

Data Analytics

Course: 18-785

Patrick McSharry

patrick@mcsharry.net

www.mcsharry.net

[Twitter: @patrickmcsharry](https://twitter.com/patrickmcsharry)

Spring 2023

ICT Center of Excellence
Carnegie Mellon University

A2Q3

- Ramps can be both positive and negative
- Upward move is a positive ramp and downward move is a negative ramp
- Separate into two groups based on sign of the ramp (positive and negative ramps) using indices: $\text{indn} = \text{find}(r < 0)$; $\text{indp} = \text{find}(r > 0)$.
- Investigate each distribution separately
- This gives the two sides of the graph showing ramps in wind power

Data Analytics

WEEK 4A

Course outline

Week	Lecture A	Lecture B
1	Data Analytics	Weather forecasting
2	Renewable energy	Wind energy
3	Solar energy	Demand forecasting
4	Risk	Extreme events
5	Health	Biomedicine
6	Early warning systems	Economic forecasting

Today's Lecture

No.	Activity	Description	Time
1	Challenge	Risk management	10
2	Discussion	Global risks	10
3	Case study	Natural disasters	10
4	Analysis	Disaster trends	20
5	Demo	Risk communication	20
6	Q&A	Questions and feedback	10

Risk definition

- Risk is potential of losing something of value.
- Examples of value are health, well-being or financial wealth.
- Value can be gained or lost as a result of particular events or actions.
- Risk is closely associated with uncertainty.
- Certain actions can lead to unpredictable outcomes that can affect value.

Poll

- For 15-49 year olds, which of the following causes the highest number of deaths globally?
 - a) Cardiovascular disease
 - b) Cancers
 - c) HIV/AIDS
 - d) Road accidents

Slido.com #60339

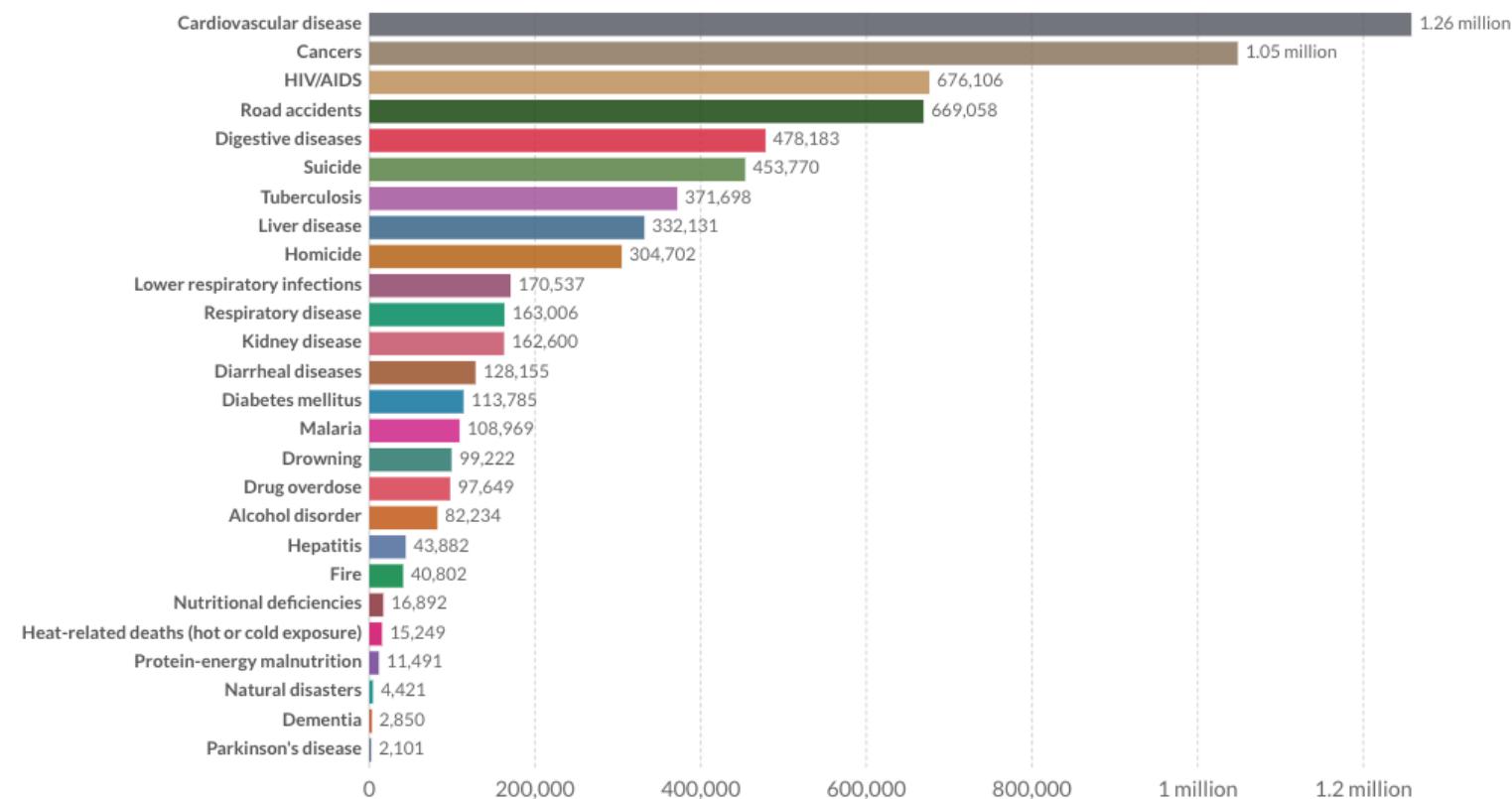
Causes of death for 15-49 yrs

Causes of deaths for 15 to 49 year olds, World, 2017

Annual number of deaths – by cause – for people aged 15 to 49 years old.

Our World
in Data

Change country



Source: IHME, Global Burden of Disease (GBD)

OurWorldInData.org/causes-of-death • CC BY

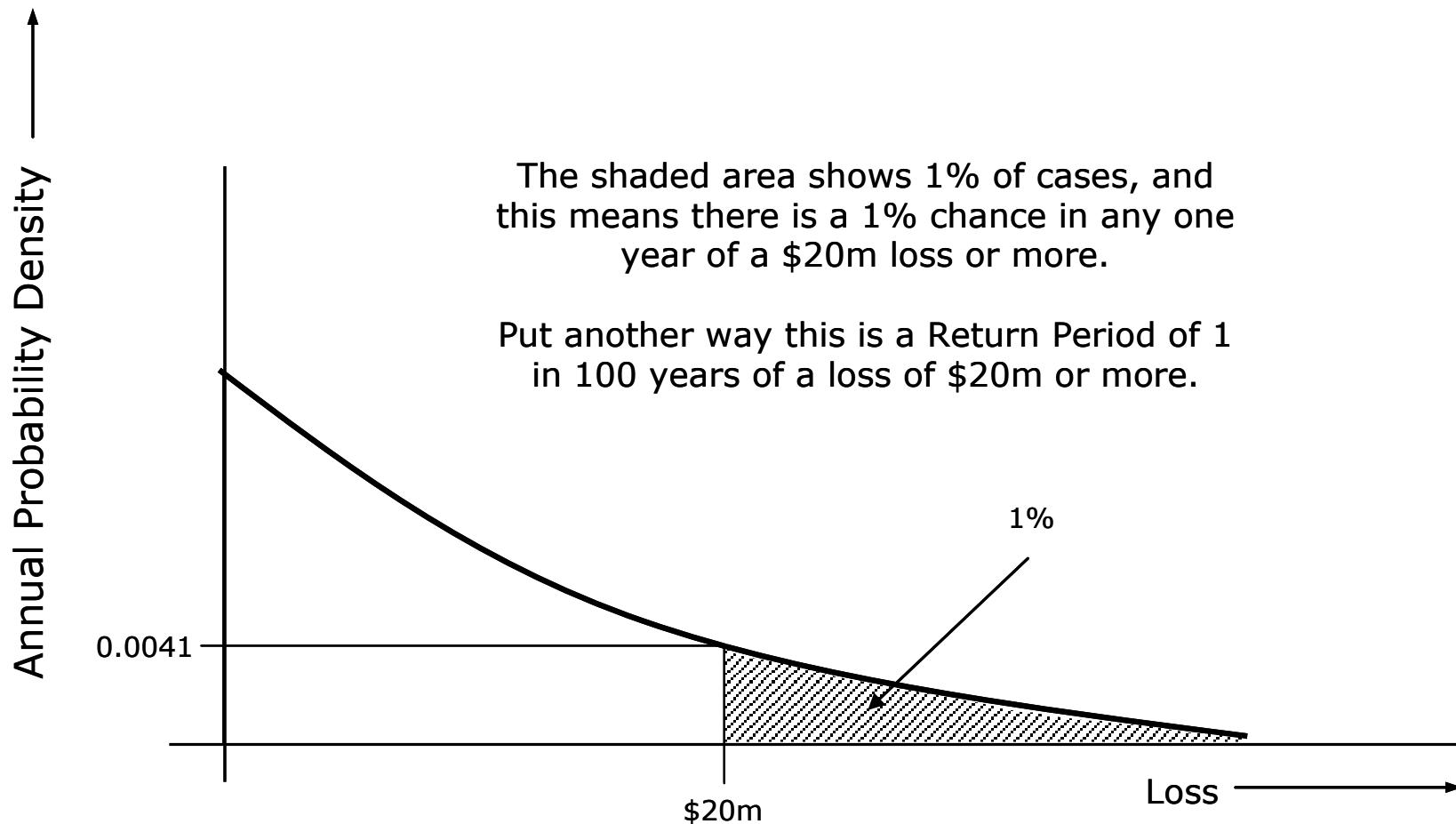
Risk measurement

- Risk is “the probability that a particular adverse event occurs during a stated period of time, or results from a particular challenge” - Royal Society 1992 *Risk: Analysis, Perception and Management*
- What is the chance of a severe drought in 2020? [Loss and chance]
- How much did you personally lose due to a particular drought in year x? [Loss from an event]
- Risk involves harm, chance and time.
- Risk quantification: hazard, exposure and vulnerability.

Regulation - acceptable risk

- Pillar I of Solvency II requires (since 2016) insurers operating in EU to use a quantitative approach for calculating their solvency capital ratio (SCR).
- The SCR is the capital an insurer should hold to meet its obligations over the next year with a probability of at least 99.5%.
- This corresponds to having sufficient capital to withstand the loss from a 1 in 200 year event

Loss Probability Curve



Risk management

- **ISO 31000 - Risk management**
- Risks affecting organizations can have consequences in terms of economic performance and professional reputation, as well as environmental, safety and societal outcomes.
- Therefore, managing risk effectively helps organizations to perform well in an environment full of uncertainty.

Risk management components



Risk and resilience

- **Resilience to shocks:** capacity to absorb, recover and transition to some preferable regime
- Disaster risk management can be cost effective
- Information: public access to flood risk maps (UK), hurricane risk (Florida), seismic fault maps (Tehran)
- Market: appropriate insurance instruments and price signals
- Sharing data: forecasts, early warning systems, evacuation and response systems

Risk assessment

- Risk associated with a particular event can be quantified by calculating the likely impact of this event and the probability of occurrence.
- The risk is then given as the product of the impact (or severity) and probability (or likelihood):

$$\text{Risk} = \text{Impact} \times \text{Probability}$$

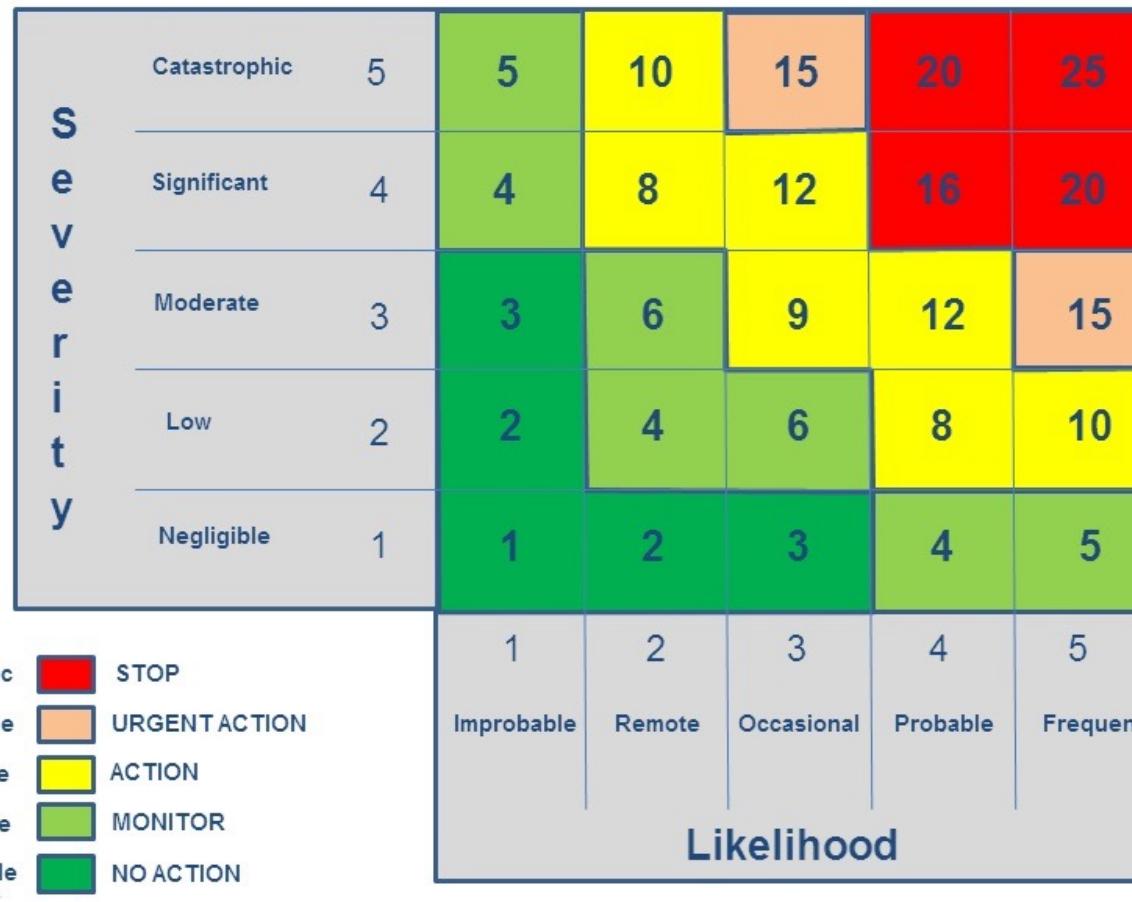
- The risk and the appropriate risk management action can then be calculated using a risk matrix.

Risk matrix

Impact\Likelihood		Rare	Unlikely	Possible	Likely	Certain
	Loss\Prob	<0.0001	0.001	0.01	0.1	1
Very high	<\$10,000,000	Yellow	Orange	Orange	Red	Red
High	<\$1,000,000	Green	Yellow	Orange	Orange	Red
Medium	<\$100,000	Green	Green	Yellow	Orange	Orange
Low	<\$10,000	Dark Green	Light Green	Light Green	Yellow	Orange
Very low	<\$1000	Dark Green	Dark Green	Light Green	Light Green	Yellow

Risk matrix

Risk Rating = Likelihood x Severity



Indonesian tsunami 2004

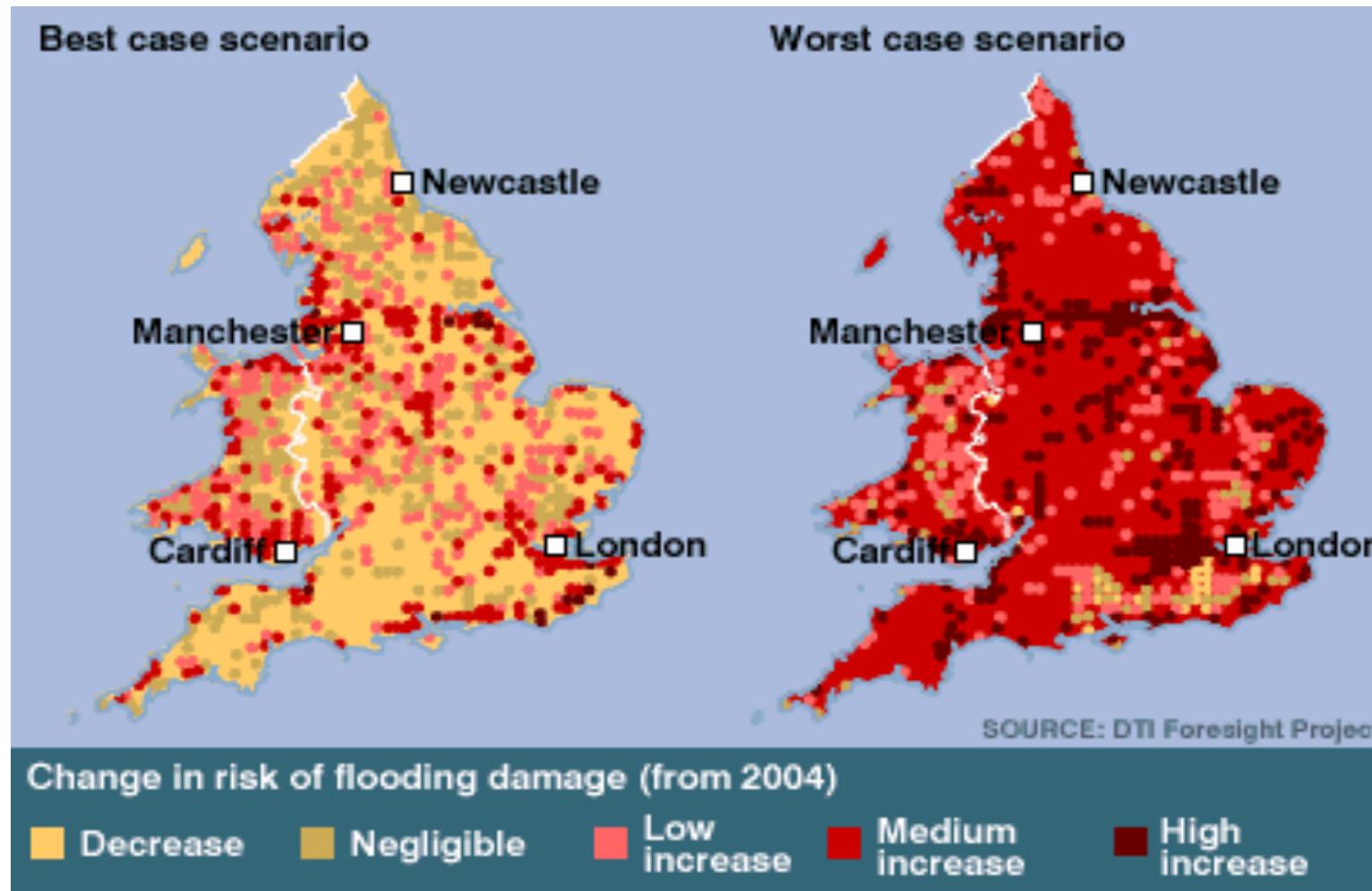


Before



After

Possible flooding in the UK by 2080



Cape Town faces Day Zero



What happens when the city turns off the taps?

