

# Latent variable models for multi-variable space-time data and applications in hydrology

Benjamin Renard<sup>1</sup>   Mark Thyer<sup>2</sup>   David McInerney<sup>2</sup>  
Dmitri Kavetski<sup>2</sup>   Michael Leonard<sup>2</sup>   Seth Westra<sup>2</sup>

<sup>1</sup>INRAE, RiverLy Research Unit, Lyon, France

<sup>2</sup>School of Civil, Environmental and Mining Engineering, University of Adelaide, Australia

AppliBUGS 10 December 2021



# Introduction

## Resources and risk management for environmental systems

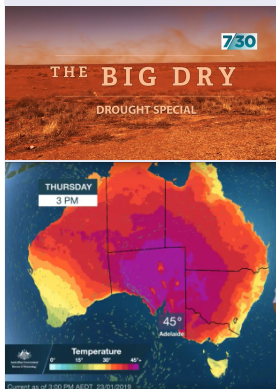
Often relies on the analysis of **several variables** measured at **many sites** and whose properties may **vary in time**.

# Introduction

## Resources and risk management for environmental systems

Often relies on the analysis of **several variables** measured at **many sites** and whose properties may **vary in time**.

## Example: Australian summers

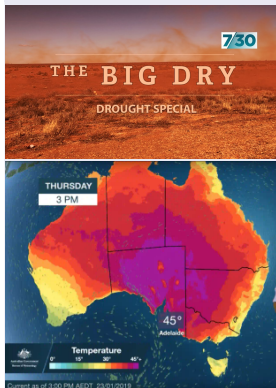


# Introduction

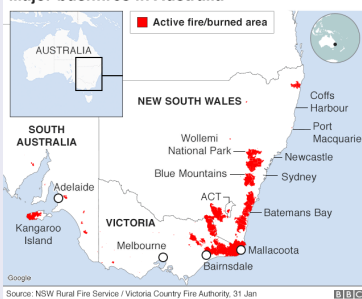
## Resources and risk management for environmental systems

Often relies on the analysis of **several variables** measured at **many sites** and whose properties may **vary in time**.

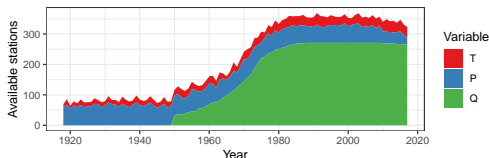
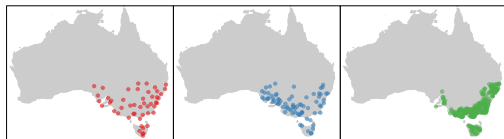
## Example: Australian summers



### Major bushfires in Australia



# Motivating dataset

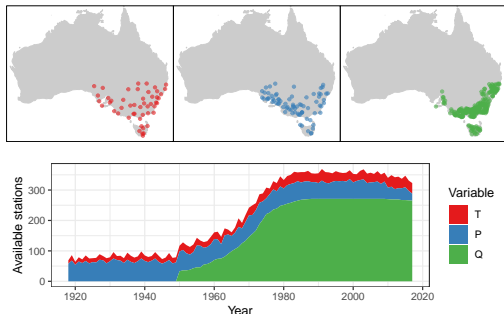


BoM's reference datasets for Temperature, Precipitation and Streamflow.

Variables of interest (DJF):

- 1 Number of heatwaves  $T_n$
- 2 Heatwave intensities  $T_x$
- 3 Dry-day duration  $P_d$
- 4 Drought duration  $Q_d$

# Motivating dataset



BoM's reference datasets for Temperature, Precipitation and Streamflow.

Variables of interest (DJF):

- ① Number of heatwaves  $T_n$
- ② Heatwave intensities  $T_x$
- ③ Dry-day duration  $P_d$
- ④ Drought duration  $Q_d$

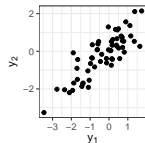
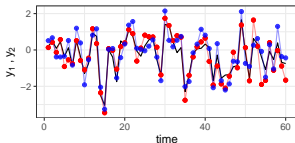
A model for  $T_n$ ,  $T_x$ ,  $P_d$  and  $Q_d$  should handle...

- Spatial dependence
- Time variability and/or trend
- Inter-variable dependence
- Data of different types, missing data

# Latent variables as Hidden Climate Indices (HCI)

Variables affected by THE SAME climate index are dependant

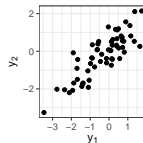
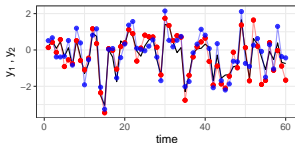
$$\begin{cases} Y_{1,t} = \lambda_1 \tau_t + \varepsilon_{1,t} \\ Y_{2,t} = \lambda_2 \tau_t + \varepsilon_{2,t} \\ \varepsilon_{1,t}, \varepsilon_{2,t} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2) \end{cases}$$



# Latent variables as Hidden Climate Indices (HCI)

Variables affected by THE SAME climate index are dependant

$$\begin{cases} Y_{1,t} = \lambda_1 \tau_t + \varepsilon_{1,t} \\ Y_{2,t} = \lambda_2 \tau_t + \varepsilon_{2,t} \\ \varepsilon_{1,t}, \varepsilon_{2,t} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2) \end{cases}$$



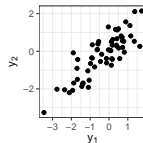
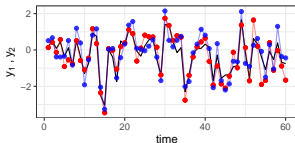
In practice, the time series  $\tau_t$  is unknown (it is *hidden*).



# Latent variables as Hidden Climate Indices (HCI)

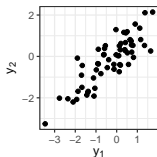
Variables affected by THE SAME climate index are dependant

$$\begin{cases} Y_{1,t} = \lambda_1 \tau_t + \varepsilon_{1,t} \\ Y_{2,t} = \lambda_2 \tau_t + \varepsilon_{2,t} \\ \varepsilon_{1,t}, \varepsilon_{2,t} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2) \end{cases}$$



In practice, the time series  $\tau_t$  is unknown (it is *hidden*).

