

# RISK ANALYSIS AND EXTREMES

Richard L. Smith

Department of Statistics and Operations Research

University of North Carolina

Chapel Hill, NC 27599-3260

[rls@email.unc.edu](mailto:rls@email.unc.edu)

Opening Workshop

SAMSI program on Risk Analysis,  
Extreme Events and Decision Theory

September 16, 2007

## REFERENCES

Finkenstadt, B. and Rootzén, H. (editors) (2003), *Extreme Values in Finance, Telecommunications and the Environment*. Chapman and Hall/CRC Press, London.

(See <http://www.stat.unc.edu/postscript/rs/semstatrls.pdf>)

Coles, S.G. (2001), *An Introduction to Statistical Modeling of Extreme Values*. Springer Verlag, New York.

Embrechts, P., Klüppelberg, C. and Mikosch, T. (1997), *Modelling Extremal Events for Insurance and Finance*. Springer, New York.

Leadbetter, M.R., Lindgren, G. and Rootzén, H. (1983), *Extremes and Related Properties of Random Sequences and Series*. Springer Verlag, New York.

Resnick, S. (1987), *Extreme Values, Point Processes and Regular Variation*. Springer Verlag, New York.

## OUTLINE OF TALK

I. Some Examples of Risk Analysis Involving Extreme Values

II. Univariate extreme value theory

- Probability Models
- Estimation
- Diagnostics

III. Insurance Extremes I

IV. Insurance Extremes II

V. Trends in Extreme Rainfall Events

# I. SOME EXAMPLES OF RISK ANALYSIS INVOLVING EXTREME VALUES

- Insurance Claims by a Large Oil Company (Smith and Goodman 2000)
- Returns on Stock Prices
- Hurricanes — are they becoming more frequent/stronger/more costly and if so, is the change in any way connected with global warming? (Tom Knutson's talk on Tuesday)
- Are extreme rainfall events becoming more frequent? (and does that have anything to do with global warming?)