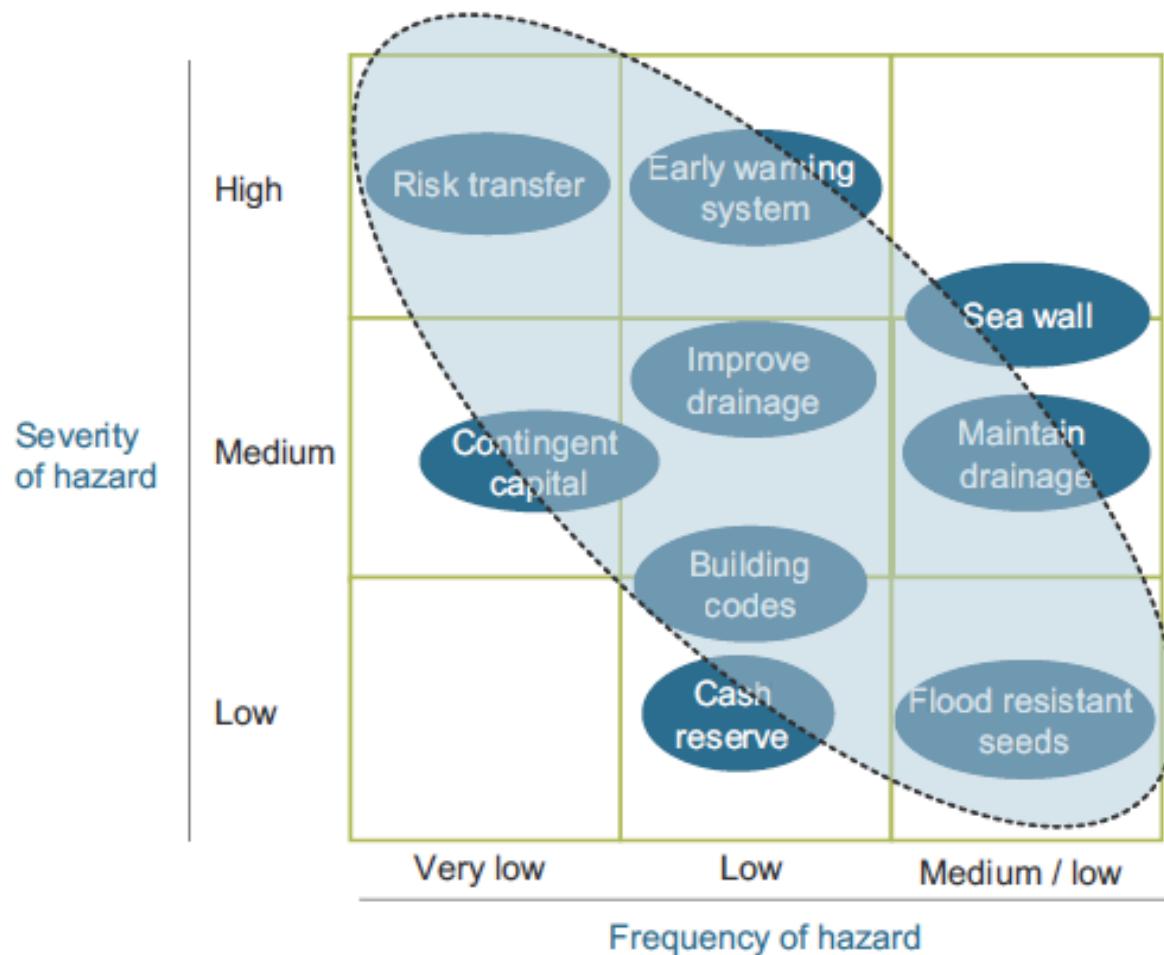
In 10 weeks engineers will turn off water for a million homes as this South African city reacts to a one-in-384-year drought. The rich are digging boreholes, more are panic-buying bottled water, and the army is on standby.

Socio-economic risk factors



- Urbanisation leading to value concentration in coastal cities
- Low quality housing and high concentrations of people living in mega cities in developing countries
- Inadequate infrastructure and institutions to deal with catastrophes
- Degradation of natural systems:
 - Reduction of mangrove swamps due to farming and pollution
 - Building on flood plains in the UK
 - Building on the coast exposed to sea level rise

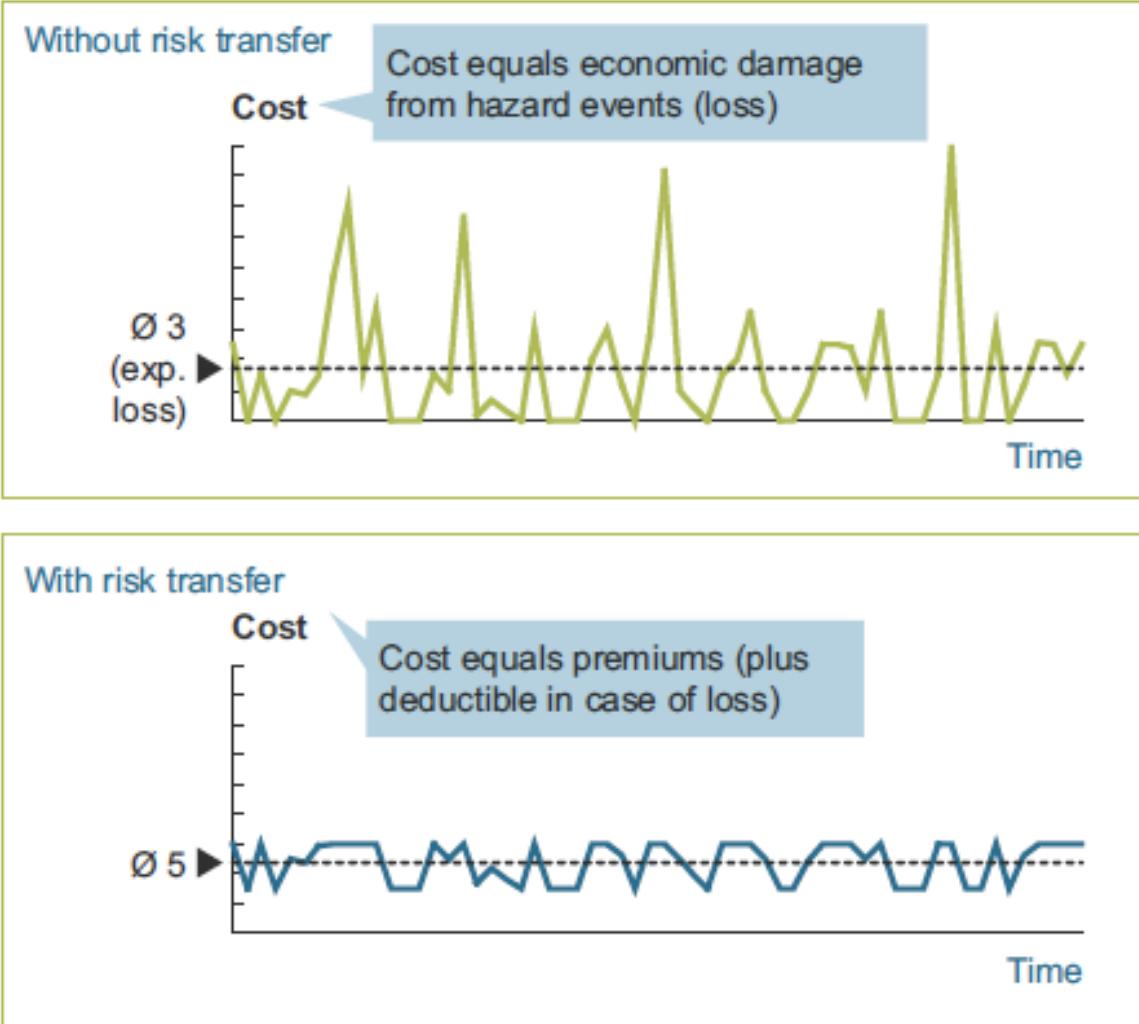
Climate risk portfolio – cost benefit analysis



Risk transfer mechanisms

- Objective is to cap losses and smooth the costs of catastrophe events for individuals, corporations and governments
- Traditional indemnity-based insurance
- Parametric index-based insurance (livestock, weather)
- Chicago Board of Trade (CBOT) has been trading futures on catastrophe loss indices and related options since Dec 1992
- Catastrophe bonds (e.g. Cat-Mex bond where Mexican government does not repay bond principal if a major earthquake were to hit)
- South Africa-based insurer African Risk Capacity is preparing for the launch of its climate change catastrophe bond facility. Extreme Climate Facility (XCF) will issue cat bonds to help African nations transfer climate change-related exposure to the capital markets.

Risk transfer



Risk transfer

Benefits

- Caps losses, protects livelihood from catastrophic events
- Smoothes costs, reduces volatility
- Increases willingness to invest
- Provides incentives ("price signals")

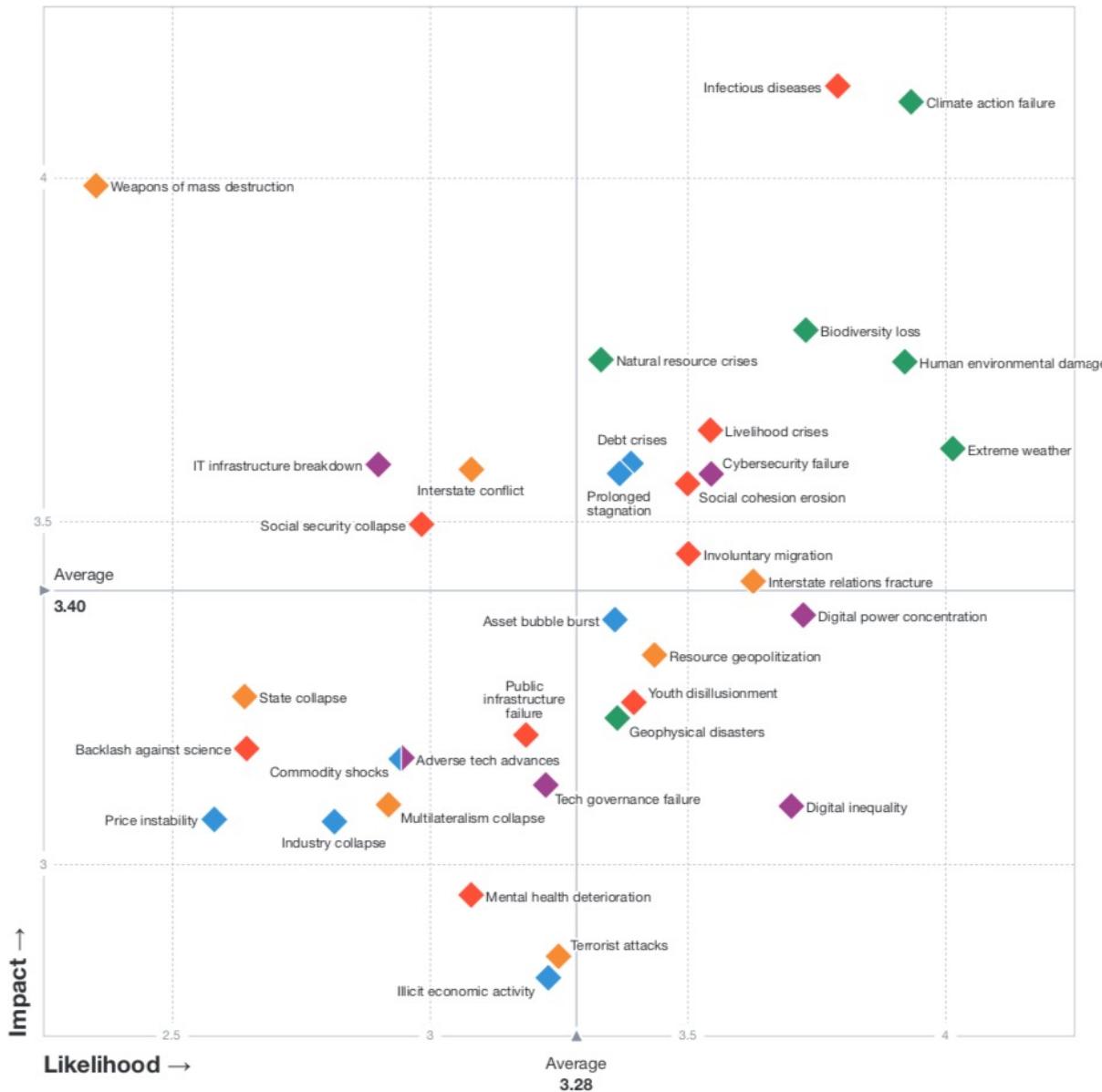
Costs

- Expected loss plus markup for production and distribution

Quiz

- According to the World Economic Forum (WEF, 2021), what is the biggest risk by likelihood facing our planet?
 - a) Livelihood crises
 - b) Climate action failure
 - c) Digital inequality
 - d) Extreme weather

Global Risks Landscape (2021)



Source: [World Economic Forum, 2021](#)

Global risks (WEF, 2021)

- Epidemics lead biggest short-term risks and this current trend could lead to future stagnation in advanced economies and lost potential in emerging and developing markets.
- Environmental threats at the top of the list for the fourth year in a row
- **Top-risks by likelihood:** Extreme Weather Event (1st place), Climate Action Failure (2nd place), Human Environmental Damage (3rd place)
- **Top-risks by impact:** Biodiversity Loss (4th place), Natural Resource Crises (5th place).
- Extreme weather events and natural disasters are often water related. In years to come, floods and droughts are set to strike harder and more often due to climate change.
- The most dramatic consequence of climate change is its impacts on the water cycle, which also means more unpredictable rainfall episodes. Therefore, extreme weather events cause the major biodiversity loss and ecosystem collapse.

Cloud Poll

- What risks would you like to quantify?
- **Slido.com #60339**

Catalogue of Disasters

- EM-DAT: International disaster database
- For a disaster to be entered into the database at least one of the following criteria must be fulfilled:
 - Ten (10) or more people reported killed.
 - Hundred (100) or more people reported affected.
 - Declaration of a state of emergency.
 - Call for international assistance.

Disaster Classification

Disaster Subgroup	Definition	Disaster Main Type
Geophysical	Events originating from solid earth	Earthquake, Volcano, Mass Movement (dry)
Meteorological	Events caused by short-lived/small to meso scale atmospheric processes (in the spectrum from minutes to days)	Storm
Hydrological	Events caused by deviations in the normal water cycle and/or overflow of bodies of water caused by wind set-up	Flood, Mass Movement (wet)
Climatological	Events caused by long-lived/meso to macro scale processes (in the spectrum from intra-seasonal to multi-decadal climate variability)	Extreme Temperature, Drought, Wildfire
Biological	Disaster caused by the exposure of living organisms to germs and toxic substances	Epidemic, Insect infestation, Animal Stampede

Natural disaster events

Number of recorded natural disaster events, All natural disasters, 1900 to 2019

The number of global reported natural disaster events in any given year. This includes those from drought, floods, extreme weather, extreme temperature, landslides, dry mass movements, wildfires, volcanic activity and earthquakes.