Approach 1: explicitly modeling dependence

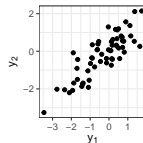
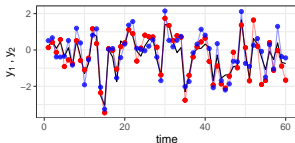


$$\begin{pmatrix} Y_{1,t} \\ Y_{2,t} \end{pmatrix} \stackrel{iid}{\sim} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \sigma^2 \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right)$$

# Latent variables as Hidden Climate Indices (HCI)

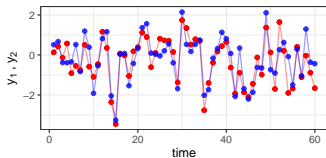
Variables affected by THE SAME climate index are dependant

$$\begin{cases} Y_{1,t} = \lambda_1 \tau_t + \varepsilon_{1,t} \\ Y_{2,t} = \lambda_2 \tau_t + \varepsilon_{2,t} \\ \varepsilon_{1,t}, \varepsilon_{2,t} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2) \end{cases}$$



In practice, the time series  $\tau_t$  is unknown (it is *hidden*).

Approach 2: uncovering the hidden climate index

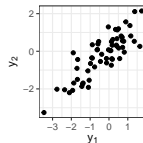
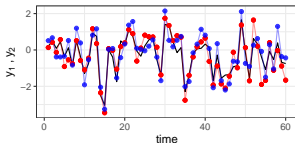


$$\begin{pmatrix} Y_{1,t} \\ Y_{2,t} \end{pmatrix} \stackrel{i}{\sim} \mathcal{N} \left( \begin{pmatrix} \lambda_1 \tau_t \\ \lambda_2 \tau_t \end{pmatrix}, \sigma^2 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right)$$

# Latent variables as Hidden Climate Indices (HCI)

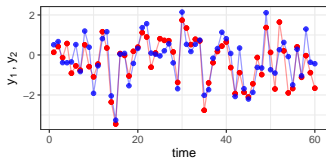
Variables affected by THE SAME climate index are dependant

$$\begin{cases} Y_{1,t} = \lambda_1 \tau_t + \varepsilon_{1,t} \\ Y_{2,t} = \lambda_2 \tau_t + \varepsilon_{2,t} \\ \varepsilon_{1,t}, \varepsilon_{2,t} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2) \end{cases}$$



In practice, the time series  $\tau_t$  is unknown (it is *hidden*).

Approach 2: uncovering the hidden climate index



$$\begin{pmatrix} Y_{1,t} \\ Y_{2,t} \end{pmatrix} \stackrel{i}{\sim} \mathcal{N} \left( \begin{pmatrix} \lambda_1 \tau_t \\ \lambda_2 \tau_t \end{pmatrix}, \sigma^2 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right)$$

**Objective:** any distribution, many sites, several variables, several HCIs

## Parent distributions of data

$$Y_v(s, t) \sim \mathcal{D}_v (\boldsymbol{\theta}_v(s, t))$$

# General modeling framework

## Parent distributions of data

$$Y_v(s, t) \sim \underbrace{\mathcal{D}_v}_{\text{any distribution, variable-specific}}(\theta_v(s, t))$$

any distribution, variable-specific

# General modeling framework

## Parent distributions of data

$$Y_v(s, t) \sim \mathcal{D}_v \left( \underbrace{\theta_v(s, t)} \right)$$

varies in space and time

# General modeling framework

## Parent distributions of data

$$Y_v(s, t) \sim \mathcal{D}_v(\theta_v(s, t))$$

## Parameters vary in space and time

$$\theta_v(s, t) = \lambda_{0,v}(s) + \lambda_{1,v}(s)\tau_1(t)$$

# General modeling framework

## Parent distributions of data

$$Y_v(s, t) \sim \mathcal{D}_v (\boldsymbol{\theta}_v(s, t))$$

## Parameters vary in space and time

$$\boldsymbol{\theta}_v(s, t) = \lambda_{0,v}(s) + \underbrace{\lambda_{1,v}(s)\tau_1(t)}$$

SAME HCI affects all sites/variables



# General modeling framework

## Parent distributions of data

$$Y_v(s, t) \sim \mathcal{D}_v(\theta_v(s, t))$$

## Parameters vary in space and time

$$\theta_v(s, t) = \lambda_{0,v}(s) + \lambda_{1,v}(s)\tau_1(t) \underbrace{+ \cdots + \lambda_{K,v}(s)\tau_K(t)}_{\text{more HClS}}$$

# General modeling framework

## Parent distributions of data

$$Y_v(s, t) \sim \mathcal{D}_v(\boldsymbol{\theta}_v(s, t))$$

## Parameters vary in space and time

$$\underbrace{g(\boldsymbol{\theta}_v(s, t))}_{\text{link function}} = \lambda_{0,v}(s) + \lambda_{1,v}(s)\tau_1(t) + \cdots + \lambda_{K,v}(s)\tau_K(t)$$

# General modeling framework

## Parent distributions of data

$$Y_v(s, t) \sim \mathcal{D}_v(\boldsymbol{\theta}_v(s, t))$$

## Parameters vary in space and time

$$g(\boldsymbol{\theta}_v(s, t)) = \lambda_{0,v}(s) + \lambda_{1,v}(s)\tau_1(t) + \cdots + \lambda_{K,v}(s)\tau_K(t)$$

## Temporal and spatial Gaussian processes for HCLs and their effects

$$\tau_k(t) \sim \mathcal{G}(\boldsymbol{\mu}_\tau, \boldsymbol{\Sigma}_\tau);$$

# General modeling framework

## Parent distributions of data

$$Y_v(s, t) \sim \mathcal{D}_v(\theta_v(s, t))$$

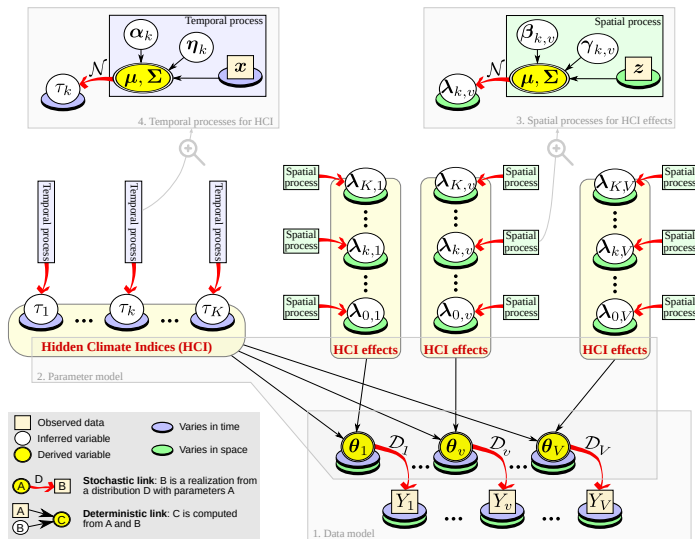
## Parameters vary in space and time

$$g(\theta_v(s, t)) = \lambda_{0,v}(s) + \lambda_{1,v}(s)\tau_1(t) + \cdots + \lambda_{K,v}(s)\tau_K(t)$$

