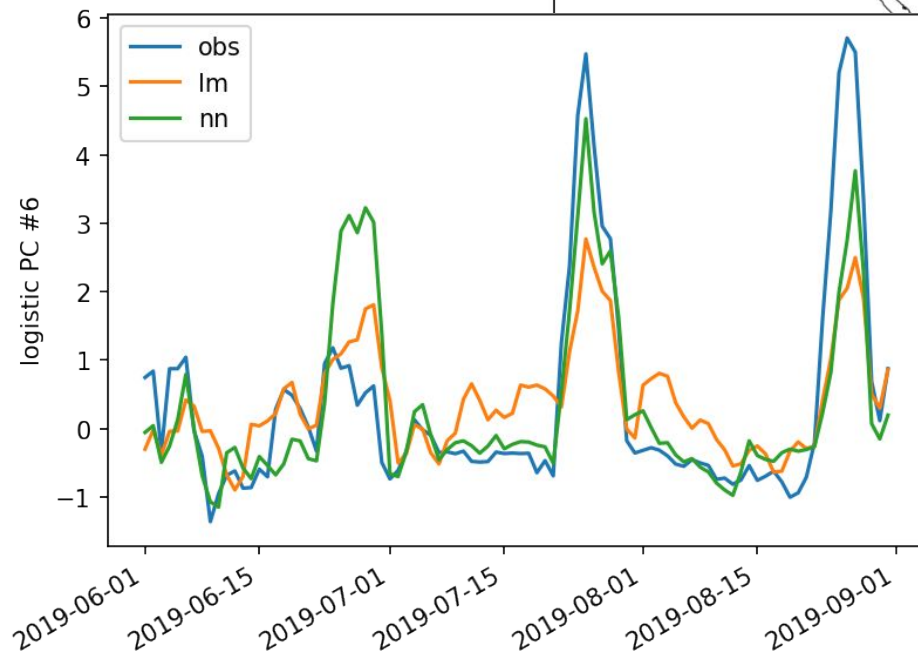
**Multivariate linear regression**

- $10 + 20 \times 10 = 210$ parameters
- Least squares fit
- **$R^2$=0.46**

**vs.**

**Simple feed forward neural net**

- One hidden layer with 40 nodes
  $\rightarrow (20+1) \times 40 + (40+1) \times 10$
  $= 1250$ parameters
- ReLu activation, 20% dropout
- Optimized with *Adam*
- **$R^2$=0.75**



*Observed and modelled heatwave PC # 6 in summer 2019*

# Modelling heatwaves in the reduced space

Predictand: 10 rotated logistic PCs of heatwaves
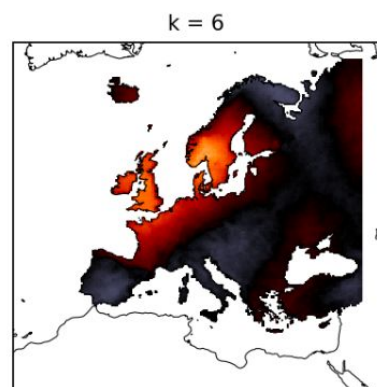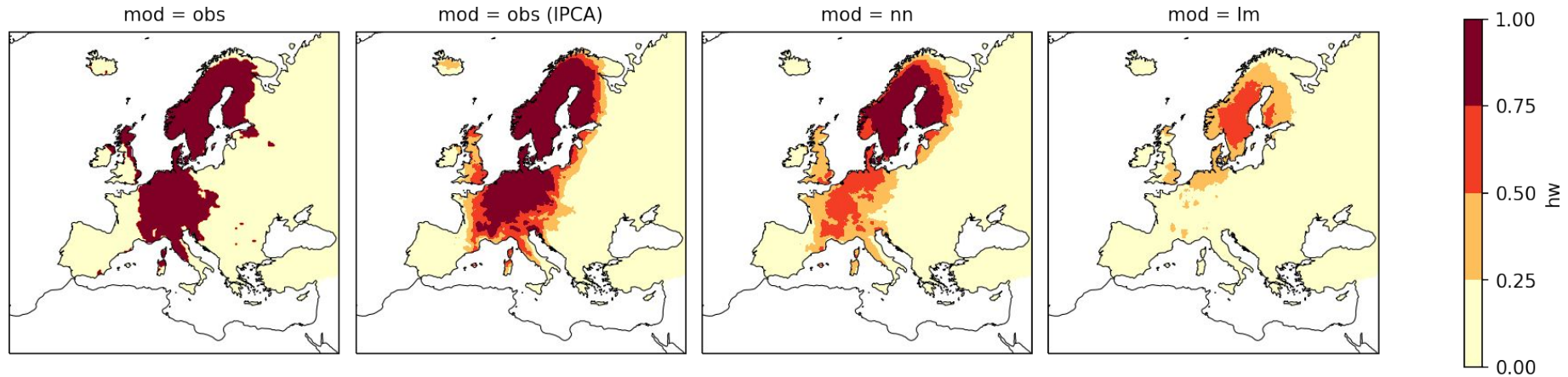Predictors: 20 PCs of soil moisture and geopotential at 1000, 800, 500, 300 and 100 hPa
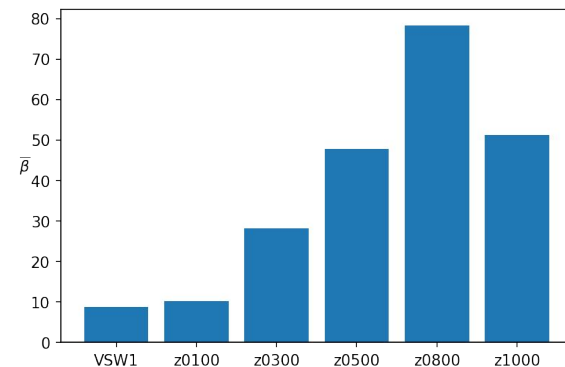