Our World
in Data

⇄ Change disaster category



Source: EMDAT (2020): OFDA/CRED International Disaster Database, Université catholique de Louvain – Brussels – Belgium

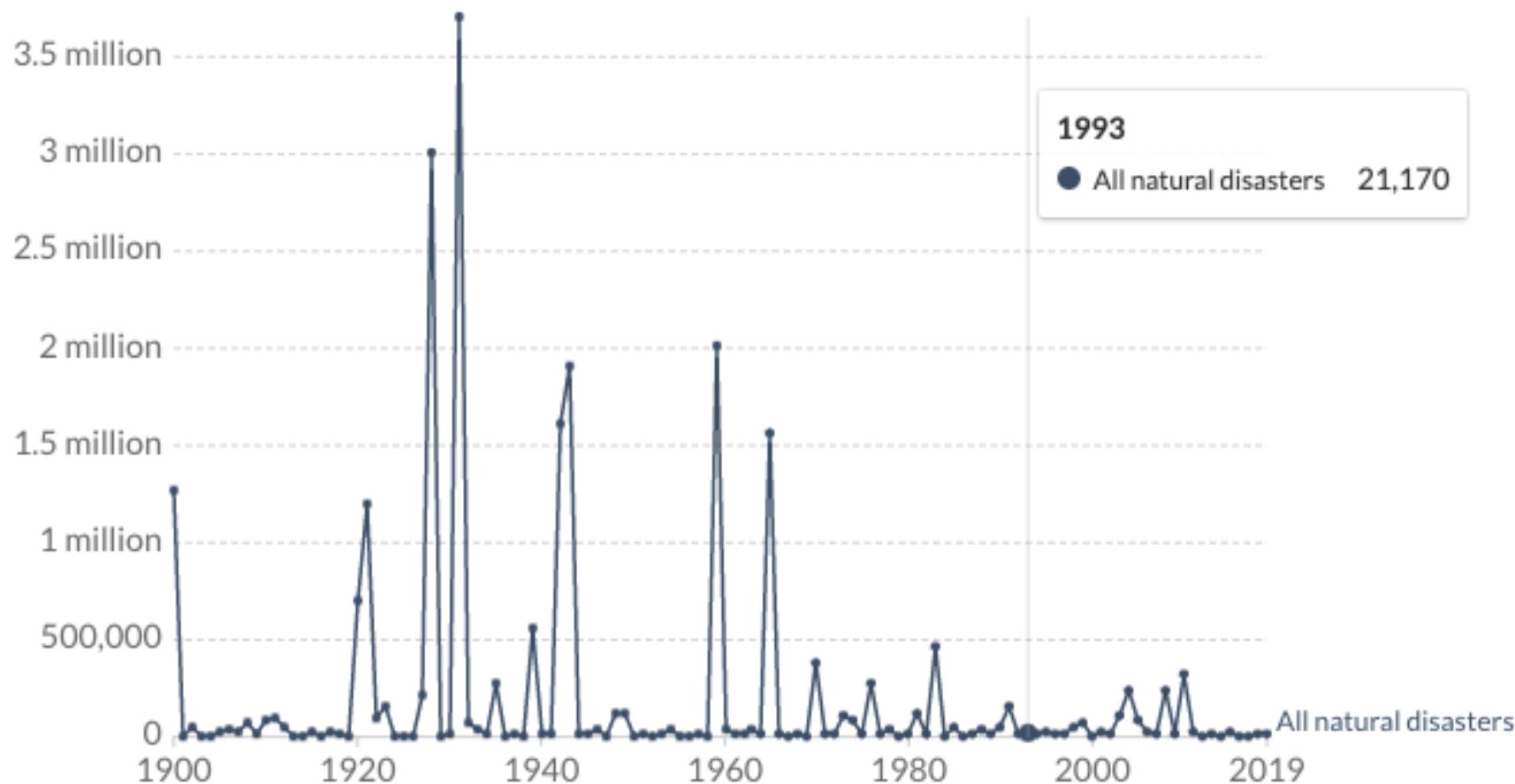
OurWorldInData.org/natural-disasters • CC BY

Global deaths from natural disasters, 1900 to 2019

Our World
in Data

Absolute number of global deaths per year as a result of natural disasters. "All natural disasters" includes those from drought, floods, extreme weather, extreme temperature, landslides, dry mass movements, wildfires, volcanic activity and earthquakes.

+ Add disaster category



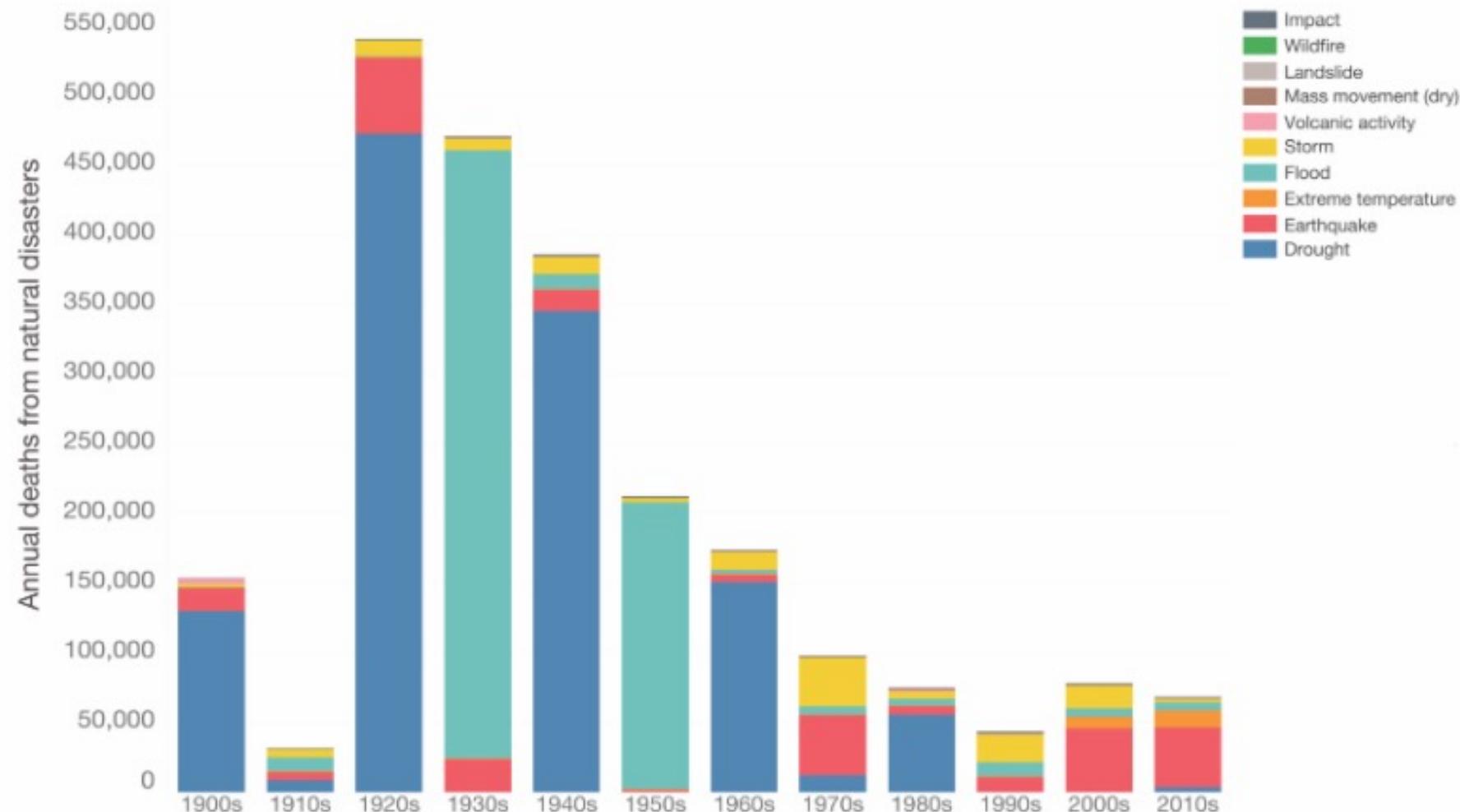
Source: EMDAT: OFDA/CRED International Disaster Database, Université catholique de Louvain – Brussels – Belgium

CC BY

Global annual deaths from natural disasters, by decade

Absolute number of global deaths from natural disasters, per year.

This is given as the annual average per decade (by decade 1900s to 2000s; and then six years from 2010-2015).

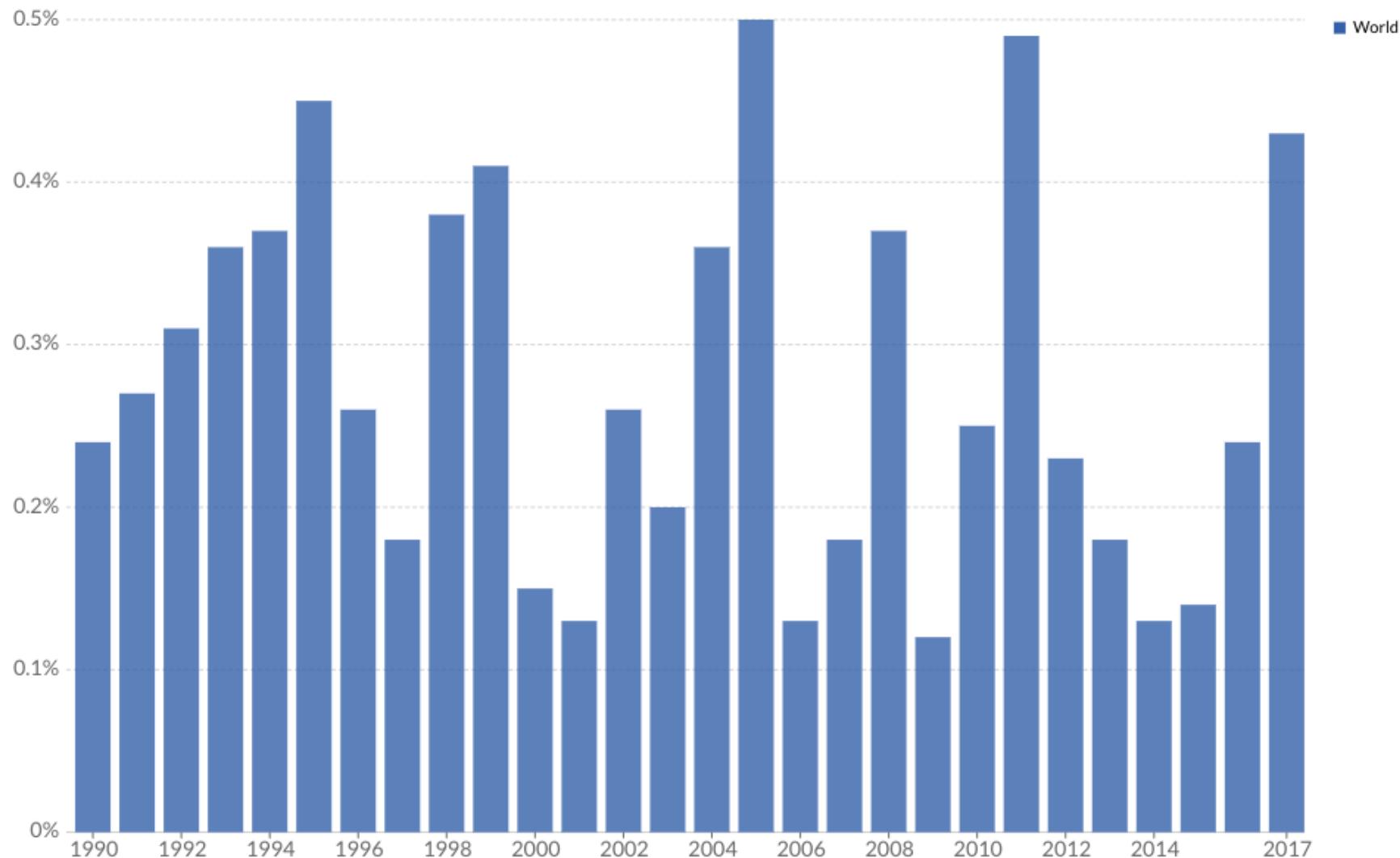


Source: EMDAT (2017); OFDA/CRED International Disaster Database, Université catholique de Louvain – Brussels – Belgium.
The data visualization is available at OurWorldInData.org. There you find research and more visualizations on this topic.

Licensed under CC-BY-SA by the authors Hannah Ritchie and Max Roser.

Global disaster losses as a share of GDP, 1990 to 2017

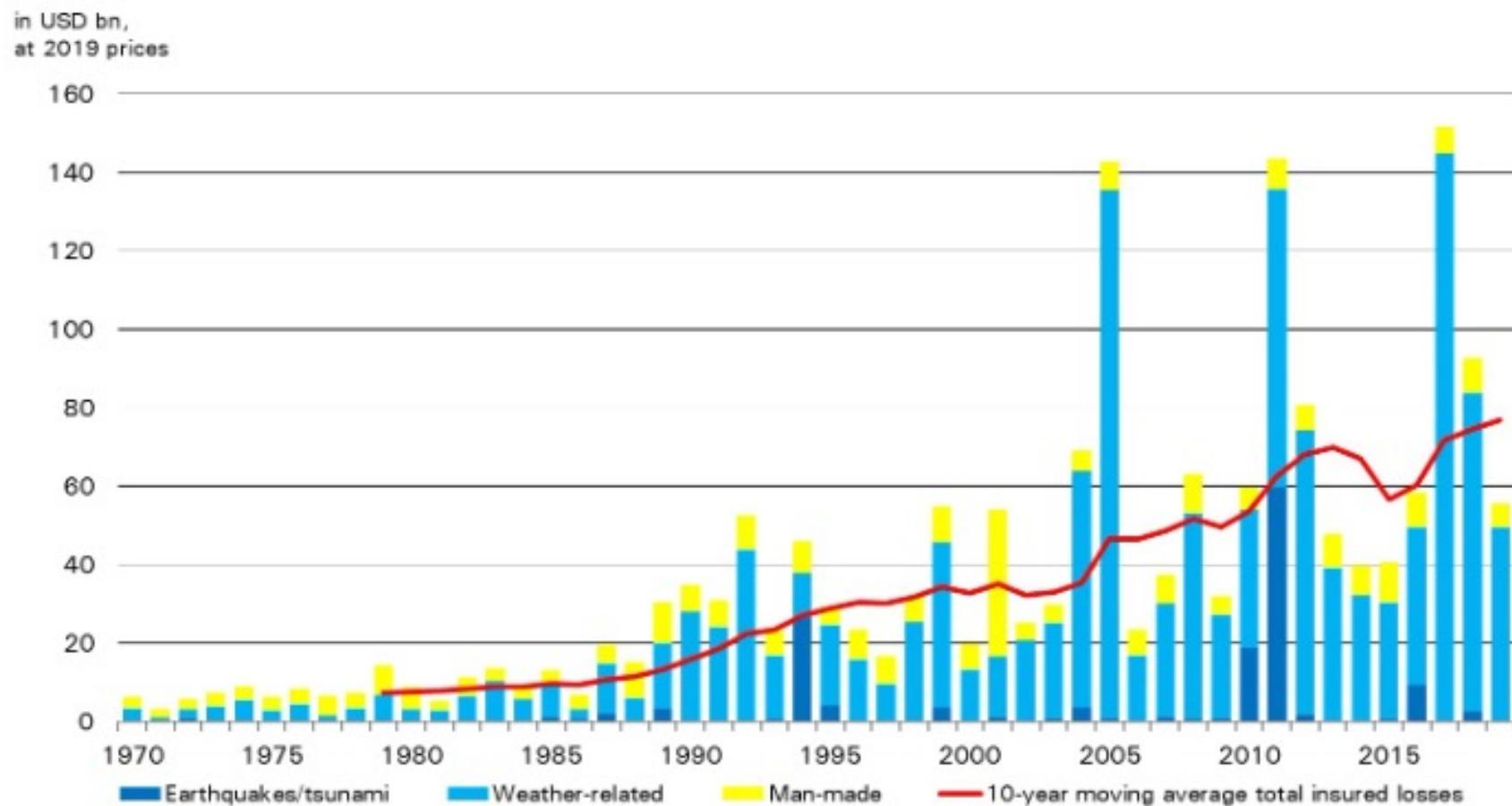
Global disaster losses (weather- and non-weather related) in economic terms, expressed as a share of global gross domestic product (GDP). Economic loss data from disasters is based on figures reported by Munich Re.



Source: Pielke, R. (2018). Tracking progress on the economic costs of disasters under the indicators of the sustainable development goals. Environmental Hazards, 1-6.

CC BY

Catastrophe-related insured losses



Insured losses continue to increase over time.