## Temporal and spatial Gaussian processes for HCLs and their effects

$$\tau_k(t) \sim \mathcal{G}(\mu_\tau, \Sigma_\tau); \quad \lambda_{k,v}(s) \sim \mathcal{G}(\mu_\lambda, \Sigma_\lambda)$$

# Schematics of an HCI model



## Possible interpretations

- $g(\theta) = \lambda_0 + \sum \lambda_k \tau_k$  is similar to GLM... but with hidden covariates!

## Possible interpretations

- $g(\theta) = \lambda_0 + \sum \lambda_k \tau_k$  is similar to GLM... but with hidden covariates!
- HCIs and their effects  $\approx$  principal components and their loadings
- HCI model  $\approx$  non-Gaussian probabilistic PCA
- see *Probabilistic Machine Learning* by Kevin P. Murphy (2022)

MCMC sampling from the posterior, Renard & Thyer (2019)



MCMC sampling from the posterior, Renard & Thyer (2019)

## Difficulty 1: identifiability constraints

- HCIs should be orthonormal
- Stepwise inference: one component at a time
- Leads to 2 simpler constraints: each HCI has mean 0 and variance 1
- Possible solution: 'Givens Representation' of Pourzanjani et al. (2021)

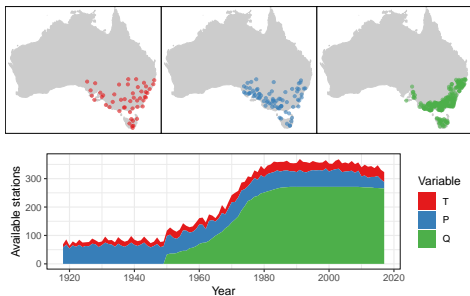
MCMC sampling from the posterior, Renard & Thyer (2019)

## Difficulty 1: identifiability constraints

- HCIs should be orthonormal
- Stepwise inference: one component at a time
- Leads to 2 simpler constraints: each HCI has mean 0 and variance 1
- Possible solution: 'Givens Representation' of Pourzanjani et al. (2021)

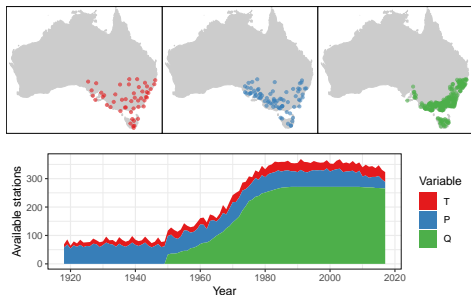
## Difficulty 2: dimensionality

- Dimension of the posterior grows with both #sites and #time steps.
- No big deal for the likelihood (thank you conditional independence!)
- Bottleneck = covariance matrices in Gaussian hyperdistributions
- One solution: nearest-neighbor Gaussian process of Datta et al. (2016)



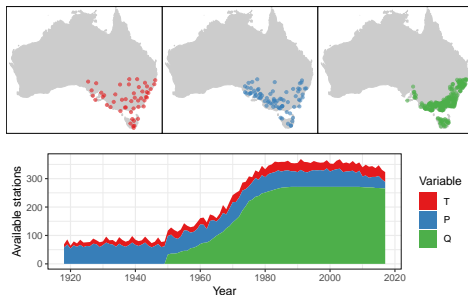
Number of heatwaves

$$\begin{cases} T_n(s, t) \sim \mathcal{P}(\mu(s, t)) \\ \log(\mu(s, t)) = \lambda_{T_n,0}(s) + \sum_{k=1}^3 \lambda_{T_n,k}(s) \tau_k(t) \end{cases}$$



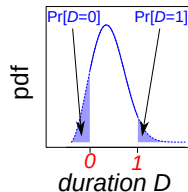
## Heatwave intensities

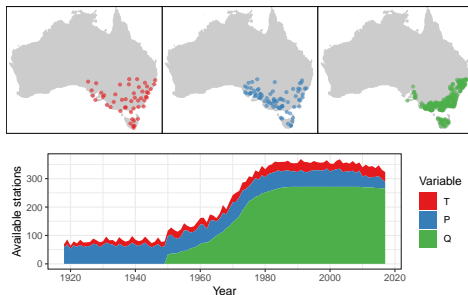
$$\begin{cases} T_x(s, t) \sim \mathcal{GPD}(0, \sigma(s, t), \xi(s)) \\ \log(\sigma(s, t)) = \lambda_{T_x, 0}(s) + \sum_{k=1}^3 \lambda_{T_x, k}(s) \tau_k(t) \end{cases}$$



Drought duration

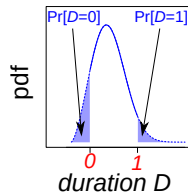
$$\begin{cases} Qd(s, t) \sim \mathcal{N}(\mu(s, t), \sigma(s)) \\ \mu(s, t) = \lambda_{Qd,0}(s) + \sum_{k=1}^3 \lambda_{Qd,k}(s) \tau_k(t) \end{cases}$$