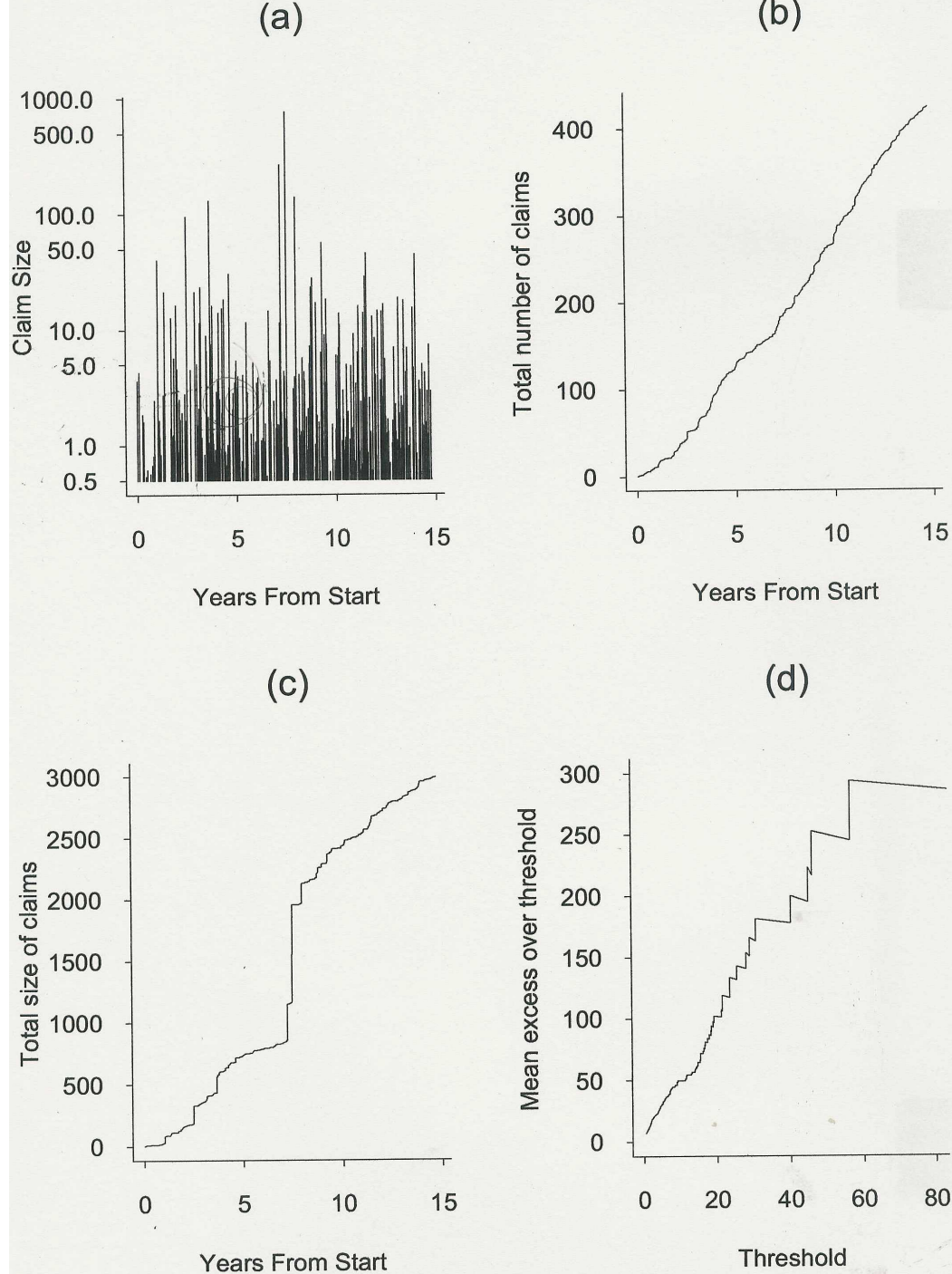
From Smith and Goodman (2000) —

The data consist of all insurance claims experienced by a large international oil company over a threshold 0.5 during a 15-year period — a total of 393 claims. Seven types:

Type	Description	Number	Mean
1	Fire	175	11.1
2	Liability	17	12.2
3	Offshore	40	9.4
4	Cargo	30	3.9
5	Hull	85	2.6
6	Onshore	44	2.7
7	Aviation	2	1.6

Total of all 393 claims: 2989.6

10 largest claims: 776.2, 268.0, 142.0, 131.0, 95.8, 56.8, 46.2, 45.2, 40.4, 30.7.



Some plots of the insurance data.

Some problems:

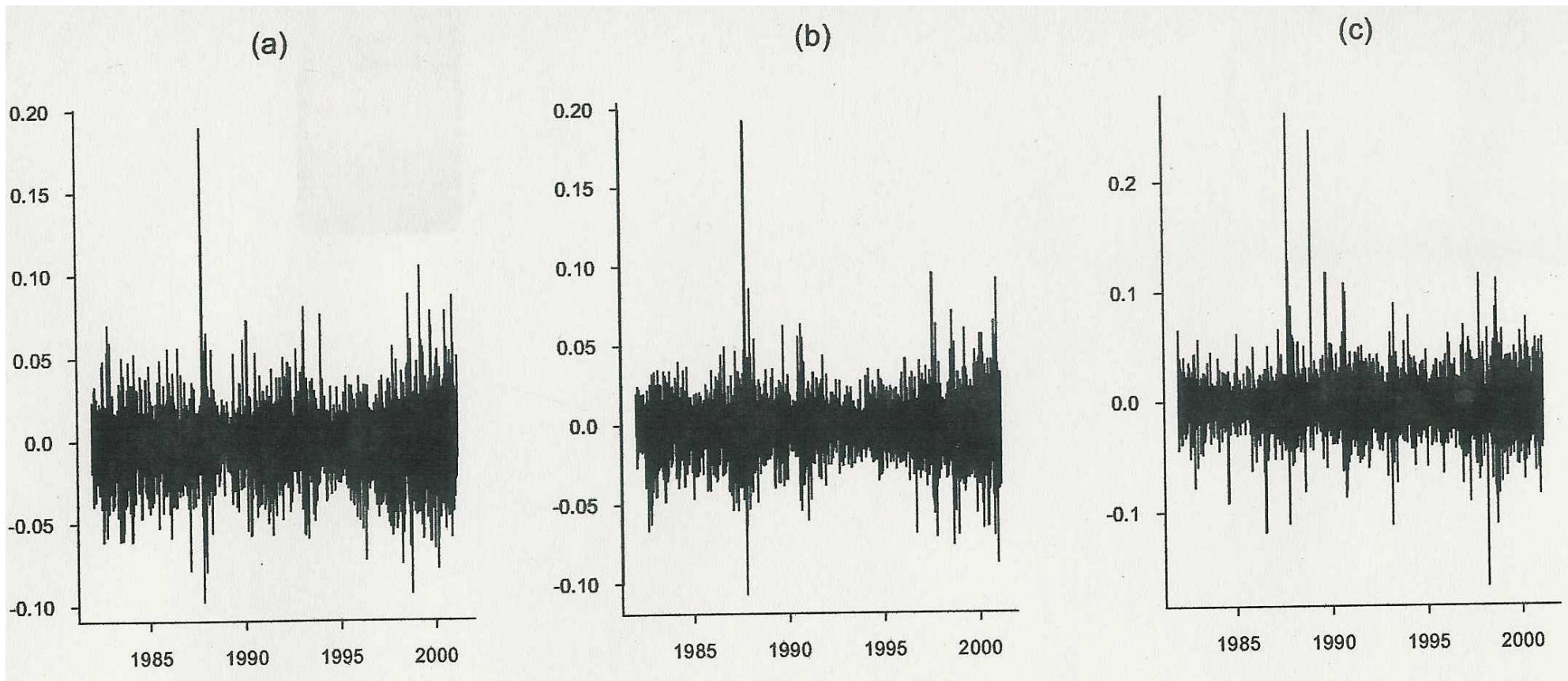
1. What is the distribution of very large claims?
2. Is there any evidence of a change of the distribution over time?
3. What is the influence of the different types of claim?
4. How should one characterize the risk to the company? More precisely, what probability distribution can one put on the amount of money that the company will have to pay out in settlement of large insurance claims over a future time period of, say, three years?