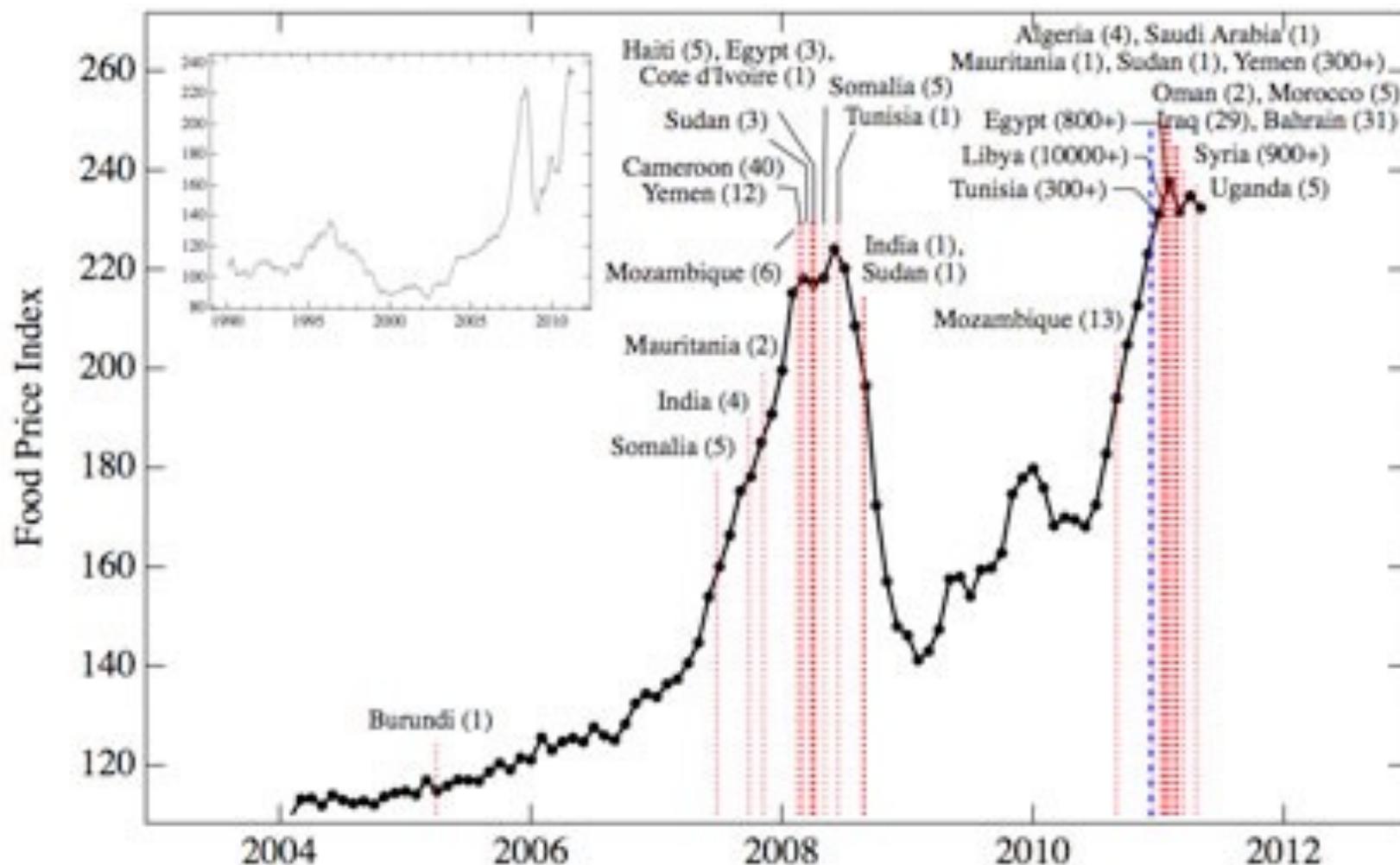
Poll

- Which of the following represent the most likely causal links?
 - a) instability > food prices > weather
 - b) instability > weather > food prices
 - c) food prices > instability > weather
 - d) food prices > weather > instability
 - e) weather > food prices > instability
 - f) weather > instability > food prices

Global Trends

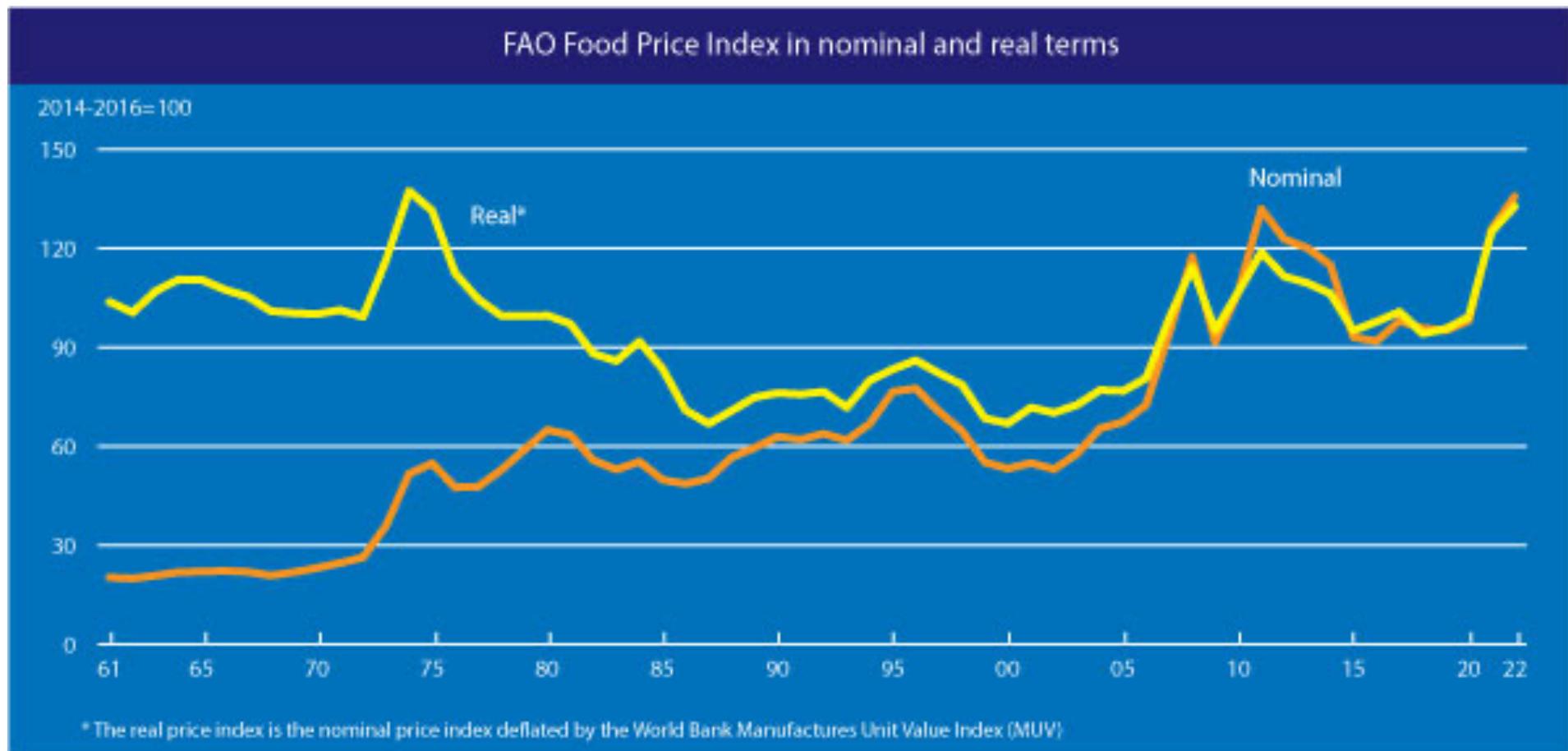
- By 2050, world population will be 9 billion with 70% living in cities
- Security of food, water and energy
- Almost a billion people experience hunger, with food price spikes and volatility threatening the sustainability of global food security
- Interconnected risks:
 - Rising food prices → Tunisian street vendor, Mohamed Bouazizi, sets himself on fire and dies → Ben Ali toppled from power → Arab Spring

Food: prices and riots



The UN FAO food price index correlates with "food riots" around 2008 and the "Arab spring" conflicts. Death tolls are reported in parentheses.
Source: New England Complex Systems Institute

Food prices



Syria and climate change

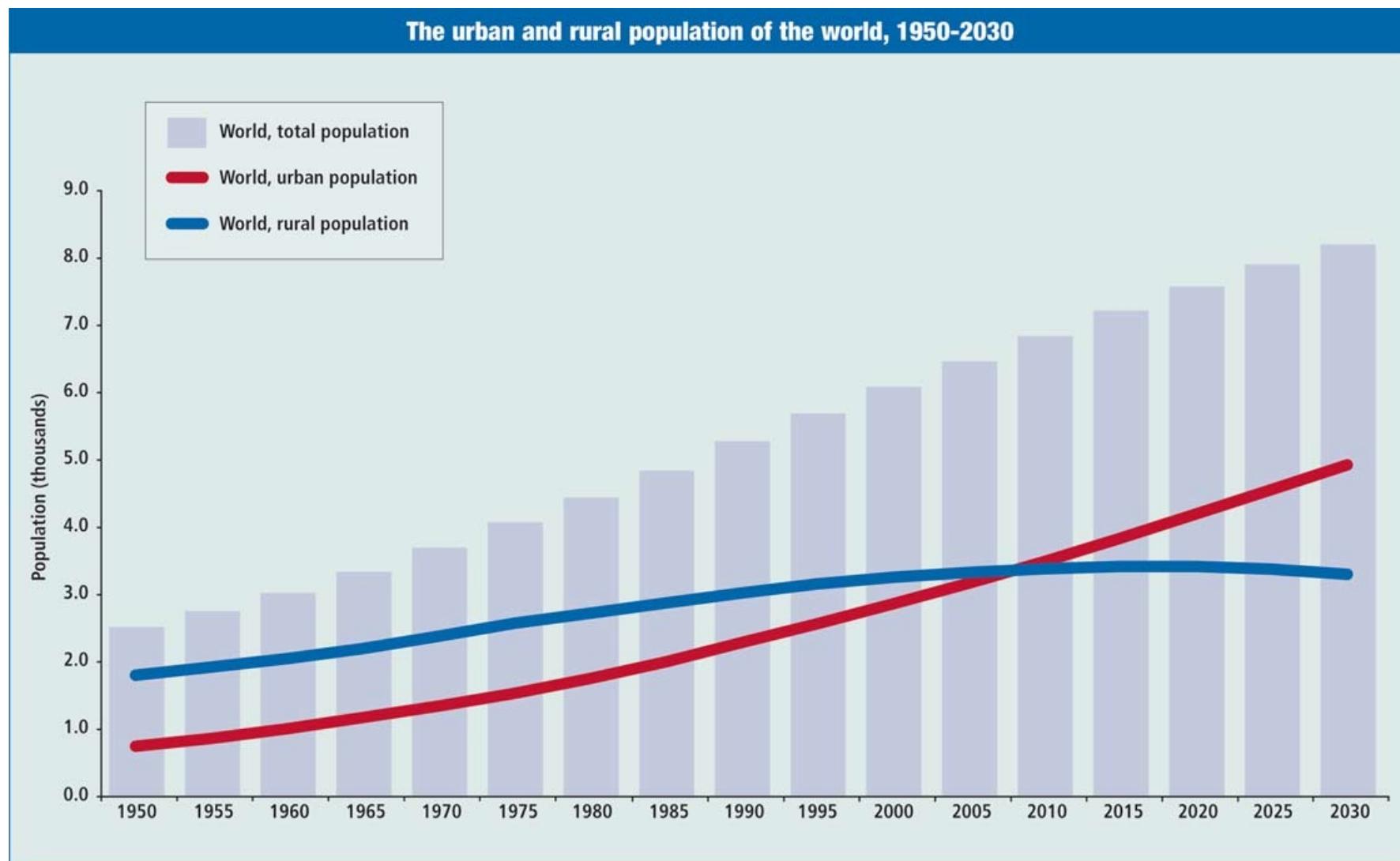
- 2006 – 2011: over half of the country suffered under the worst drought on record
- Drought was more intense and lasted longer than can be explained by natural variability
- Nearly 85% of livestock died.
- Since the uprising against Assad in March 2011, over 240,000 people have been killed, 4 million Syrians have fled their country, and over 7 million have been displaced.
- Syria is on path to lose nearly 50% more of its agricultural capacity by 2050.

Climate change could lead to 1m migrants a year entering EU by 2100



Researchers plotted temperature rises against the number of asylum applications and are predicting that as the southern hemisphere heats up the number of people migrating to the EU each year will triple.

Urban and Rural Population



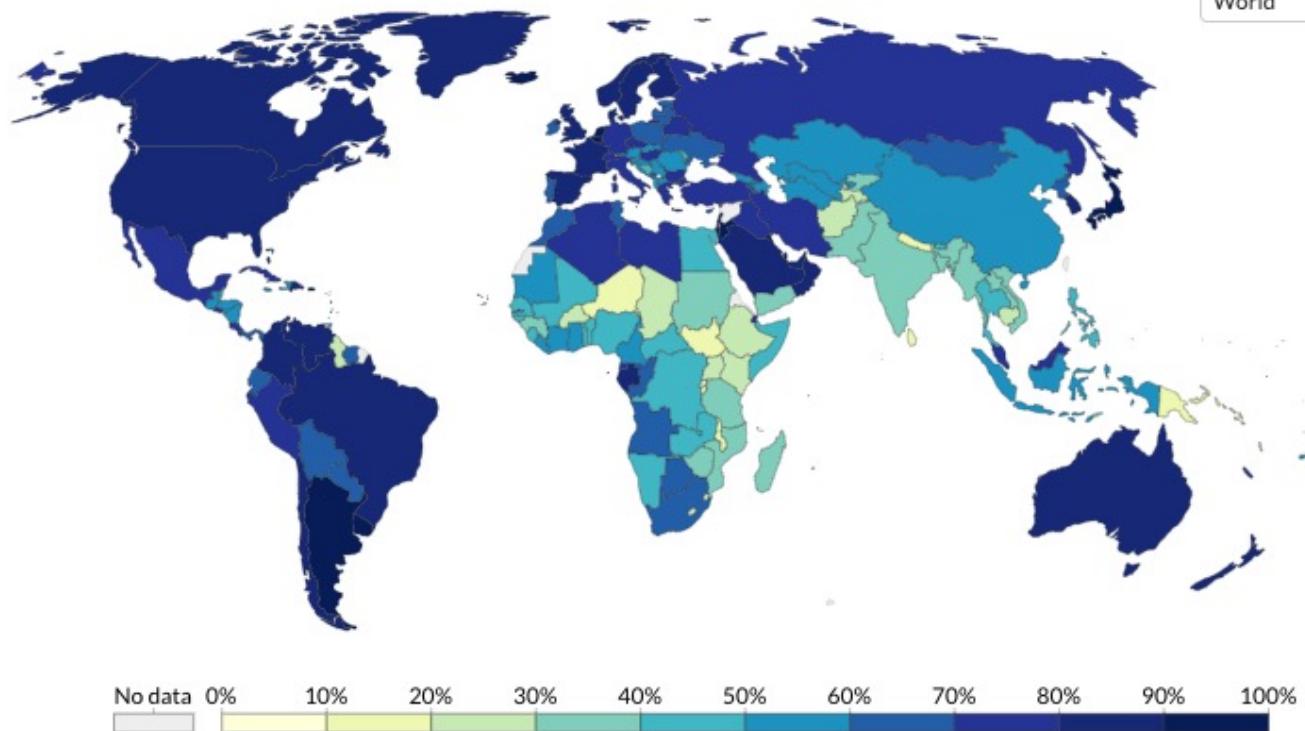
Source: UN

World Urban Population

Share of people living in urban areas, 2017

Our World
in Data

World



Source: UN World Urbanization Prospects (2018)

Note: Urban populations are defined based on the definition of urban areas by national statistical offices.

CC BY

► 1960

2017

Risk communication

- There is a need to standardised measurements of risk across society and to ensure that the units are meaningful and easy to communicate and comprehend.
- In this way, it will be possible to compare risks relative to a known and understood benchmark such as cycling to work.
- A micromort is a unit of risk measuring a one-in-a-million probability of death and the term comes from combining the words micro and mortality.
- Micromorts provide a useful means of measuring the risk arising from various day-to-day activities.

Size of acute risks (immediate death)

Activity	Cause of death
Drinking 0.5 liter of wine	cirrhosis of the liver
Smoking 1.4 cigarettes	cancer, heart disease
Spending 1 hour in a coal mine	black lung disease
Spending 3 hours in a coal mine	accident
Living 2 days in New York or Boston in 1979	air pollution
Living 2 months with a smoker	cancer, heart disease
Drinking Miami water for 1 year	cancer from chloroform
Eating 100 charcoal-broiled steaks	cancer from benzopyrene
Eating 40 tablespoons of peanut butter	liver cancer from aflatoxin B
Travelling 6000 miles (10,000 km) by jet	cancer due to increased background radiation

Data Analytics

WEEK 4B

Course outline

Week	Lecture A	Lecture B
1	Data Analytics	Weather forecasting
2	Renewable energy	Wind energy
3	Solar energy	Demand forecasting
4	Risk	Extreme events
5	Health	Biomedicine
6	Early warning systems	Economic forecasting

Today's Lecture

No.	Activity	Description	Time
1	Challenge	Floods	10
2	Discussion	Extreme values	10
3	Case study	Climate change and flood risk	10
4	Analysis	Wind farm risk metric	20
5	Demo	Estimating V50	20
6	Q&A	Questions and feedback	10

Poll: extreme events

- Give examples of extreme events:
- **Slido.com #98798**

Extreme events

- Heatwaves
- Droughts
- Crop failure
- Wildfires
- Floods
- Hurricanes
- Tropical cyclones

Extreme weather

- Extreme weather threatens African society and economy
- Droughts, floods and storms are hitting Africa with increasing frequency.
- With so many people dependent on subsistence agriculture, the results can be devastating, and the future looks uncertain.