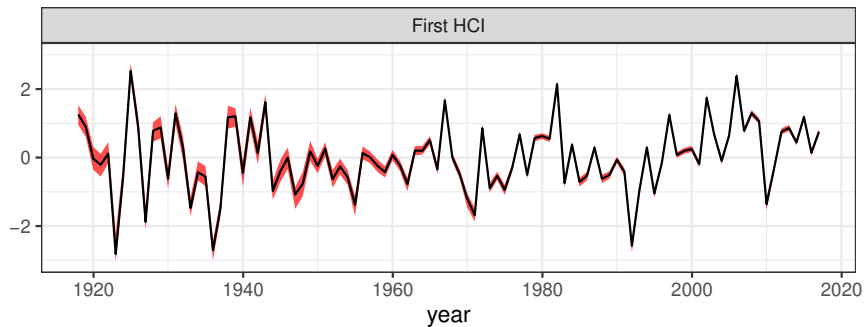


Dry-day duration

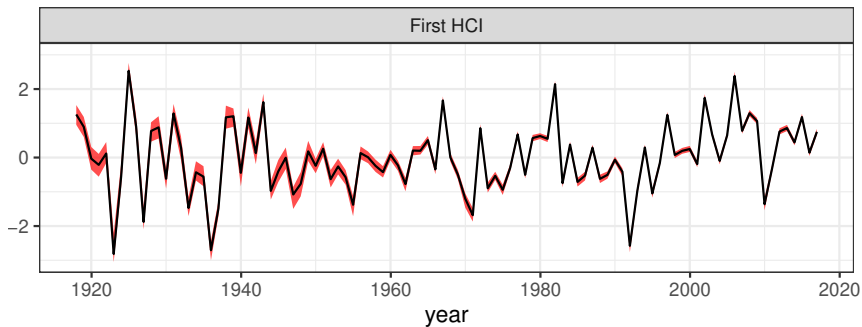
$$\begin{cases} Pd(s, t) \sim \mathcal{N}(\mu(s, t), \sigma(s)) \\ \mu(s, t) = \lambda_{Pd,0}(s) + \sum_{k=1}^3 \lambda_{Pd,k}(s) \tau_k(t) \end{cases}$$



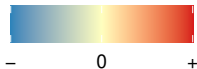
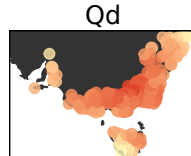
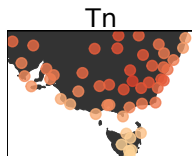
# First HCI and its effects



# First HCI and its effects

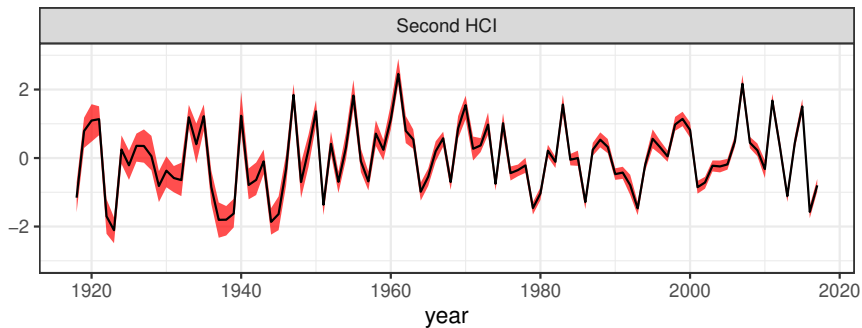


Effect on...

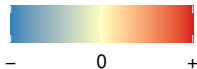
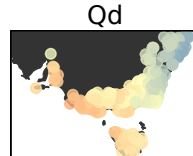
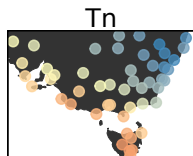




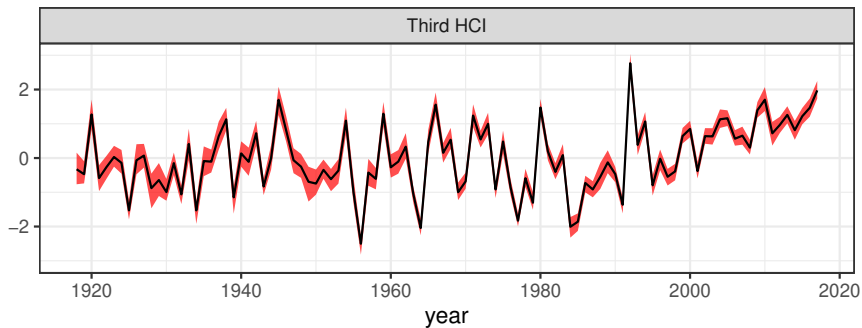
# Second HCI and its effects



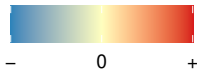
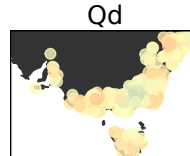
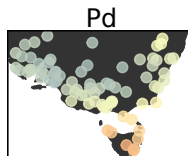
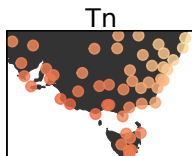
Effect on...



# Third HCI and its effects

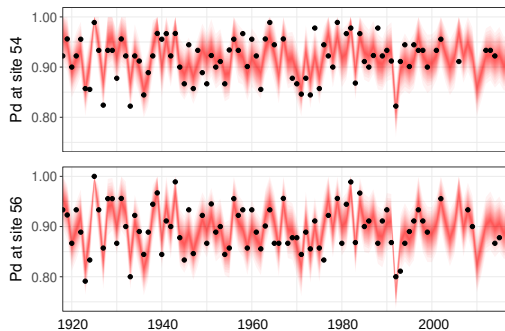


Effect on...