## Stock Market Returns

The next page shows negative daily returns from closing prices of 1982-2001 stock prices in three companies, Pfizer, GE and Citibank. Typical questions here are

1. How to determine *Value at Risk*, i.e. the amount which might be lost in a portfolio of assets over a specified time period with a specified small probability,
2. Dependence among the extremes of different series, and application to the portfolio management problem,
3. Modeling extremes in the presence of volatility.





Daily returns of Pfizer, GE and Citibank (thanks to Zhengjun Zhang)

## II. UNIVARIATE EXTREME VALUE THEORY

# EXTREME VALUE DISTRIBUTIONS

$X_1, X_2, \dots$ , i.i.d.,  $F(x) = \Pr\{X_i \leq x\}$ ,  $M_n = \max(X_1, \dots, X_n)$ ,  
 $\Pr\{M_n \leq x\} = F(x)^n$ .

For non-trivial results must *renormalize*: find  $a_n > 0, b_n$  such that

$$\Pr\left\{\frac{M_n - b_n}{a_n} \leq x\right\} = F(a_n x + b_n)^n \rightarrow H(x).$$

The *Three Types Theorem* (Fisher-Tippett, Gnedenko) asserts that if nondegenerate  $H$  exists, it must be one of three types:

$$\begin{aligned} H(x) &= \exp(-e^{-x}), \text{ all } x && \text{(Gumbel)} \\ H(x) &= \begin{cases} 0 & x < 0 \\ \exp(-x^{-\alpha}) & x > 0 \end{cases} && \text{(Fréchet)} \\ H(x) &= \begin{cases} \exp(-|x|^\alpha) & x < 0 \\ 1 & x > 0 \end{cases} && \text{(Weibull)} \end{aligned}$$

In Fréchet and Weibull,  $\alpha > 0$ .

The three types may be combined into a single *generalized extreme value* (GEV) distribution:

$$H(x) = \exp \left\{ - \left( 1 + \xi \frac{x - \mu}{\psi} \right)_+^{-1/\xi} \right\},$$

( $y_+ = \max(y, 0)$ )

where  $\mu$  is a location parameter,  $\psi > 0$  is a scale parameter and  $\xi$  is a shape parameter.  $\xi \rightarrow 0$  corresponds to the Gumbel distribution,  $\xi > 0$  to the Fréchet distribution with  $\alpha = 1/\xi$ ,  $\xi < 0$  to the Weibull distribution with  $\alpha = -1/\xi$ .

$\xi > 0$ : “long-tailed” case,  $1 - F(x) \propto x^{-1/\xi}$ ,

$\xi = 0$ : “exponential tail”

$\xi < 0$ : “short-tailed” case, finite endpoint at  $\mu - \xi/\psi$

# EXCEEDANCES OVER THRESHOLDS

Consider the distribution of  $X$  conditionally on exceeding some high threshold  $u$ :

$$F_u(y) = \frac{F(u + y) - F(u)}{1 - F(u)}.$$

As  $u \rightarrow \omega_F = \sup\{x : F(x) < 1\}$ , often find a limit

$$F_u(y) \approx G(y; \sigma_u, \xi)$$

where  $G$  is *generalized Pareto distribution* (GPD)

$$G(y; \sigma, \xi) = 1 - \left(1 + \xi \frac{y}{\sigma}\right)_+^{-1/\xi}.$$

Equivalence to three types theorem established by Pickands (1975).