*Observed, PCA reduced, and simulated heatwaves on 2019-07-26*

# Variable importance



*Mean absolute regression coefficients*
$\rightarrow$ *is* $\Phi_{800}$ *the most important predictor?*

Models are not bad, but how do they identify heatwaves?
Linear model: just look at coefficients (?)

What to do for the neural net? The coefficients tell us
nothing !

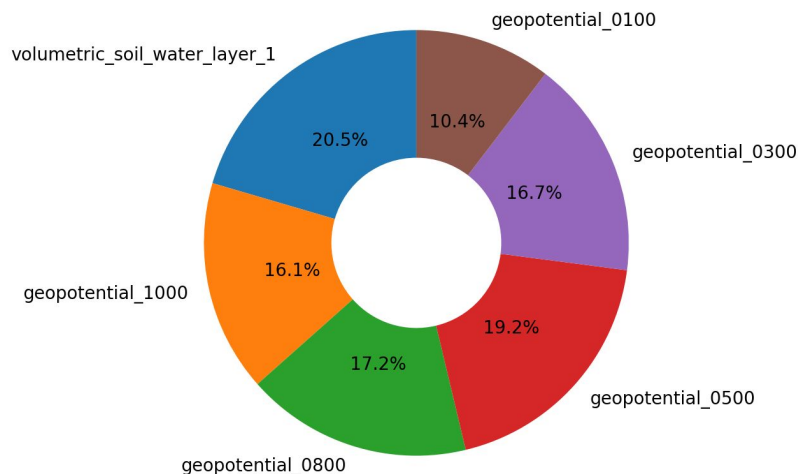Idea (Shapley 1952, Lipovetsky and Conklin 2001): Split up the overall $R^2$ as follows:

$$R^2 = f(X_1, ..., X_n) = \varphi_1 + \varphi_2 + ... + \varphi_n$$

$$\varphi_i = \underbrace{n^{-1} \sum_{j=0}^{n-1}}_{\text{mean over set sizes}} \underbrace{\binom{n-1}{j}^{-1} \sum_{\text{all } S \text{ with } |S|=j, \, X_i \notin S}}_{\text{mean over sets of size } j \text{ missing } X_i} \underbrace{f(S \cup X_i) - f(S)}_{\text{change if } i \text{ were added}}$$

$\rightarrow$ train all possible $2^6$ models, compare their $R^2$ to get the Shapley values!