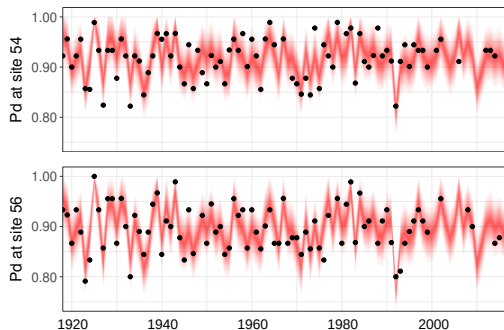
# Probabilistic predictions

## Time-varying distributions

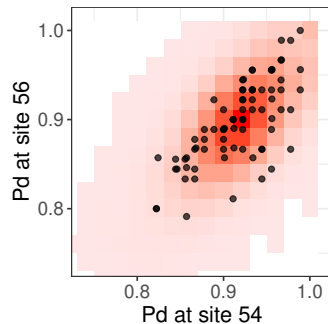


# Probabilistic predictions

## Time-varying distributions

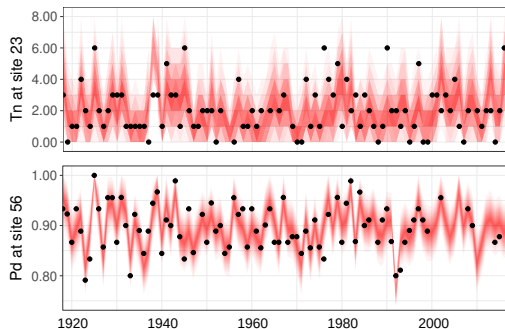


## Joint bivariate distribution

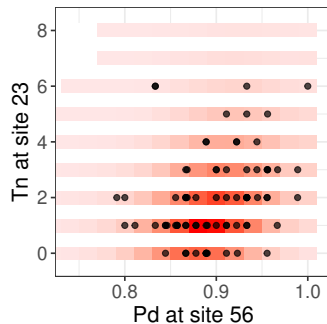


# Probabilistic predictions

## Time-varying distributions

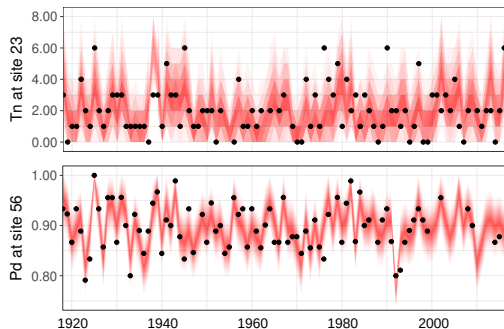


## Joint bivariate distribution

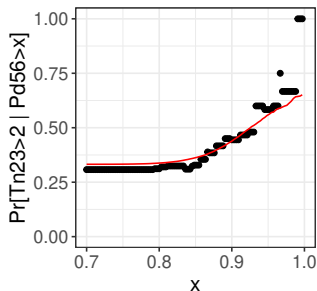


# Probabilistic predictions

## Time-varying distributions

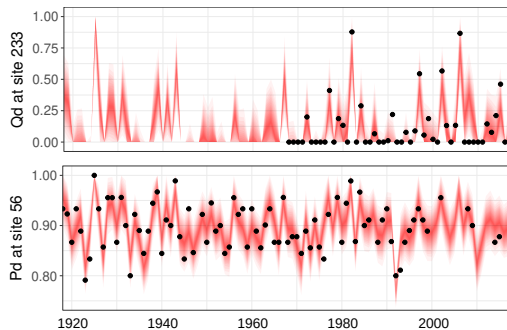


## Conditional probabilities



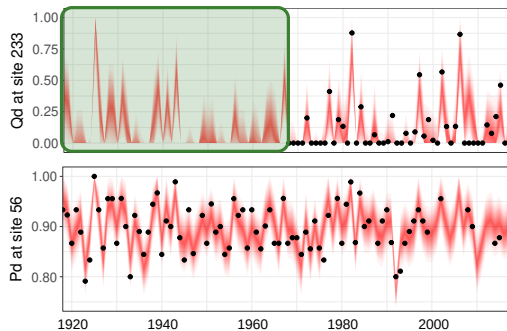
# Probabilistic predictions

## Time-varying distributions



# Probabilistic predictions

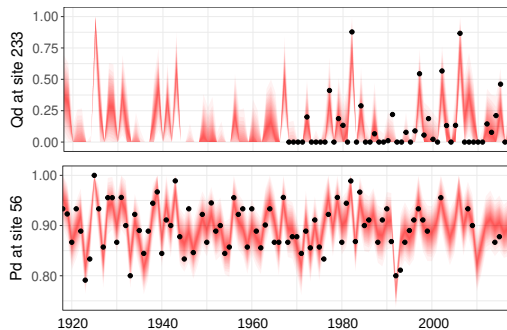
## Time-varying distributions



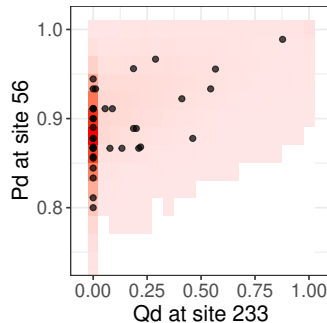


# Probabilistic predictions

## Time-varying distributions

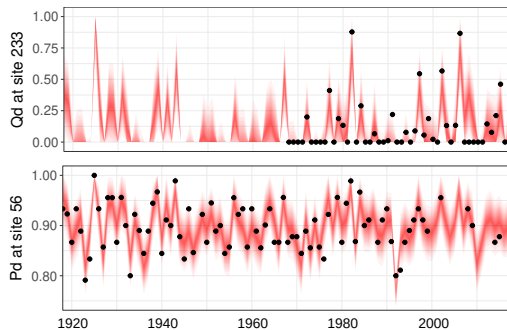


## Joint bivariate distribution

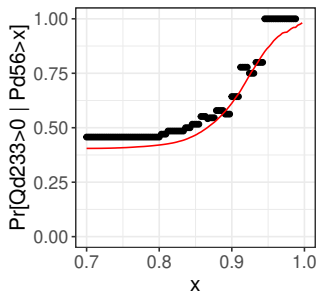


# Probabilistic predictions

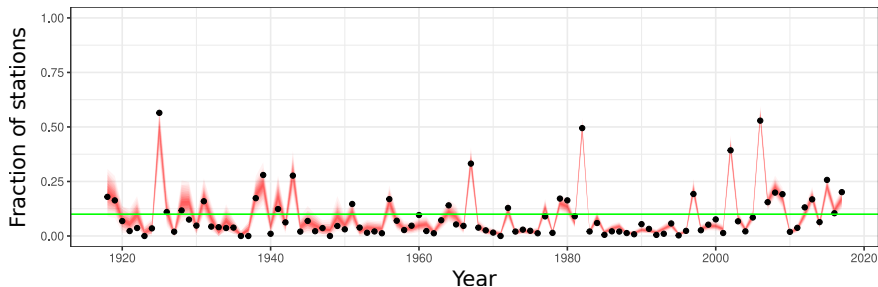
## Time-varying distributions



## Conditional probabilities



## Fraction of stations exceeding a 10-year event



⇒ Consequences in terms of risk management

## Summary: the HCI modeling framework

- A general probabilistic model for multi-variable space-time data
- Based on hidden climate indices extracted from the target data
- Flexible: can handle a wide range of hydro-meteorological datasets

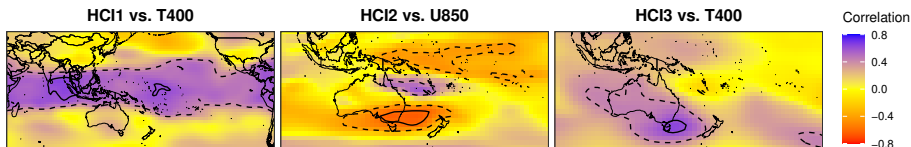
## Summary: the HCI modeling framework

- A general probabilistic model for multi-variable space-time data
- Based on hidden climate indices extracted from the target data
- Flexible: can handle a wide range of hydro-meteorological datasets

## Other noticeable results [not shown here]

- Probabilistic predictions are reliable, including in cross-validation
- HCIs  $\neq$  standard indices such as NINO, SAM, IOD, etc.
- Replacing HCIs with standard indices underestimates dependence

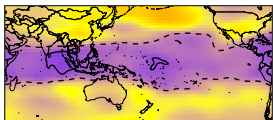
## Predicting HCIs from large-scale climate variables?