<https://www.dw.com/en/extreme-weather-threatens-african-society-and-economy/a-45713490>

Floods



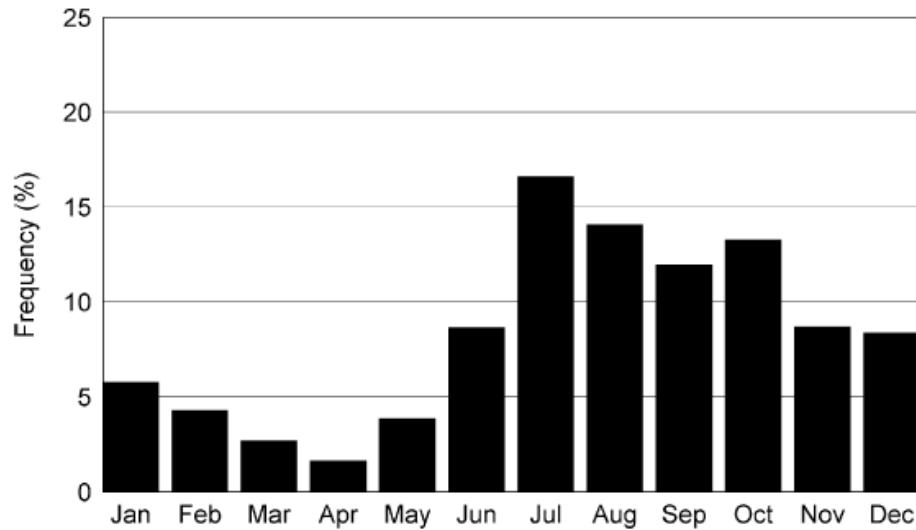
Source: www.abc.net.au; Thailand, 2011

Flash Floods

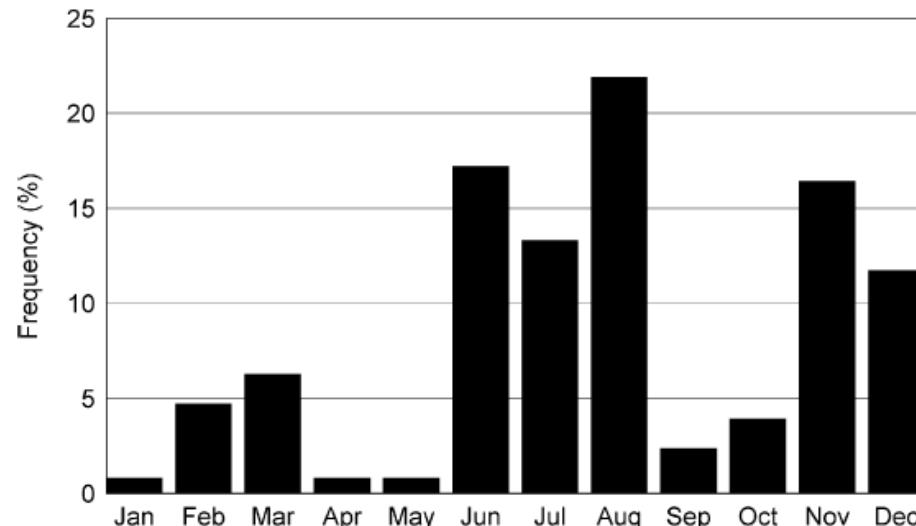


Source: Met Office, Boscastle, 2004

British rainfall extremes



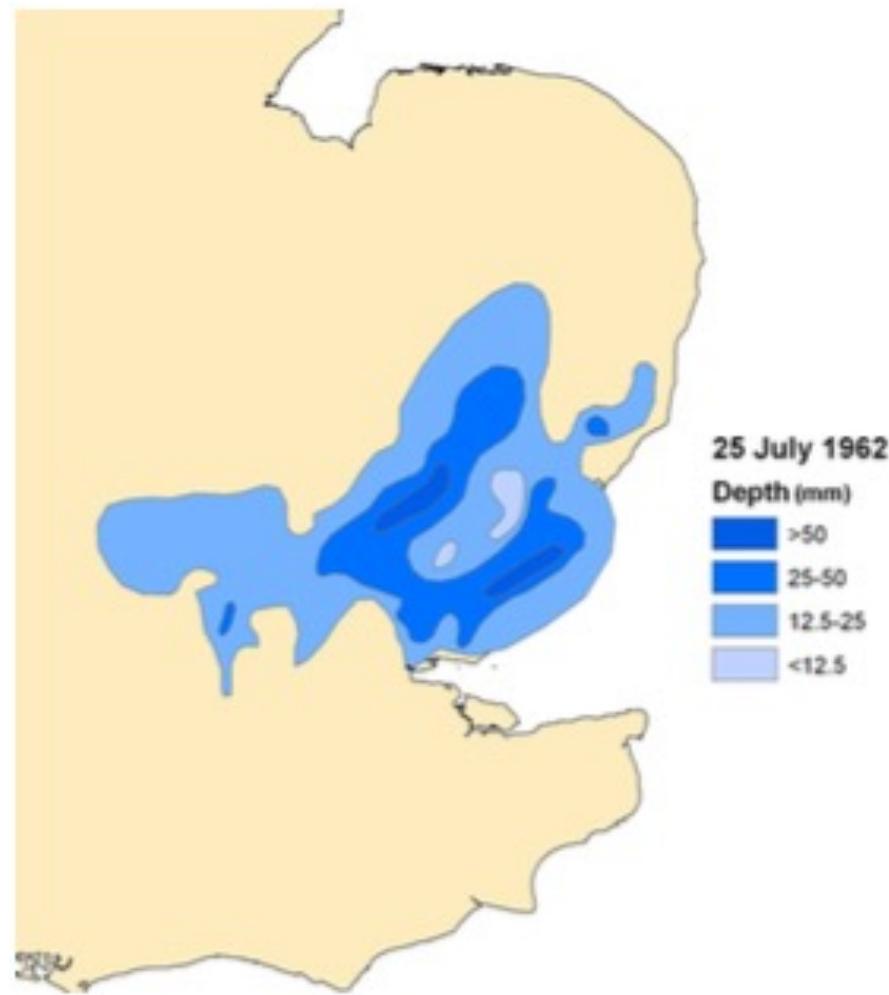
All recorded extremes



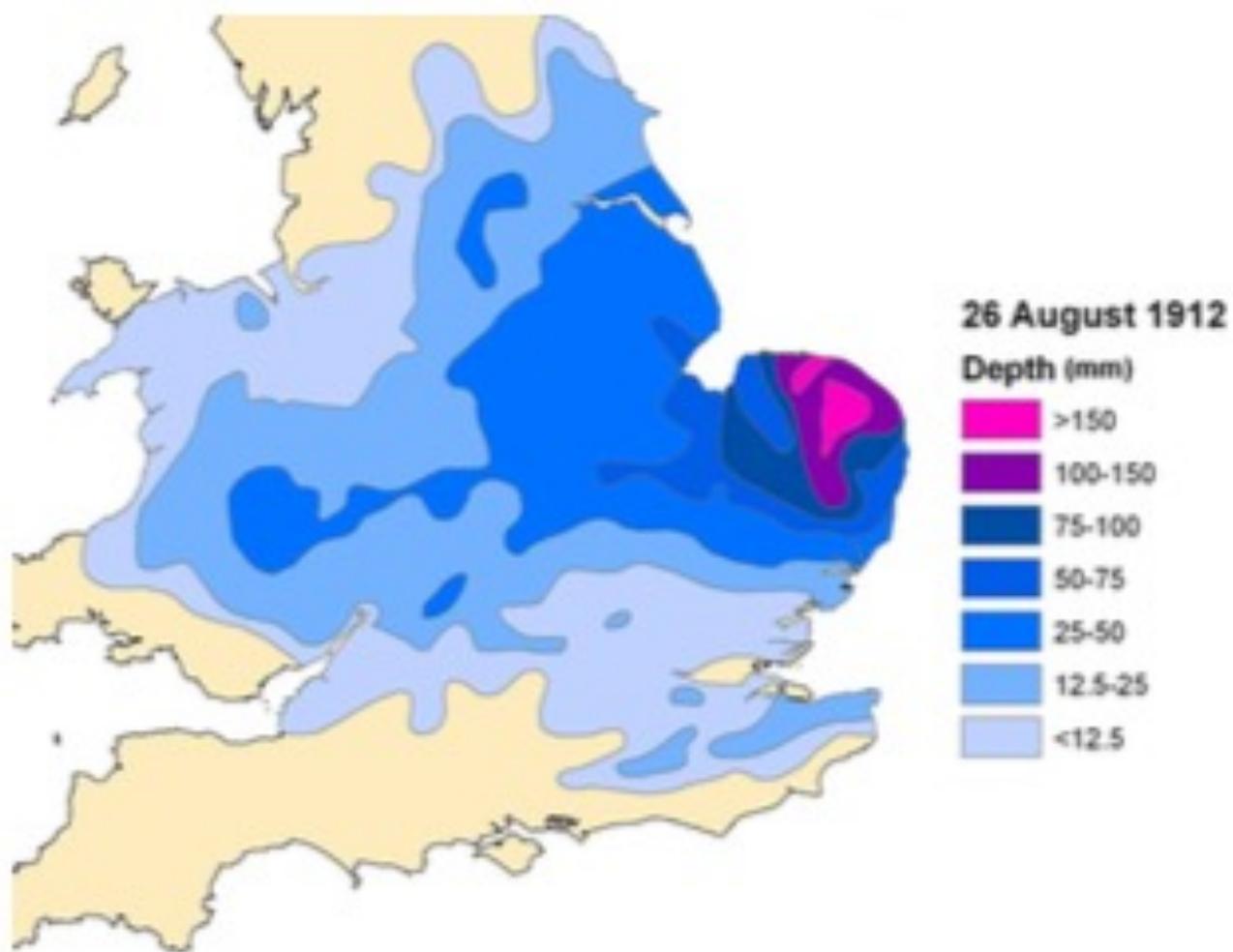
Daily observations of 150mm and above

- British Rainfall (Symons, G. J., 1885) lists all observed 24 hour rainfall depths which exceeded 2.5 inches (63.5mm).

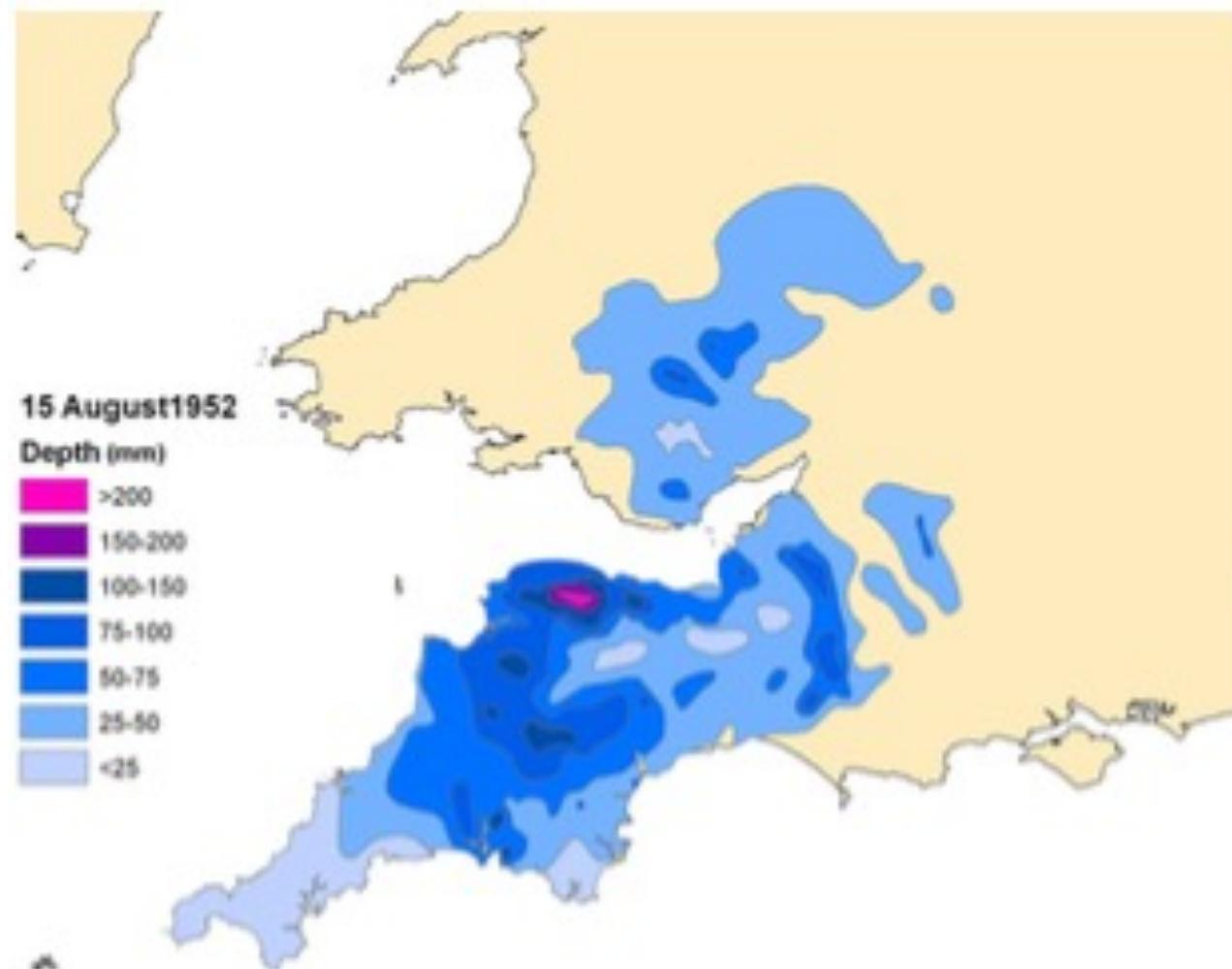
Thunderstorm



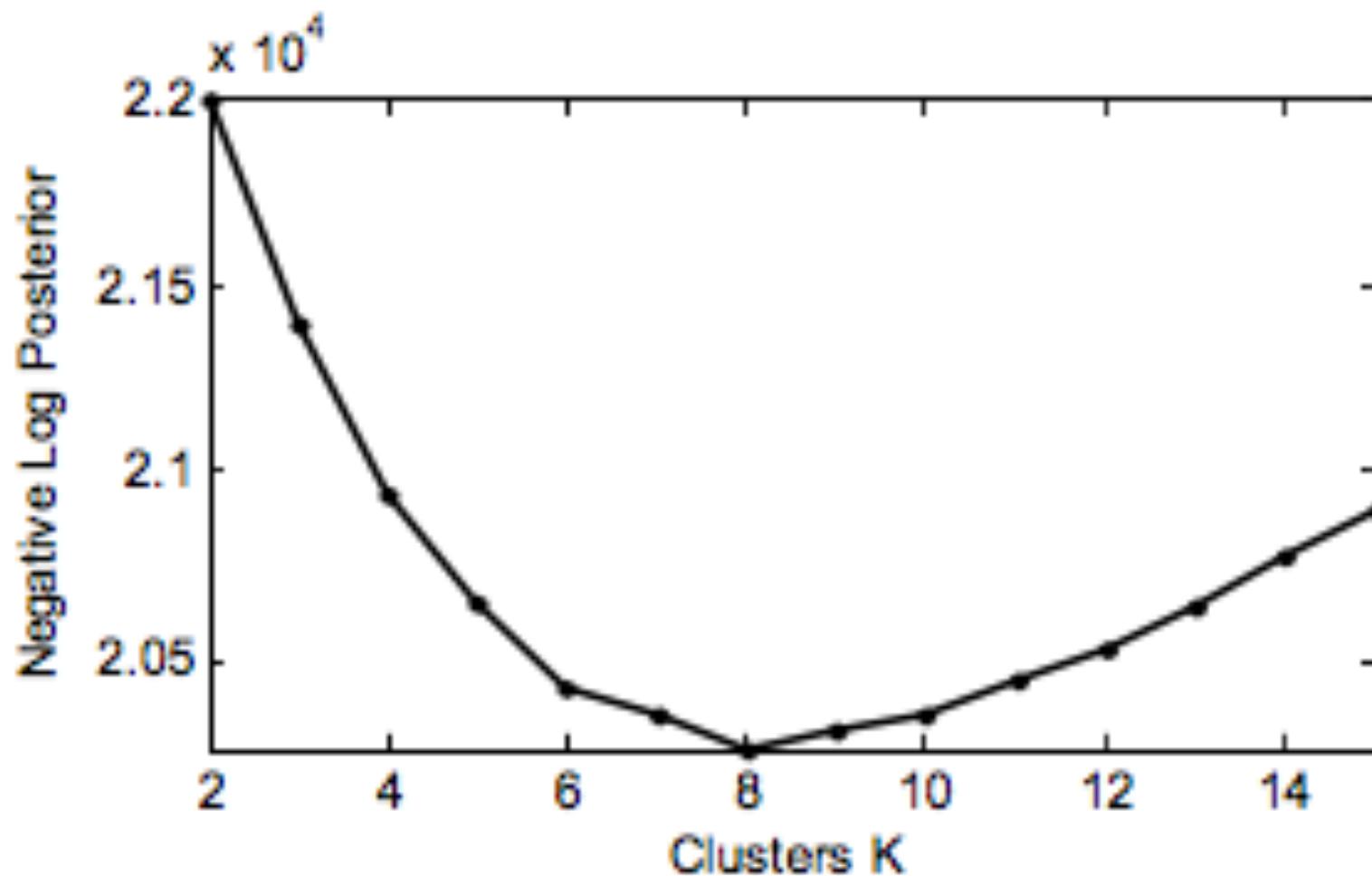
East coast extreme event



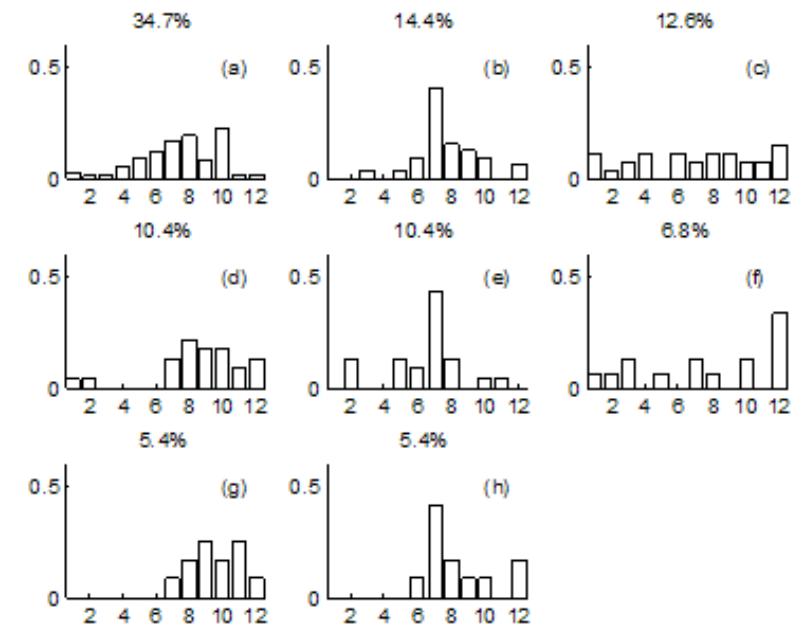
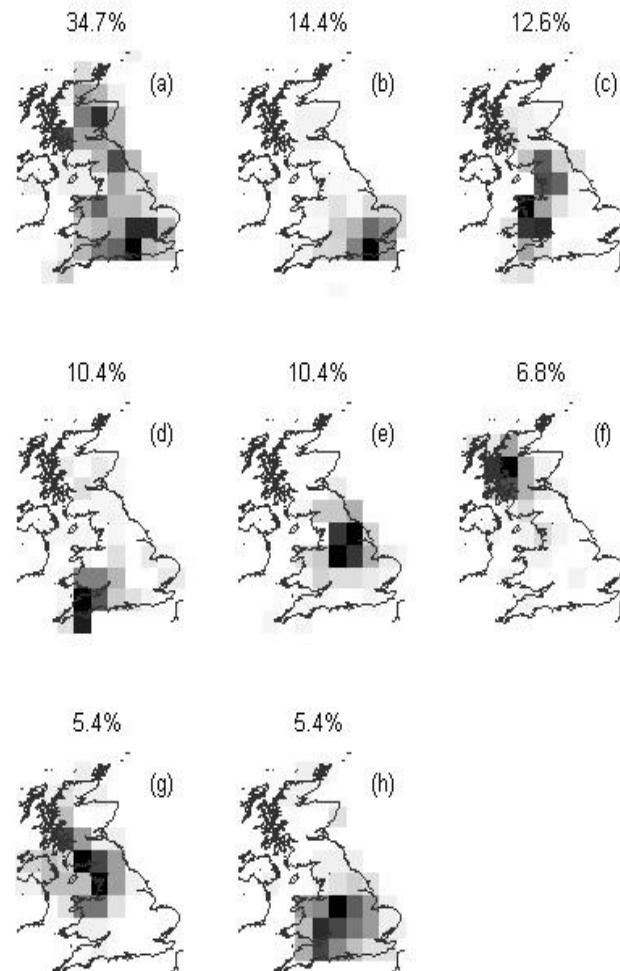
Mesoscale convective complex (MCC)



How many extreme rainfall patterns?