# Variable importance: Shapley values

$$\varphi_i = \underbrace{n^{-1}\sum_{j=0}^{n-1}}_{\text{mean over set sizes}} \underbrace{\binom{n-1}{j}^{-1}\sum_{\text{all } S \text{ with } |S|=j,\, X_i \notin S}}_{\text{mean over sets of size } j \text{ missing } X_i} \underbrace{f(S \cup X_i) - f(S)}_{\text{change if } i \text{ were added}}$$
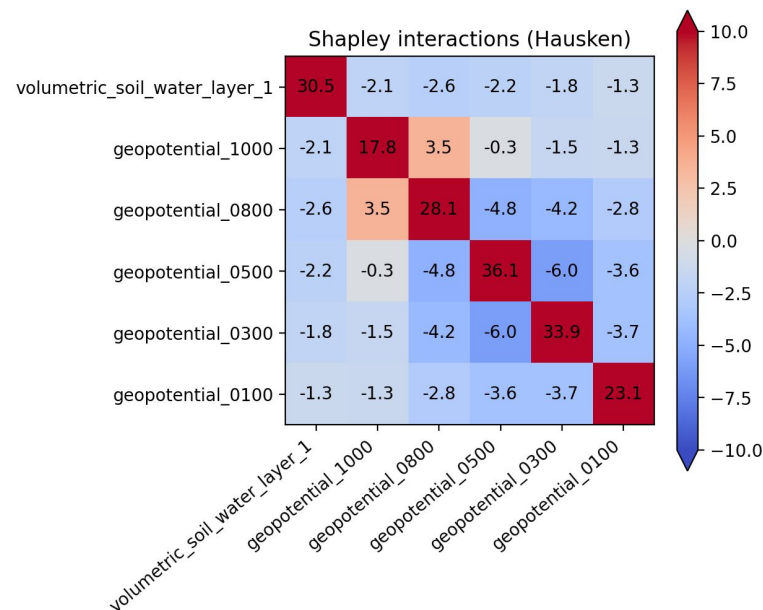


*Percentage contributions to the overall model performance $R^2$ for the neural net*



*Recursive Shapley "interactions" (Hausken 2001) for the neural net*

# Variable importance: Shapley values

$$\varphi_i = \underbrace{n^{-1}\sum_{j=0}^{n-1}}_{\text{mean over set sizes}} \underbrace{\binom{n-1}{j}^{-1} \sum_{\text{all } S \text{ with } |S|=j,\ X_i \notin S}}_{\text{mean over sets of size } j \text{ missing } X_i} \underbrace{f(S \cup X_i) - f(S)}_{\text{change if } i \text{ were added}}$$
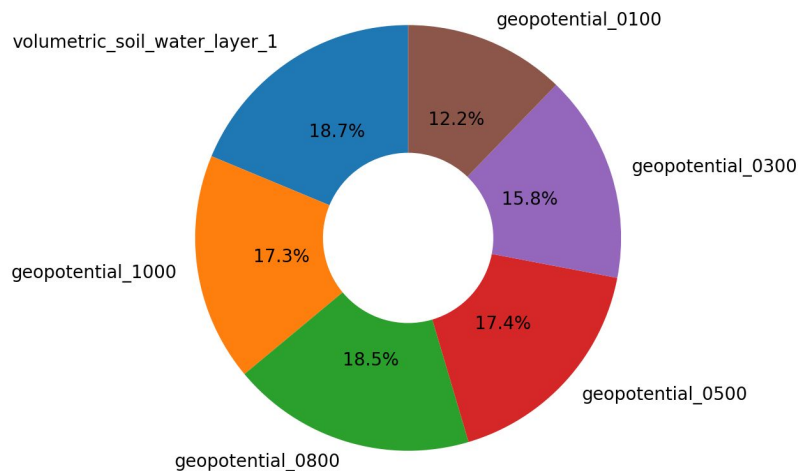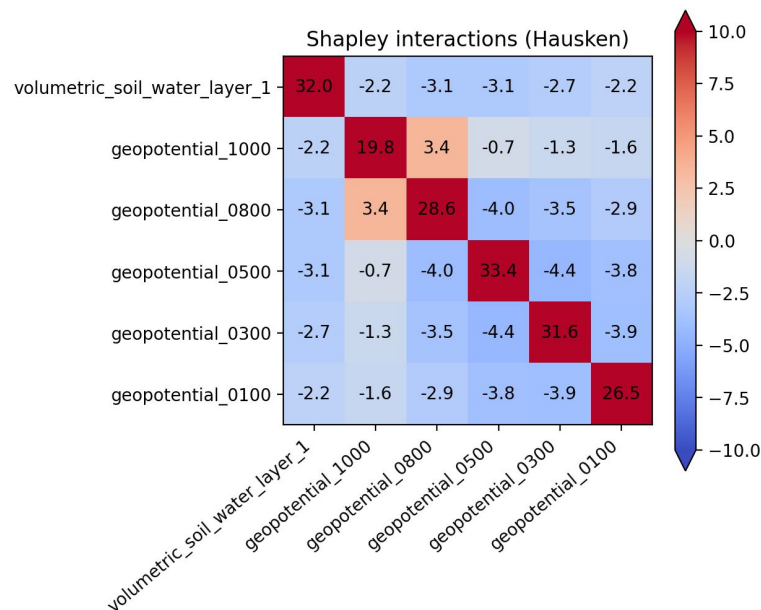


*Percentage contributions to the overall model performance R² for the **linear model***



*Recursive Shapley "interactions" (Hausken 2001) for the **linear model***

# Summary

- There is a **PCA for binary variables** → reduced version of any binary event you want !