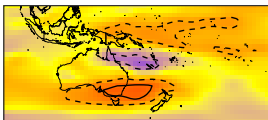
⇒ Downscaling device for past reconstructions, seasonal forecasting or future projections

## Predicting HCIs from large-scale climate variables?

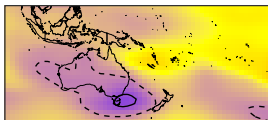
HCI1 vs. T400



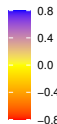
HCI2 vs. U850



HCI3 vs. T400



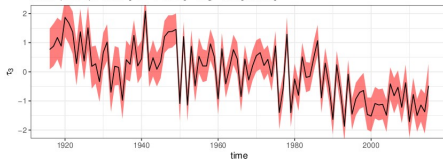
Correlation



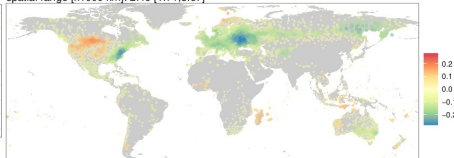
⇒ Downscaling device for past reconstructions, seasonal forecasting or future projections

## Application to global floods and extreme precipitation

trend (x0.01): -2.17 [-2.7;-1.63] – lag-1: 0 [0;0.12]



spatial range [x1000 km]: 2.46 [1.71;3.67]



# Thank you!

Renard & Thyer (2019). Revealing Hidden Climate Indices from the Occurrence of Hydrologic Extremes. *Water Resources Research*.

Renard et al. (2021?). A Hidden Climate Indices Modeling Framework for Multi-Variable Space-Time Data. *Submitted to Water Resources Research*.



This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 835496



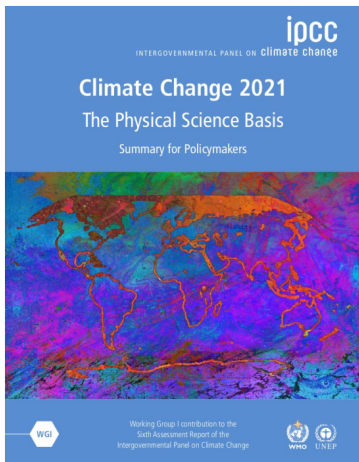
<https://globxblog.inrae.fr/>



<https://github.com/STooDs-tools>

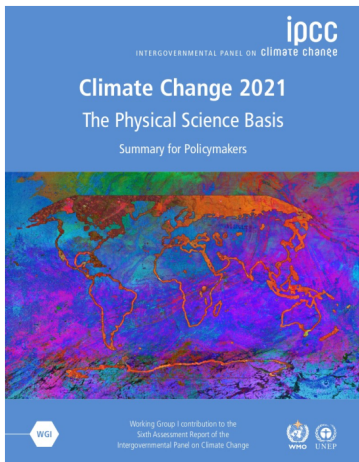


# Application 2: context



*" The frequency and intensity of heavy precipitation events have increased since the 1950s over most land area for which observational data are sufficient for trend analysis (high confidence), and human-induced climate change is likely the main driver"*

# Application 2: context



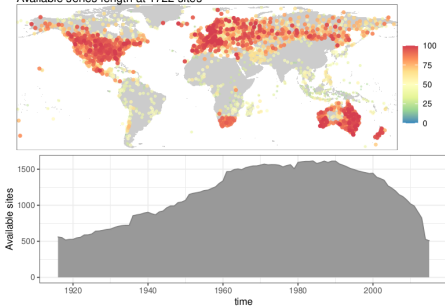
*" The frequency and intensity of heavy precipitation events have increased since the 1950s over most land area for which observational data are sufficient for trend analysis (high confidence), and human-induced climate change is likely the main driver"*

*" Confidence about peak flow trends over past decades on the global scale is low, but there are regions experiencing increases [...] and regions experiencing decreases [...]"*

# Global datasets for hydrologic extremes

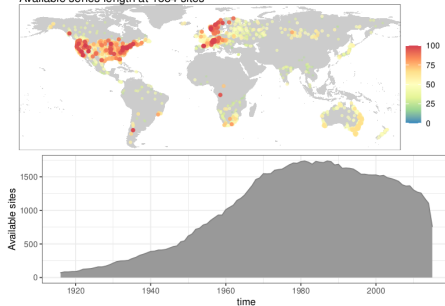
## Precipitation: Hadex 2+3

Available series length at 1722 sites



## Streamflow: GSIM

Available series length at 1884 sites



## Objectives

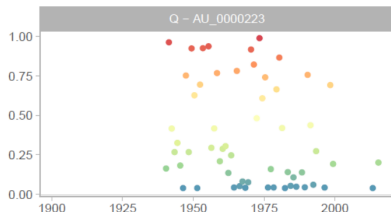
- Look for trends, low-frequency variability and teleconnections for both P and Q extremes
- Analyze long 100-year period (vs. a typical 50-year)
- Attempt at predicting extreme P/Q from large-scale climate

# Analyzed variables

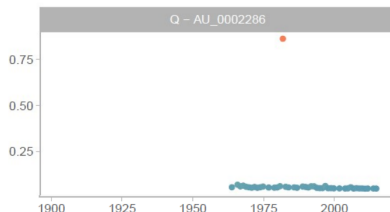
Non-exceedance probability ( $\Leftrightarrow$  return period) of the largest event of the season

**Example:** Maximum streamflow in December-January-February for 2 Australian stations

Barker Creek at Brooklands (QLD)