Bayesian objective classification



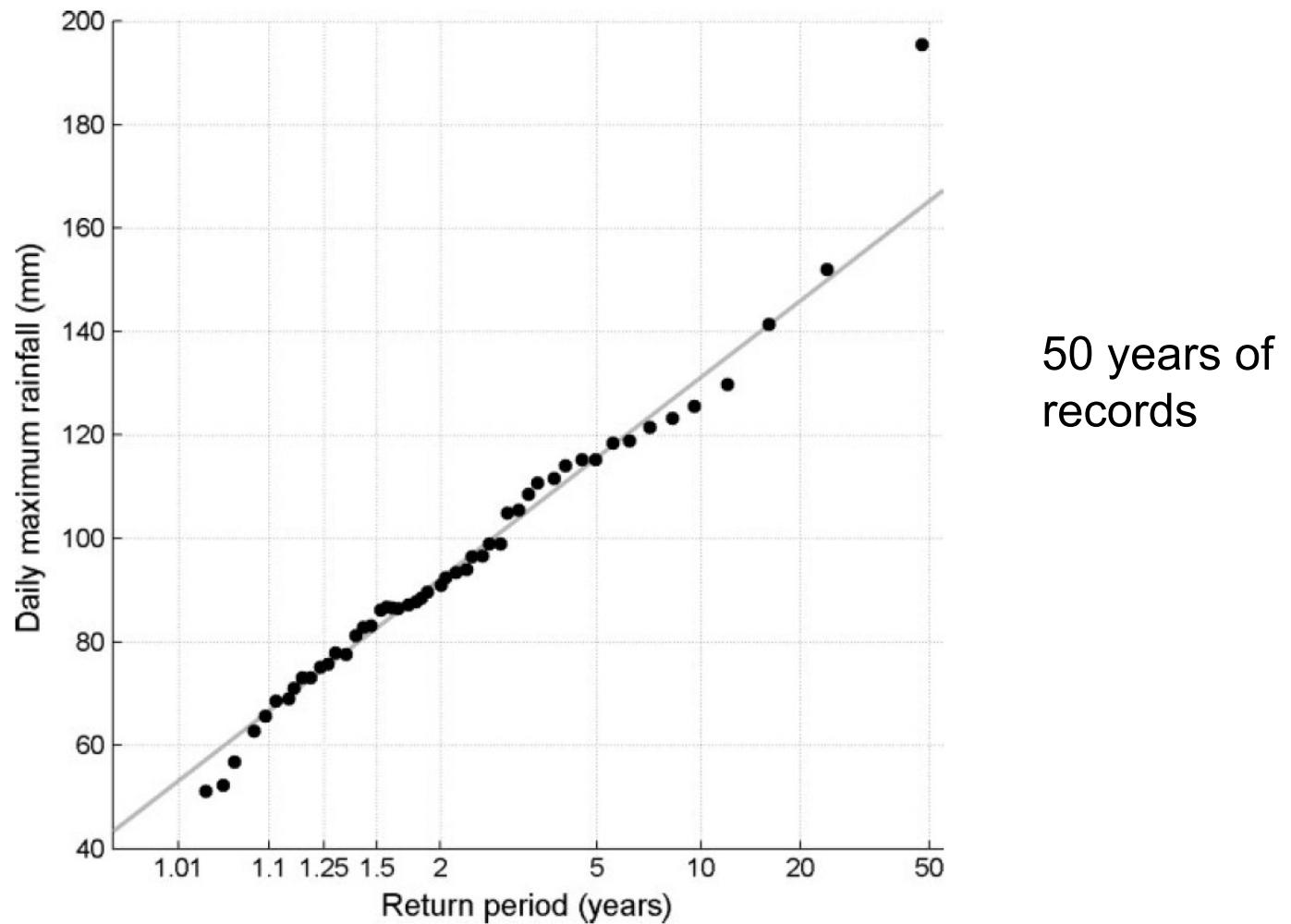
Poll: Identifying extreme events

- From a time series, how would you identify extreme events?
- **Slido.com #98798**

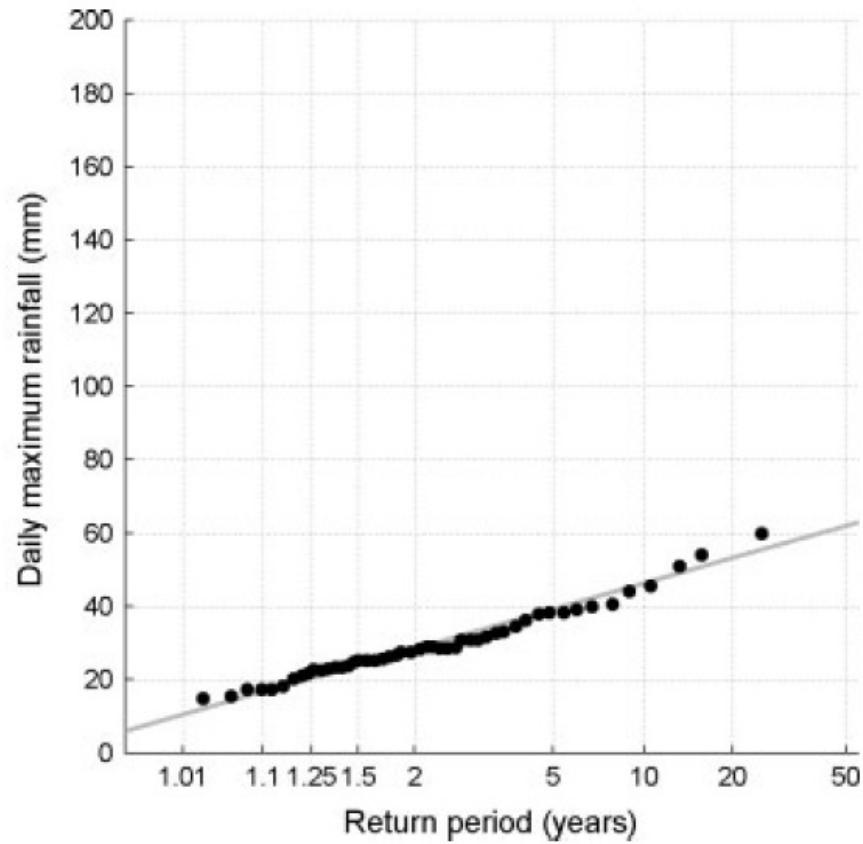
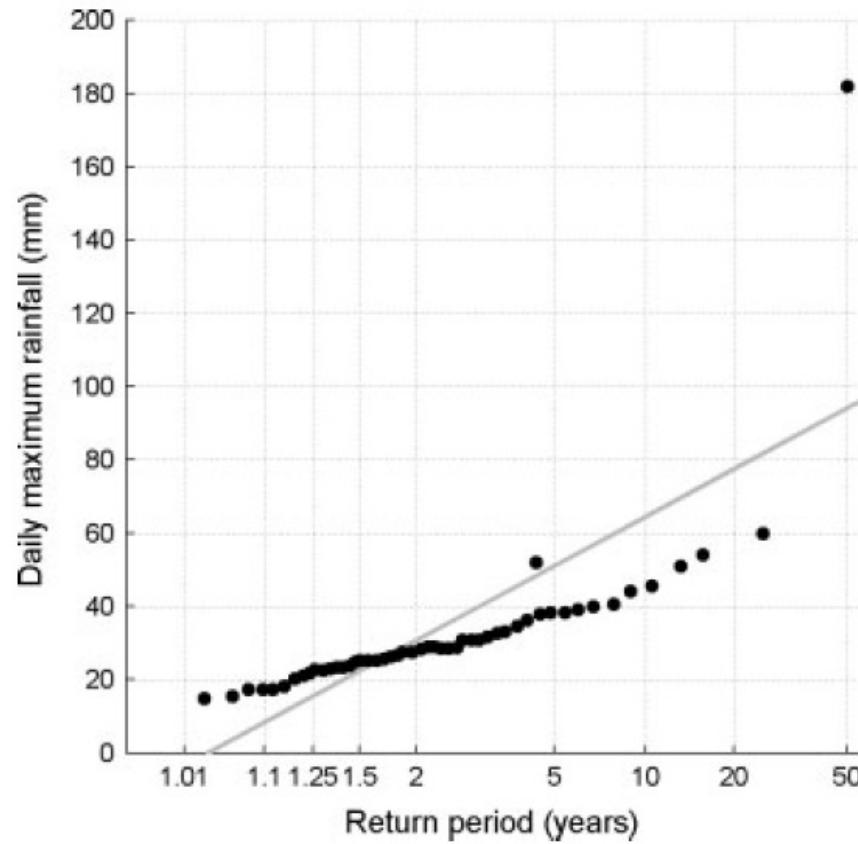
Extreme value theory

- 1. Block maxima, e.g. annual maxima series
 - Assumes only one relevant data point per block and that an appropriate time-scale exists for separating to obtain independent events
- 2. Point over threshold (POT)
 - Relies on selecting a threshold and fitting two distributions: one for the frequency of events and another for the size of exceedances

Gumbel extreme value analysis for Blaenau Ffestiniog



Sensitivity to extremes



Gumbel analyses for Sprowston with and without the highest maximum fall

Risk management

- Selecting an appropriate probability level for defining extreme events
- Compromise between accuracy of statistical estimates and application
- Regulator (e.g. Solvency II for insurers): 1 in 200 year event
- Flood: 1 in 100 year event
- Wind turbines: 1 in 50 year event

Climatology

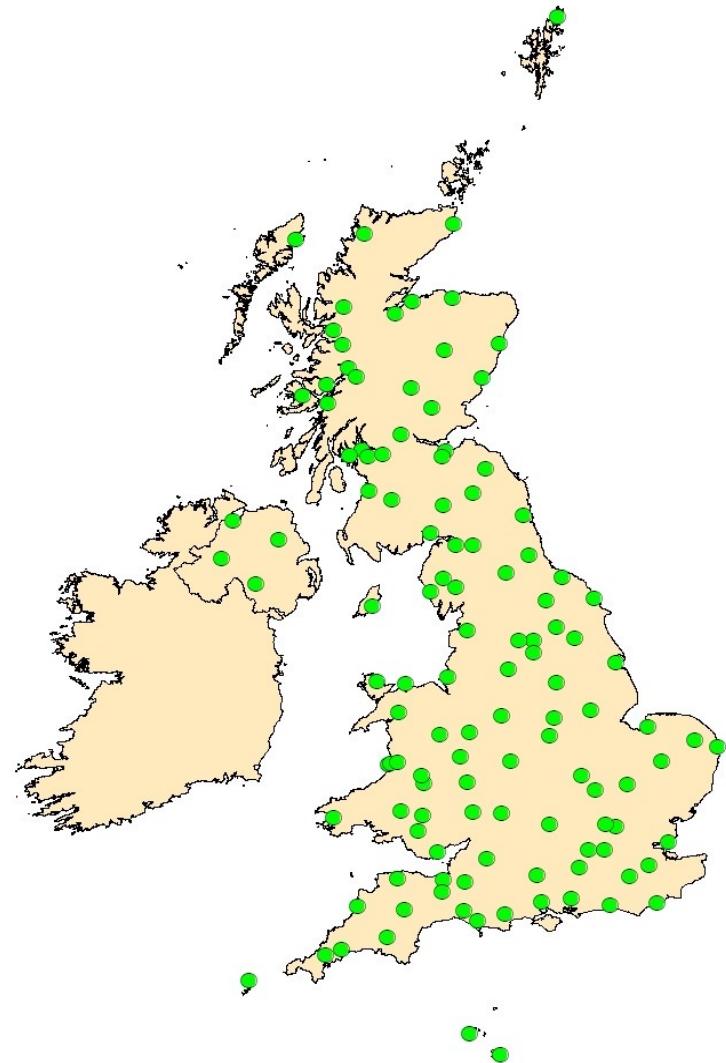
- Study of the magnitude, frequency and distribution of intense rainfall in the UK
- Comparison of the same 120 daily read rain gauges using the same methodology over two periods:
 - Period 1: 1911 to 1960
 - Period 2: 1961 to 2006
- Increases up to 20% have occurred in the north-west of the country and in parts of East Anglia
- There have also been changes in other areas, including decreases of the same magnitude over central England.

Extreme rainfalls 1911-1960

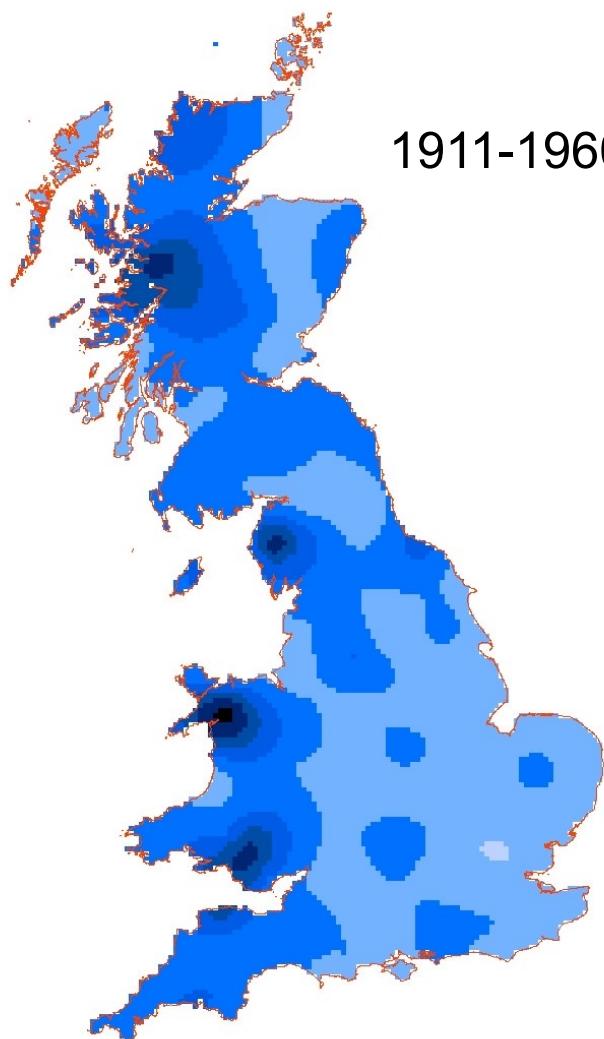
1911-1960 – using results as published in 1966
(Rodda)

1961-2006 – using data from the Met Office
MIDAS dataset

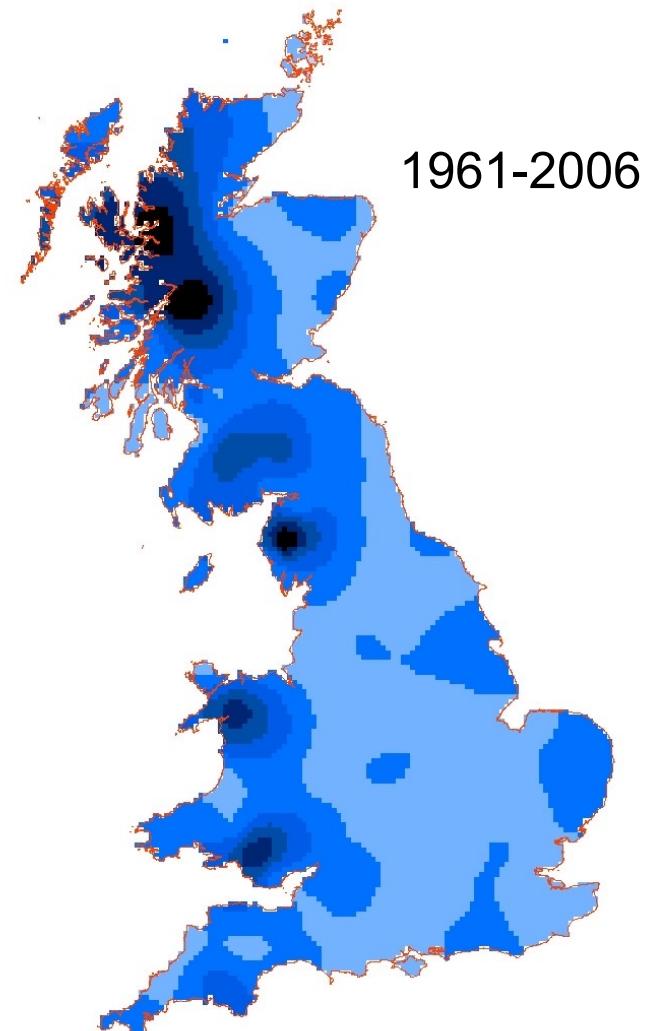
- Used annual maxima from 120 rain-gauges across the British Isles
- Return period 24 hour rainfall calculated using the EV1 distribution
- Point interpolation using GIS