## The Generalized Pareto Distribution

$$G(y; \sigma, \xi) = 1 - \left(1 + \xi \frac{y}{\sigma}\right)_+^{-1/\xi}.$$

$\xi > 0$ : long-tailed (equivalent to usual Pareto distribution), tail like  $x^{-1/\xi}$ ,

$\xi = 0$ : take limit as  $\xi \rightarrow 0$  to get

$$G(y; \sigma, 0) = 1 - \exp\left(-\frac{y}{\sigma}\right),$$

i.e. exponential distribution with mean  $\sigma$ ,

$\xi < 0$ : finite upper endpoint at  $-\sigma/\xi$ .

# POISSON-GPD MODEL FOR EXCEEDANCES

1. The number,  $N$ , of exceedances of the level  $u$  in any one year has a Poisson distribution with mean  $\lambda$ ,
2. Conditionally on  $N \geq 1$ , the excess values  $Y_1, \dots, Y_N$  are IID from the GPD.

*Relation to GEV for annual maxima:*

Suppose  $x > u$ . The probability that the annual maximum of the Poisson-GPD process is less than  $x$  is

$$\begin{aligned} \Pr\{\max_{1 \leq i \leq N} Y_i \leq x\} &= \Pr\{N = 0\} + \sum_{n=1}^{\infty} \Pr\{N = n, Y_1 \leq x, \dots, Y_n \leq x\} \\ &= e^{-\lambda} + \sum_{n=1}^{\infty} \frac{\lambda^n e^{-\lambda}}{n!} \left\{ 1 - \left( 1 + \xi \frac{x - u}{\sigma} \right)^{-1/\xi} \right\}^n \\ &= \exp \left\{ -\lambda \left( 1 + \xi \frac{x - u}{\sigma} \right)^{-1/\xi} \right\}. \end{aligned}$$

This is GEV with  $\sigma = \psi + \xi(u - \mu)$ ,  $\lambda = \left( 1 + \xi \frac{u - \mu}{\psi} \right)^{-1/\xi}$ . Thus the GEV and GPD models are entirely consistent with one another above the GPD threshold, and moreover, shows exactly how the Poisson-GPD parameters  $\sigma$  and  $\lambda$  vary with  $u$ .



# ALTERNATIVE PROBABILITY MODELS

## 1. The $r$ largest order statistics model

If  $Y_{n,1} \geq Y_{n,2} \geq \dots \geq Y_{n,r}$  are  $r$  largest order statistics of IID sample of size  $n$ , and  $a_n$  and  $b_n$  are EVT normalizing constants, then

$$\left( \frac{Y_{n,1} - b_n}{a_n}, \dots, \frac{Y_{n,r} - b_n}{a_n} \right)$$

converges in distribution to a limiting random vector  $(X_1, \dots, X_r)$ , whose density is

$$h(x_1, \dots, x_r) = \psi^{-r} \exp \left\{ - \left( 1 + \xi \frac{x_r - \mu}{\psi} \right)^{-1/\xi} - \left( 1 + \frac{1}{\xi} \right) \sum_{j=1}^r \log \left( 1 + \xi \frac{x_j - \mu}{\psi} \right) \right\}.$$

## 2. Point process approach (Smith 1989)

Two-dimensional plot of exceedance times and exceedance levels forms a nonhomogeneous Poisson process with

$$\begin{aligned}\Lambda(A) &= (t_2 - t_1)\Psi(y; \mu, \psi, \xi) \\ \Psi(y; \mu, \psi, \xi) &= \left(1 + \xi \frac{y - \mu}{\psi}\right)^{-1/\xi}\end{aligned}$$

$(1 + \xi(y - \mu)/\psi > 0)$ .