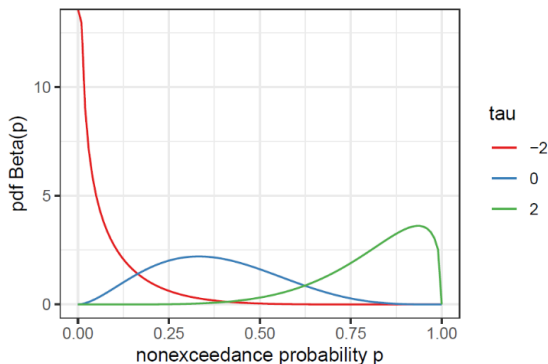
Clarke Brooke at Hillview Farm (WA)



# Model

Beta distribution reparameterized in terms of mean  $\mu$  and precision  $\gamma$

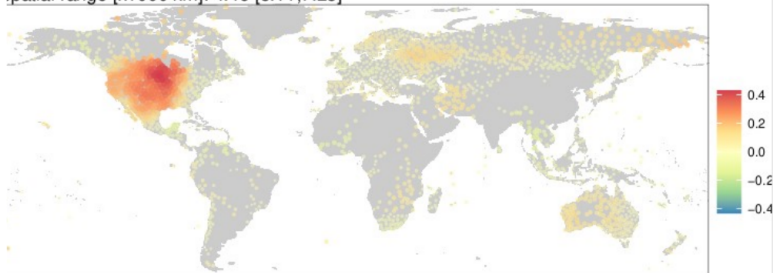
$$\begin{cases} Y_v(s, t) \sim \text{Beta}_v(\mu_v(s, t), \gamma_v(s)) \\ \text{logit}(\mu_v(s, t)) = \lambda_{v,0}(s) + \lambda_{v,1}(s)\tau_1(t) + \dots + \lambda_{v,K}(s)\tau_K(t) \end{cases}$$



Beta distribution reparameterized in terms of mean  $\mu$  and precision  $\gamma$

$$\begin{cases} Y_v(s, t) \sim \text{Beta}_v(\mu_v(s, t), \gamma_v(s)) \\ \text{logit}(\mu_v(s, t)) = \lambda_{v,0}(s) + \lambda_{v,1}(s)\tau_1(t) + \dots + \lambda_{v,K}(s)\tau_K(t) \\ \lambda_k \sim \text{NNGP}(\beta, V); V_{i,j} = \nu_0^2 \exp(-d_{i,j}/\nu_1) \end{cases}$$

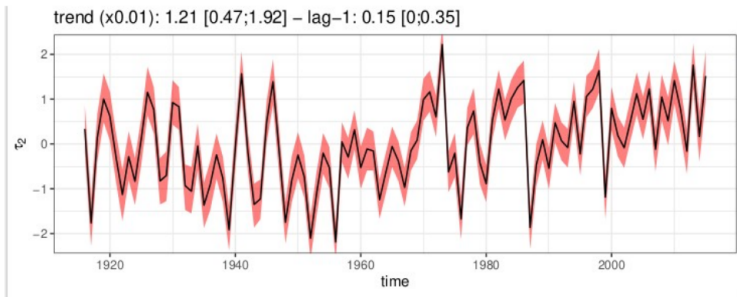
spatial range [x1000 km]: 4.46 [3.11;7.23]



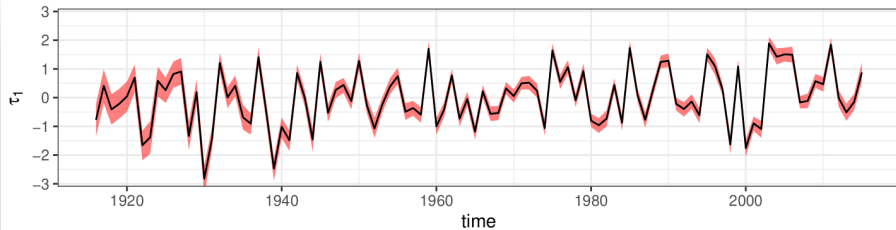
# Model

Beta distribution reparameterized in terms of mean  $\mu$  and precision  $\gamma$

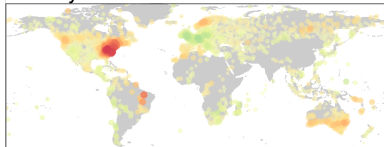
$$\begin{cases} Y_v(s, t) \sim \text{Beta}_v(\mu_v(s, t), \gamma_v(s)) \\ \text{logit}(\mu_v(s, t)) = \lambda_{v,0}(s) + \lambda_{v,1}(s)\tau_1(t) + \dots + \lambda_{v,K}(s)\tau_K(t) \\ \lambda_k \sim \text{NNGP}(\beta, V); V_{i,j} = \nu_0^2 \exp(-d_{i,j}/\nu_1) \\ \tau_k \sim \text{NNGP}(m, W); W_{i,j} = u_0^2 \exp(-d_{i,j}/u_1); m_t = \eta_0 + \eta_1 t \end{cases}$$