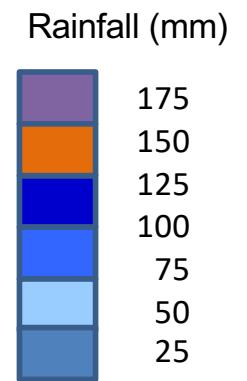
24-hour rainfall 100 year return period



1911-1960

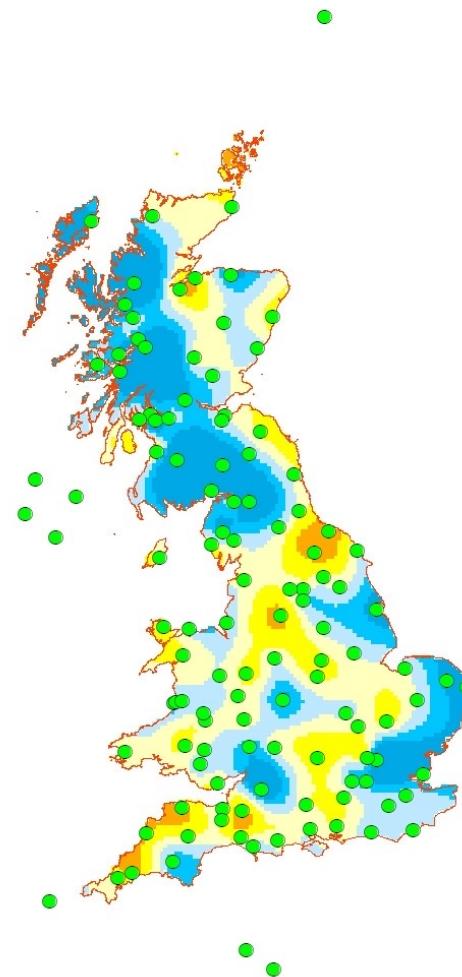
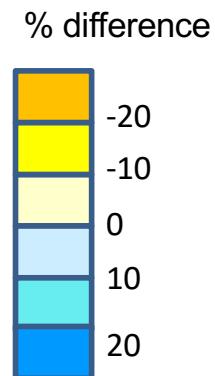


1961-2006



% difference between 1911-1960 and 1961-2006 return period magnitudes

- Increase in NW
- Increase in SE
- Decrease in SW
- Smaller changes in other areas



Description of UK Flood 2015

- More than half the 5,000 properties flooded over the weekend were still under water last night, with families facing the prospect of eight months away from their homes.
- Farmers ignored police health and safety warnings as they cleared a river of debris to stop its banks bursting and to prevent further damage.
- The body of a man swept away by a swollen river was recovered by police.
- About 40 schools remained closed and hospitals cancelled routine appointments.
- Households in flood-hit areas face further rises in insurance bills, while industry experts warned that clean-up costs could reach £100 million.
- Storm Desmond set a new record for rainfall in a 48-hour period, with 15.9in in 38 hours at Thirlmere, Cumbria. The previous record was 15.6in, in November 2009. Water levels in Carlisle reached about 20in above the level seen in 2005, which was itself 20in above the previous record in 1853.

UK Floods (2015) - reaction

- Ministers announced a full review of Britain's storm defences yesterday as they blamed global warming for the floods that have devastated Cumbria.
- Liz Truss, the environment secretary, told MPs that Storm Desmond was consistent with climate change and the government's predictions for future flooding would now have to be updated.
- Last night, however, an Oxford academic said that the government had been warned five years ago of the risk of higher rainfall in the area possibly created by climate change.
- There was also criticism of the government for cutting spending on flood defences by 14 per cent in the past year, with Labour accusing David Cameron of breaking promises to protect homes.

The warning

- Some scientists are worried that the agency may have underestimated the dangers posed by global warming and a cluster of record-breaking weather. Yet the government was warned of a sharp rise in the risk of flooding in the northwest five years ago.
- Oxford University researchers and analysts from a private company found in 2010 that the amount of extreme rainfall in the region had risen by at least 20 per cent over five decades. Patrick McSharry, the head of forecasting and risk analysis at Oxford's Smith School of Enterprise and the Environment, said the Environment Agency's calculations of the odds might well be out of date.
- “I would imagine that many at the Environment Agency are concerned about climate change and how the usual calculation of return periods may be flawed,” he said. “This paper was written with a meteorologist and an environmental consultant so I am sure that the paper would have been seen by members of the Environment Agency.”

Weather extremes

- In order to study weather extremes, which of the following offers the most promising resource?
 - a) Five years of weather station observations
 - b) Satellite imagery for 20 years
 - c) Reanalysis weather data for 50 years
 - d) Historical archives
- **Slido.com #98798**

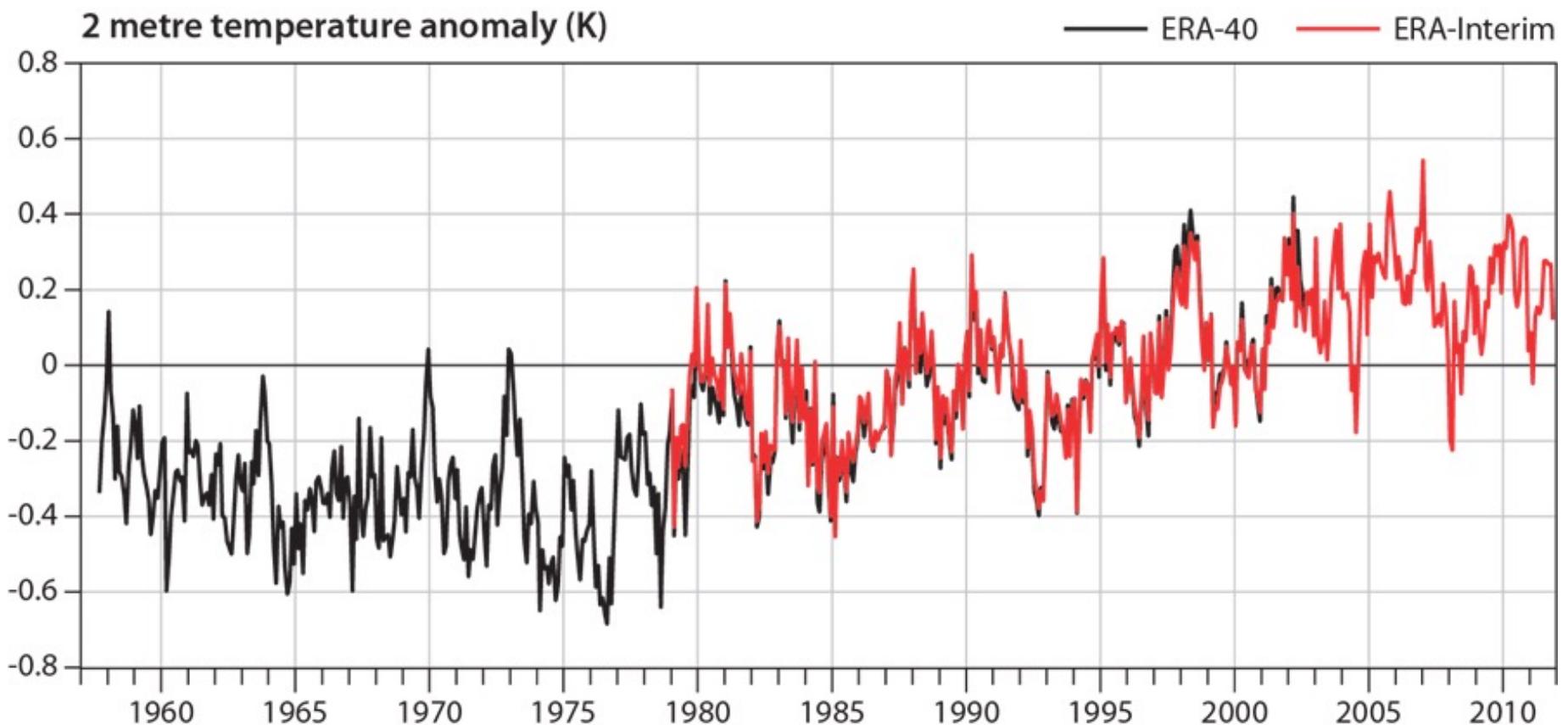
Reanalysis Data

- ECMWF periodically uses its forecast models and data assimilation systems to 'reanalyse' archived observations, creating global data sets describing the recent history of the atmosphere, land surface, and oceans.
- Reanalysis data are used for monitoring climate change, for research and education, and for commercial applications.

What is a climate reanalysis?

- A climate reanalysis gives a numerical description of the recent climate, produced by combining models with observations.
- It contains estimates of atmospheric parameters such as air temperature, pressure and wind at different altitudes, and surface parameters such as rainfall, soil moisture content, and sea-surface temperature.
- The estimates are produced for all locations on earth, and they span a long time period that can extend back by decades or more.
- Climate reanalyses generate large datasets that can take up many terabytes of space.

Temperature reanalysis



ERA-40: 1957 – 2002; ERA-Interim: 1979 – present

Wind turbine fires

- What is the approximate probability of a wind turbine going on fire in a given year?
 - a) 1 in 20
 - b) 1 in 200
 - c) 1 in 2000
 - d) 1 in 20000
- Slido.com #98798

Wind turbine fire estimation

- The Caithness Windfarm Information Forum (CWIF) identified 200 reported fires between 1995 and 2012, an average of 11.7 fires per year.
- This makes fires the second most common cause of reported accidents in wind turbines.
- According to the study, the three most important causes of accidents in wind turbines are:
 - Blade failure (19%)
 - Fire (15%)
 - Structural failure (9.7%)
- 91% of all fires are not reported!
- In 2011 an estimated 200,000 wind turbines were operated worldwide.
- Based on the data in the IAFSS report, we can assume that there were 117 fires (both reported and unreported) in the same year. This means that in 2011, 1 out of 1,710 turbines caught fire.



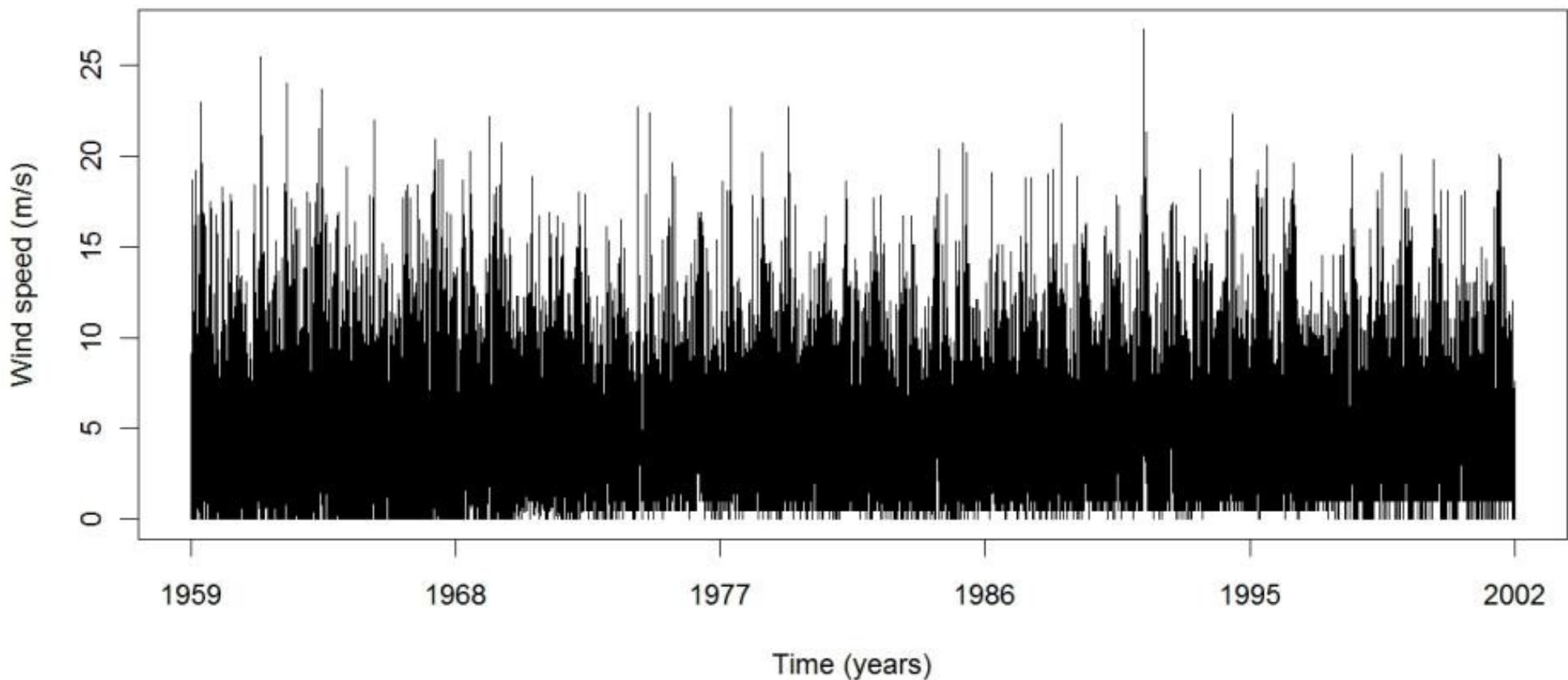
Wind extremes

- The 50-year return level, V_{50} , is defined as the wind speed value that is exceeded, on average, once in a year with a probability of 1/50
- It is usually calculated from observations using a long record of historical data
- Most potential wind farms sites have datasets of a few years
- Can we use reanalysis data to get trustworthy estimates of V_{50} , when we have only a few years of real data?

Anastasiades & McSharry (2014). [Extreme value analysis for estimating 50 year return wind speeds from reanalysis data](#). Wind Energy.

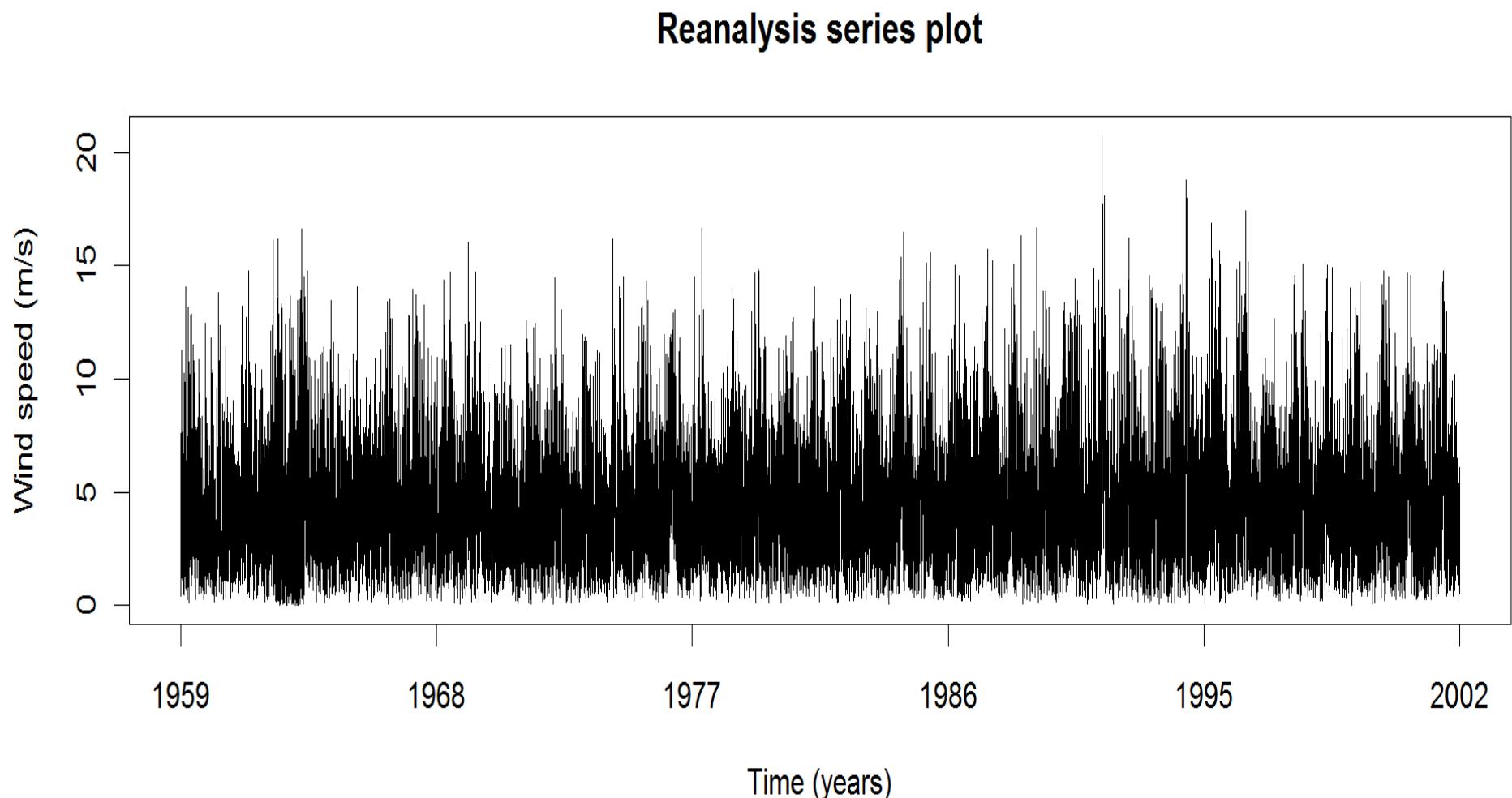
Schiphol wind data

Actual series plot



Hourly data for 45 years (01-Sep-1957 to 31-Aug-2002)

ECMWF Reanalysis data



Data sampled every six hours for 45 years (01-Sep-1957 to 31-Aug-2002)

Block maxima method

- Assume we have k years of wind speed data
- Divide the time series into n independent blocks of m days long, and fit a Generalized Extreme Value (GEV) distribution to the block maxima
- Evaluate the 50-year return level, $V_{50}(k)$, by using extrapolation: Calculate the $1 - \frac{m}{50 \cdot 365}$ quantile of the fitted GEV

Two challenges

- The 50-year return value from the reanalysis data, $V_{50}^R(k)$, is underestimating the corresponding value of actual data, $V_{50}^A(k)$, due to the difference in temporal resolution of the two series
- When $k < 20$, the $V_{50}(k)$ estimates are not reliable

Solution- Calibration method

- Solution: $V_{50}^R(k)$ is underestimated by a constant factor for all k . Hence
$$\frac{V_{50}^A(k)}{V_{50}^R(k)} \approx \text{constant}$$
- Therefore $V_{50} \approx \frac{V_{50}^A(45)}{V_{50}^R(45)} \cdot V_{50}^R(45) \approx \frac{V_{50}^A(k)}{V_{50}^R(k)} \cdot V_{50}^R(45)$
- But $V_{50}^A(k)$ and $V_{50}^R(k)$ are highly depended on the goodness of fit of the GEV distribution and hence this ratio is not robust
- Use a non-parametric method to approximate this ratio

Solution- Calibration method

- From two series of length k , we can use interpolation to find estimates of $V_{50}^A(k)$ and $V_{50}^R(k)$ by just calculating the $1 - \frac{1}{50 \cdot 365 \cdot l}$ quantiles for $l=24,4$ respectively.