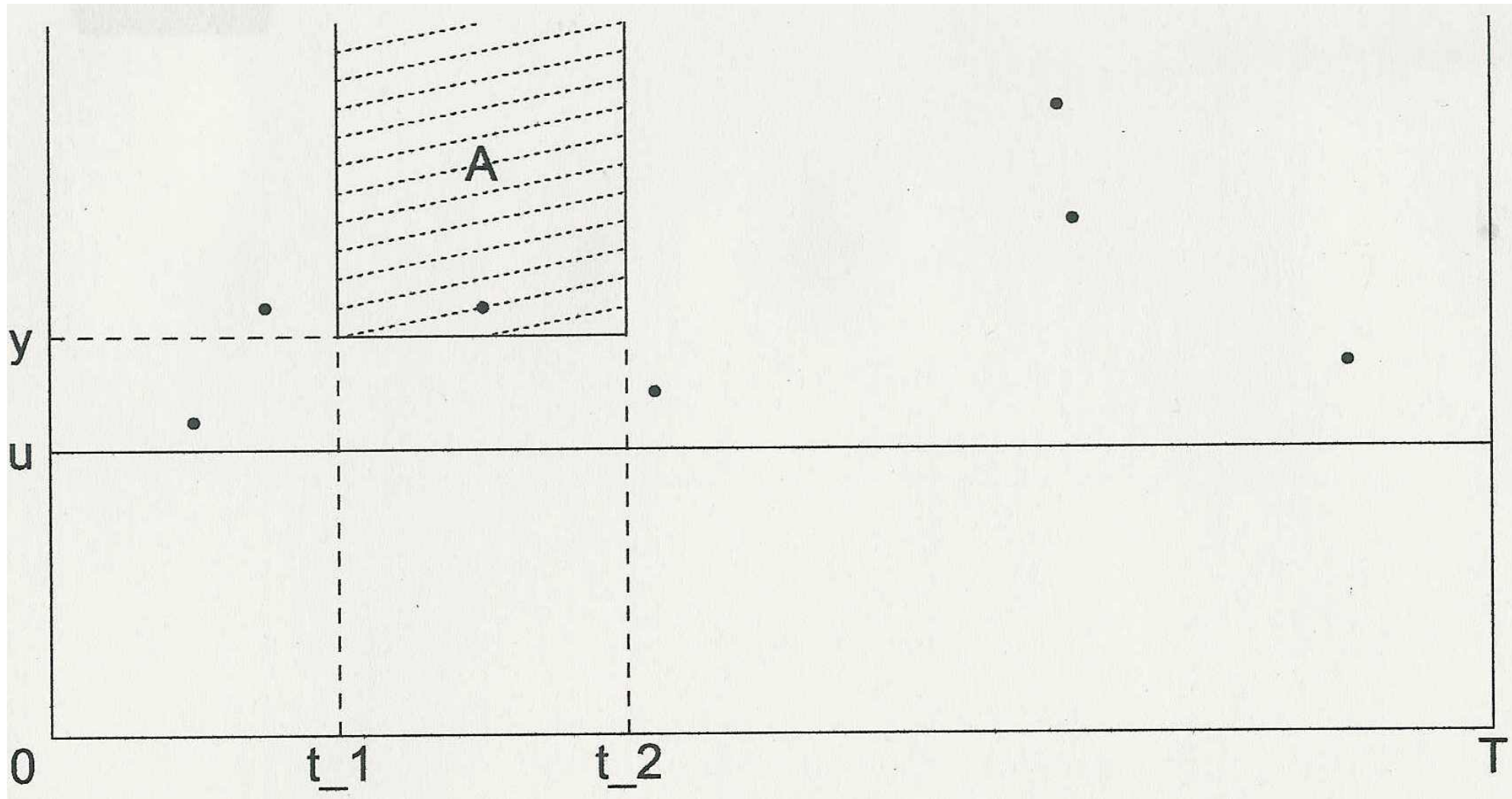
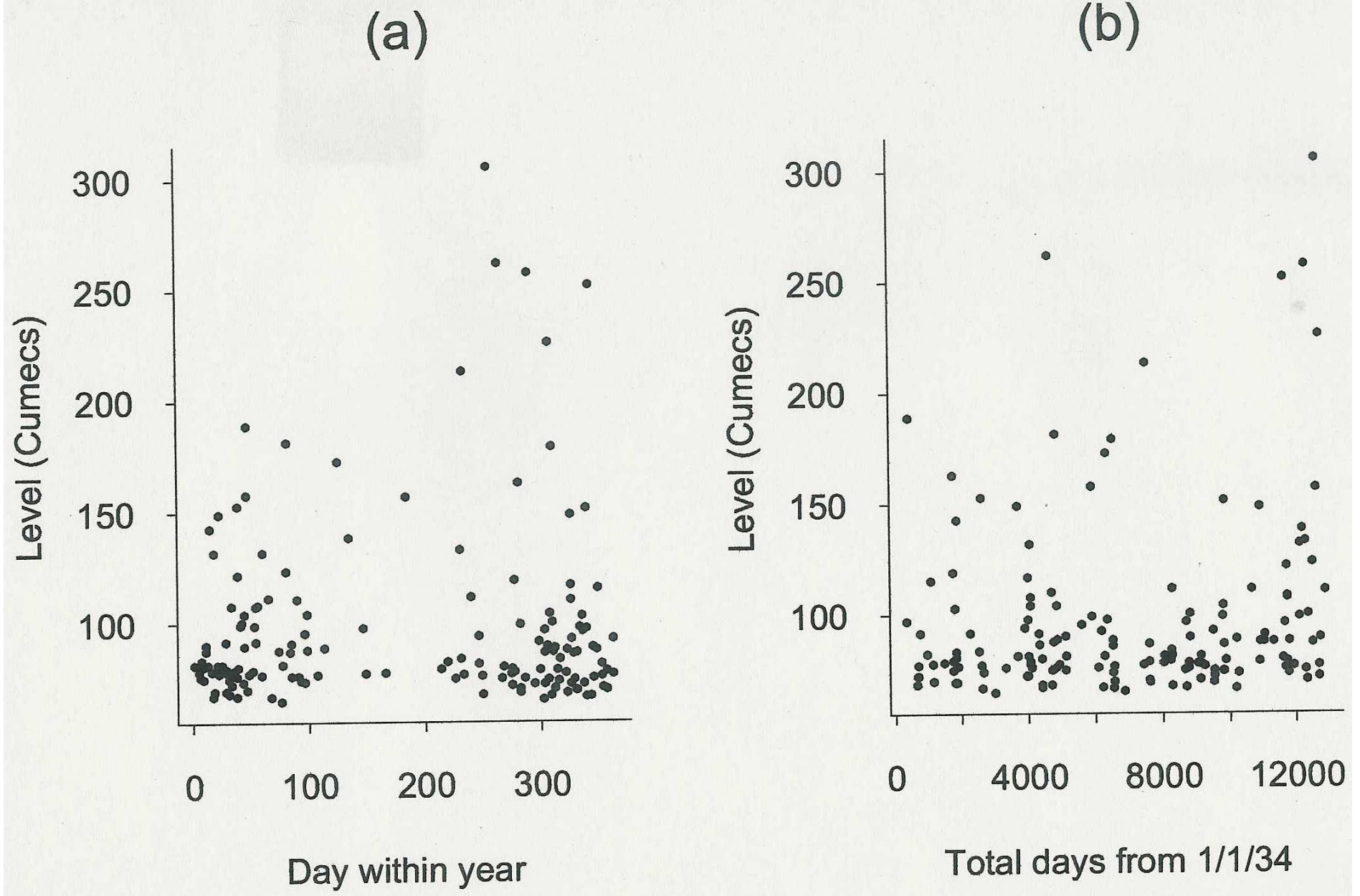