- A simple neural net can explain 75% of the reduced heatwave variability

- **Shapley values** and interactions reveal how much *can* be learned from each predictor, lm and neural net are not so different after all

- It seems that the model has learned $\Phi_{800}$ - $\Phi_{1000}$ ~ $T_{900}$ (hydrostatic relation)

---

# References

Hausken, K., & Mohr, M. (2001). The value of a player in n-person games. *Social Choice and Welfare*, *18*(3), 465-483
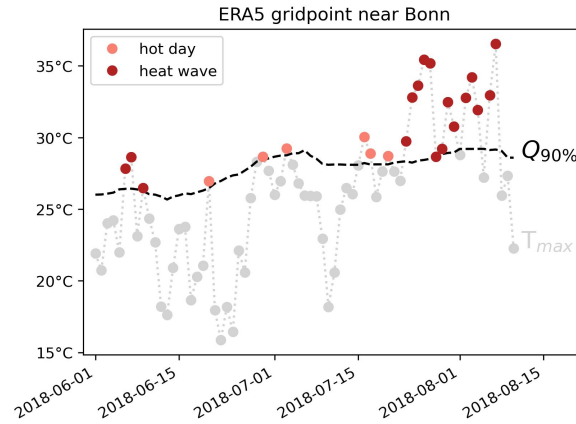
Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz‑Sabater, J., ... & Thépaut, J. N. (2020). The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, *146*(730), 1999-2049.

Landgraf, A. J., & Lee, Y. (2020). Dimensionality reduction for binary data through the projection of natural parameters. *Journal of Multivariate Analysis*, *180*, 104668.
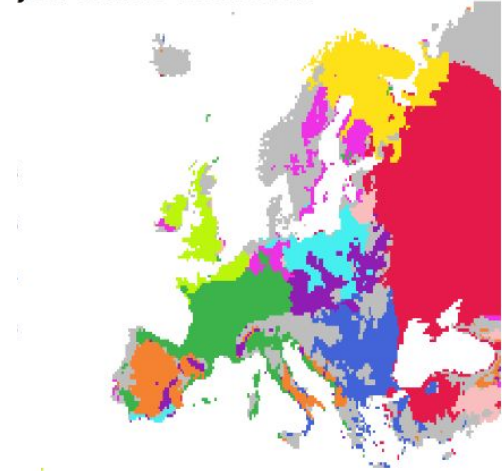
Lipovetsky, S., & Conklin, M. (2001). Analysis of regression in game theory approach. *Applied Stochastic Models in Business and Industry*, *17*(4), 319-330.

Shapley, L. S. (1997). A value for n-person games. *Classics in game theory*, *69*.
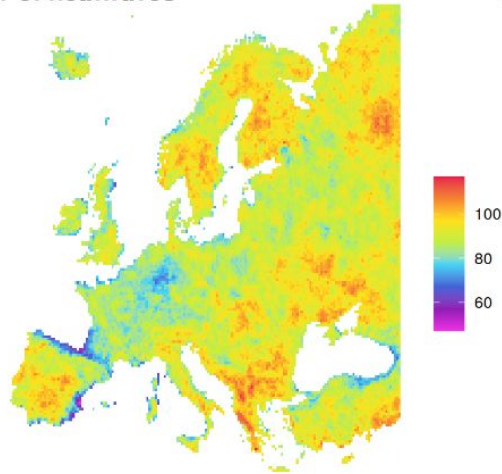
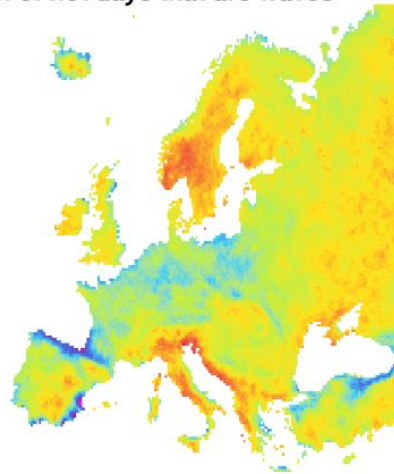# Heatwave definition and basic statistics



ERA5 gridpoint near Bonn

year of most intense HW

number of heatwaves

fraction of hot days that are waves

average duration