- These estimates will be highly biased but their ratio should be equal to the optimal ratio:

$$\hat{R}(k) = \frac{V_{50}^A(k)}{V_{50}^R(k)} \approx \text{constant}$$

- This ratio will be more robust since it does not depend on any parametric method.

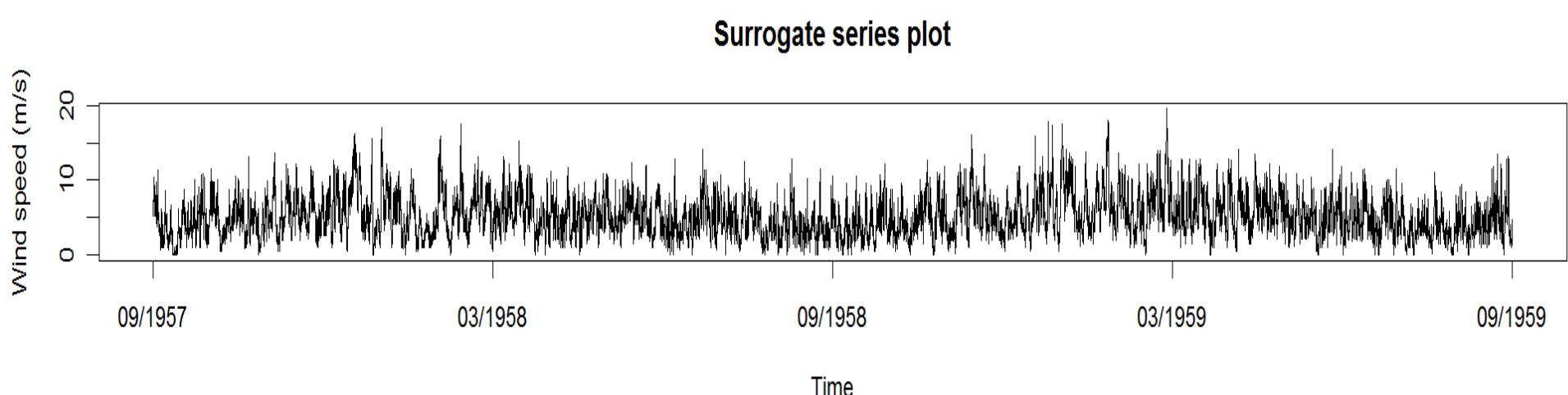
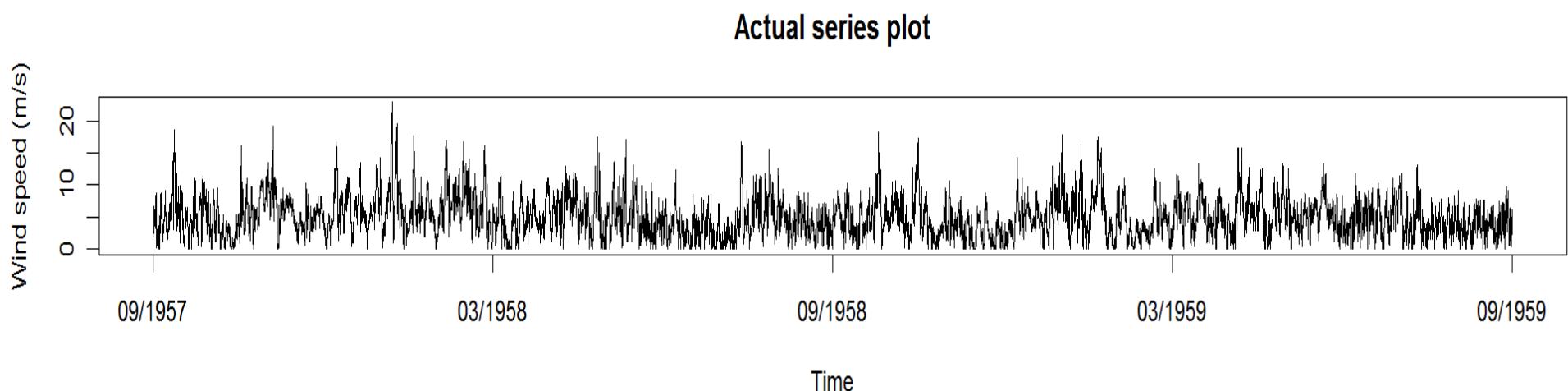
Bootstrap technique

- Create 1500 surrogates of the actual data, with 1500 corresponding reanalysis series
- The surrogates are created so they preserve the autocorrelation (seasonality of our actual data) and the unconditional distribution
- For each surrogate of the two series we randomly select from 1 to 45 years (in steps of four years) of data and calculate the corresponding $V_{50}^A(k)$, $V_{50}^R(k)$, $\hat{R}(k)$ and calibrated 50-year return level

$$V_{50}^{cal}(k) = \hat{R}(k) \cdot V_{50}^R(45)$$

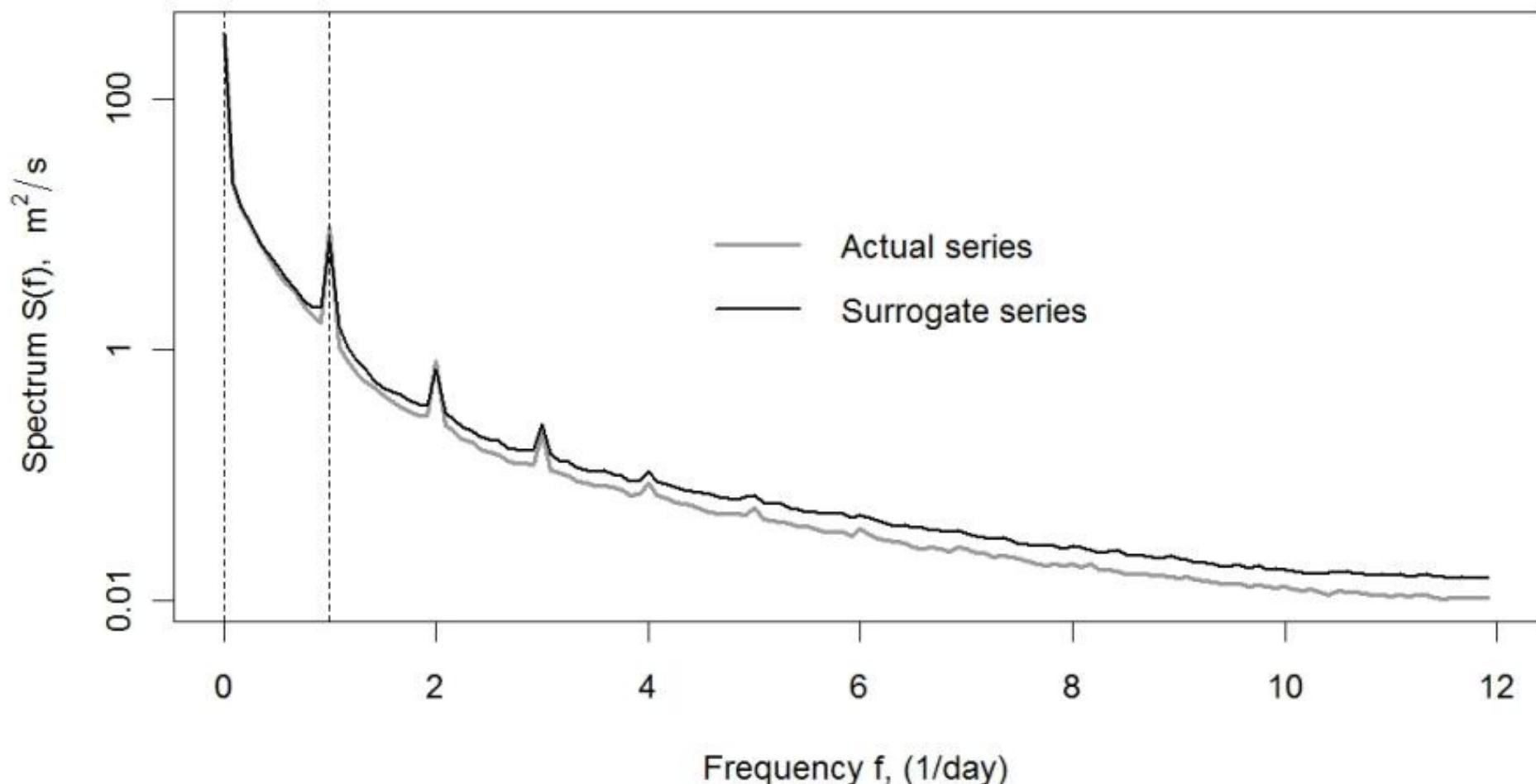
for $k=1, \dots, 45$.

Surrogate series

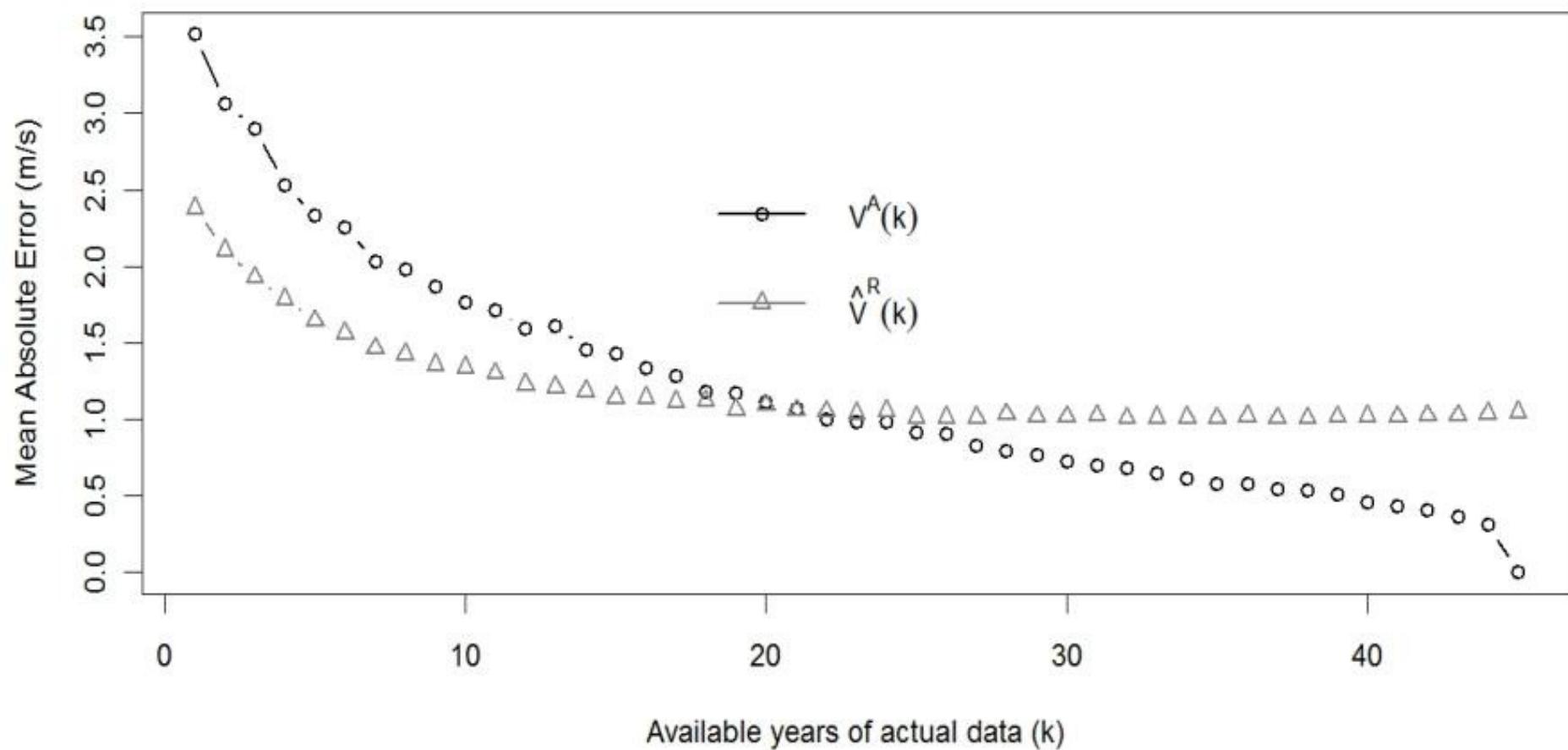


Surrogate series

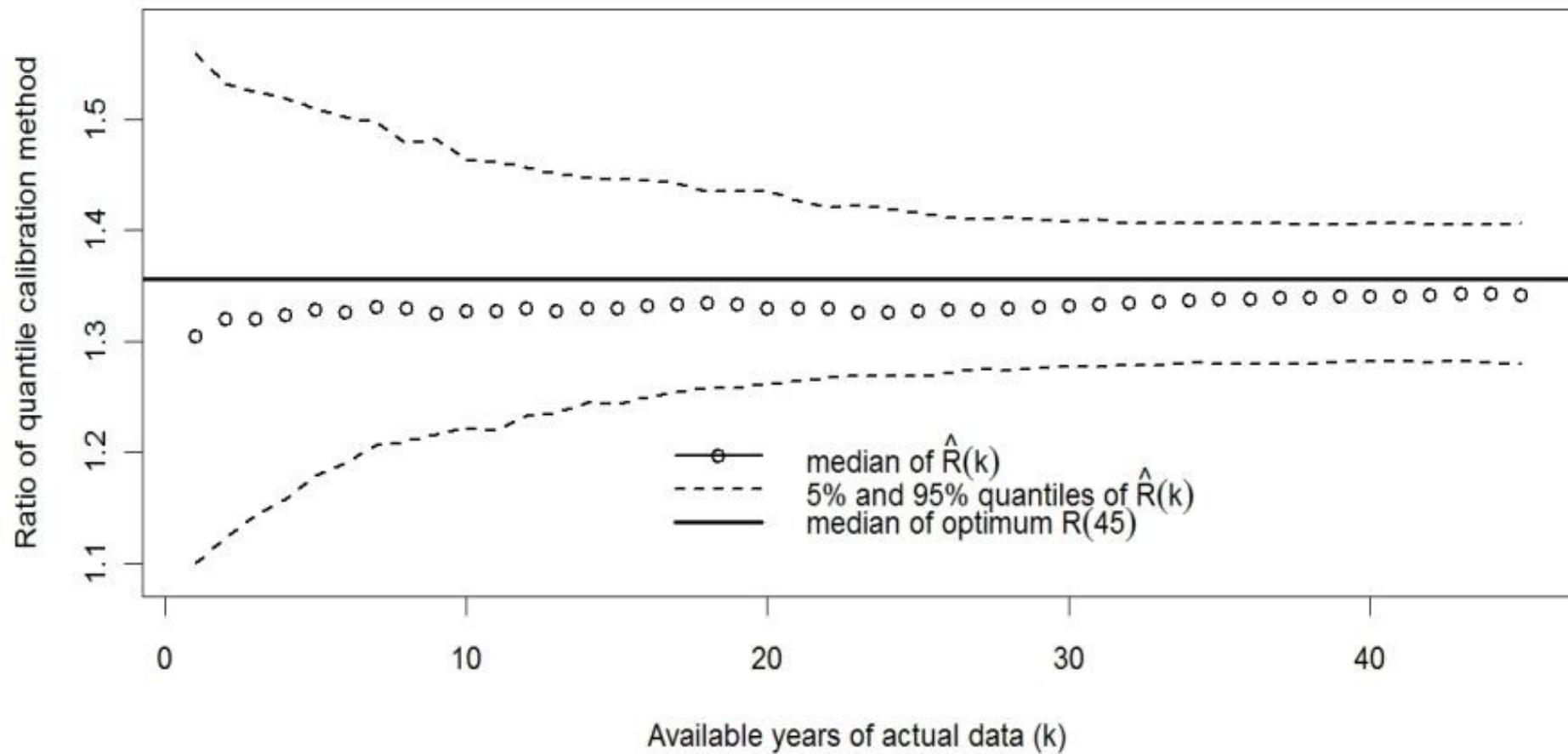
Smoothed periodogram



Results: MAE



Results: Calibrated mean ratio



Variability of 50-year returns