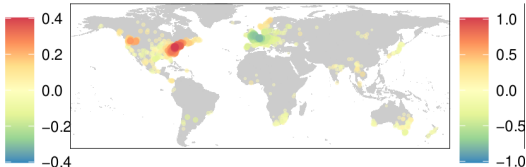
# Step 1: identify components



Rx1day

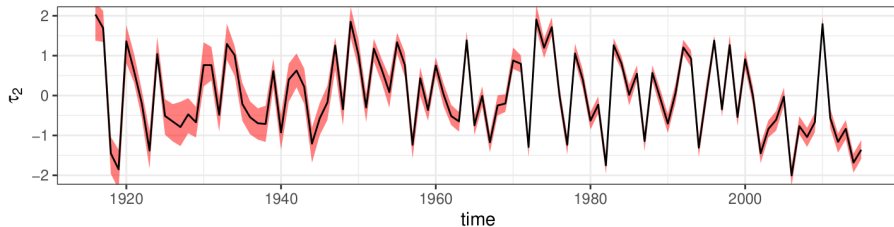


Qx

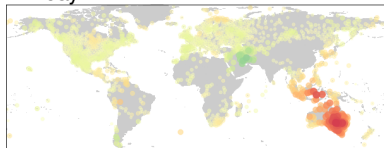




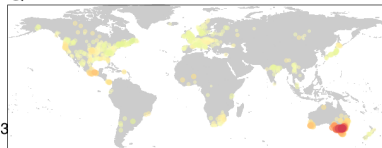
# Step 1: identify components



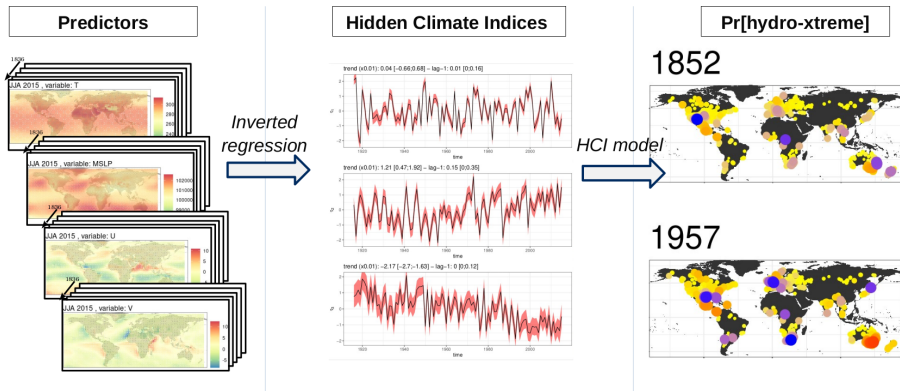
Rx1day



Qx



# Step 2: past reconstructions from 1836



# Thank you!

Renard & Thyer (2019). Revealing Hidden Climate Indices from the Occurrence of Hydrologic Extremes. *Water Resources Research*.

Renard et al. (2021?). A Hidden Climate Indices Modeling Framework for Multi-Variable Space-Time Data. *Submitted to Water Resources Research*.



This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 835496



<https://globxblog.inrae.fr/>



<https://github.com/STooDs-tools>

## Method: inverted regression

**Step 1:**  $w(s, t)$ : climate field at time  $t$  and location  $s$

$\hat{\tau}_k(t)$ : estimated HCI's (from previous analysis)

Goal: estimate  $\psi_k(s)$ 's in:

$$w(s, t) = \psi_0(s) + \psi_1(s)\hat{\tau}_1(t) + \dots + \psi_K(s)\hat{\tau}_K(t) + \varepsilon(s, t)$$

**Step 2:**  $w(s, t^*)$ : climate field at time  $t^*$  and location  $s$

$\hat{\psi}_k(s)$ : estimated from previous step

Goal: estimate  $\tau_k(t^*)$ 's in:

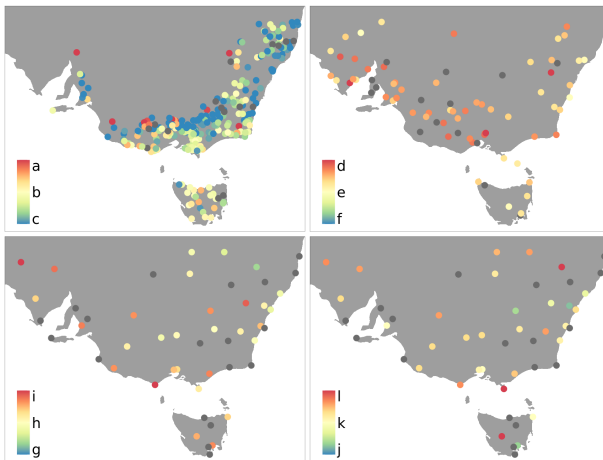
$$w(s, t^*) = \psi_0(s) + \hat{\psi}_1(s)\tau_1(t^*) + \dots + \hat{\psi}_K(s)\tau_K(t^*) + \varepsilon(s, t^*)$$

**Alternatives:** LASSO, RIDGE and other form of penalised regression, but first attempts inconclusive

# Motivating dataset: a few years

## Droughts and heatwaves in South-East Australia during summer 1981

Top-left: river drought duration Top-right: number of dry days Bottom-left: heatwave intensity Bottom-right: number of heatwaves

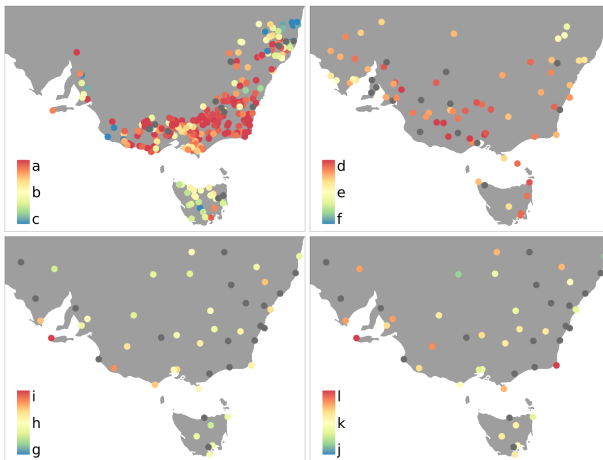


Source: Bureau of Meteorology

# Motivating dataset: a few years

## Droughts and heatwaves in South-East Australia during summer 1982

Top-left: river drought duration Top-right: number of dry days Bottom-left: heatwave intensity Bottom-right: number of heatwaves

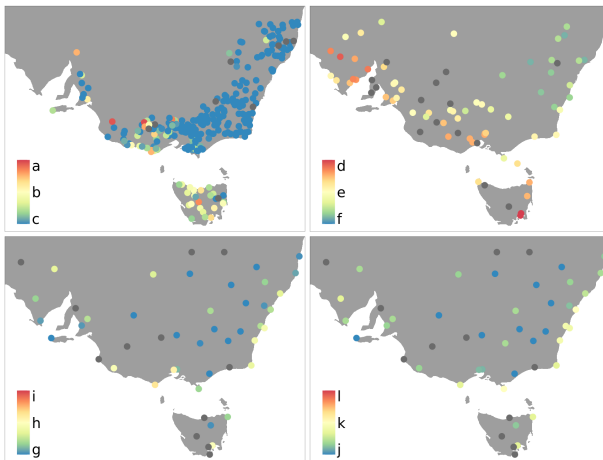


Source: Bureau of Meteorology

# Motivating dataset: a few years

## Droughts and heatwaves in South-East Australia during summer 1983

Top-left: river drought duration Top-right: number of dry days Bottom-left: heatwave intensity Bottom-right: number of heatwaves



Source: Bureau of Meteorology