Illustration of point process model.

An extension of this approach allows for nonstationary processes in which the parameters  $\mu$ ,  $\psi$  and  $\xi$  are all allowed to be time-dependent, denoted  $\mu_t$ ,  $\psi_t$  and  $\xi_t$ .

This is the basis of the extreme value regression approaches introduced later



Plots of exceedances of River Nidd, (a) against day within year, (b) against total days from January 1, 1934. Adapted from Davison and Smith (1990).

# ESTIMATION

GEV log likelihood:

$$\begin{aligned} \ell_Y(\mu, \psi, \xi) = & -N \log \psi - \left(\frac{1}{\xi} + 1\right) \sum_i \log \left(1 + \xi \frac{Y_i - \mu}{\psi}\right) \\ & - \sum_i \left(1 + \xi \frac{Y_i - \mu}{\psi}\right)^{-1/\xi} \end{aligned}$$

provided  $1 + \xi(Y_i - \mu)/\psi > 0$  for each  $i$ .

Poisson-GPD model:

$$\ell_{N,Y}(\lambda, \sigma, \xi) = N \log \lambda - \lambda T - N \log \sigma - \left(1 + \frac{1}{\xi}\right) \sum_{i=1}^N \log \left(1 + \xi \frac{Y_i}{\sigma}\right)$$

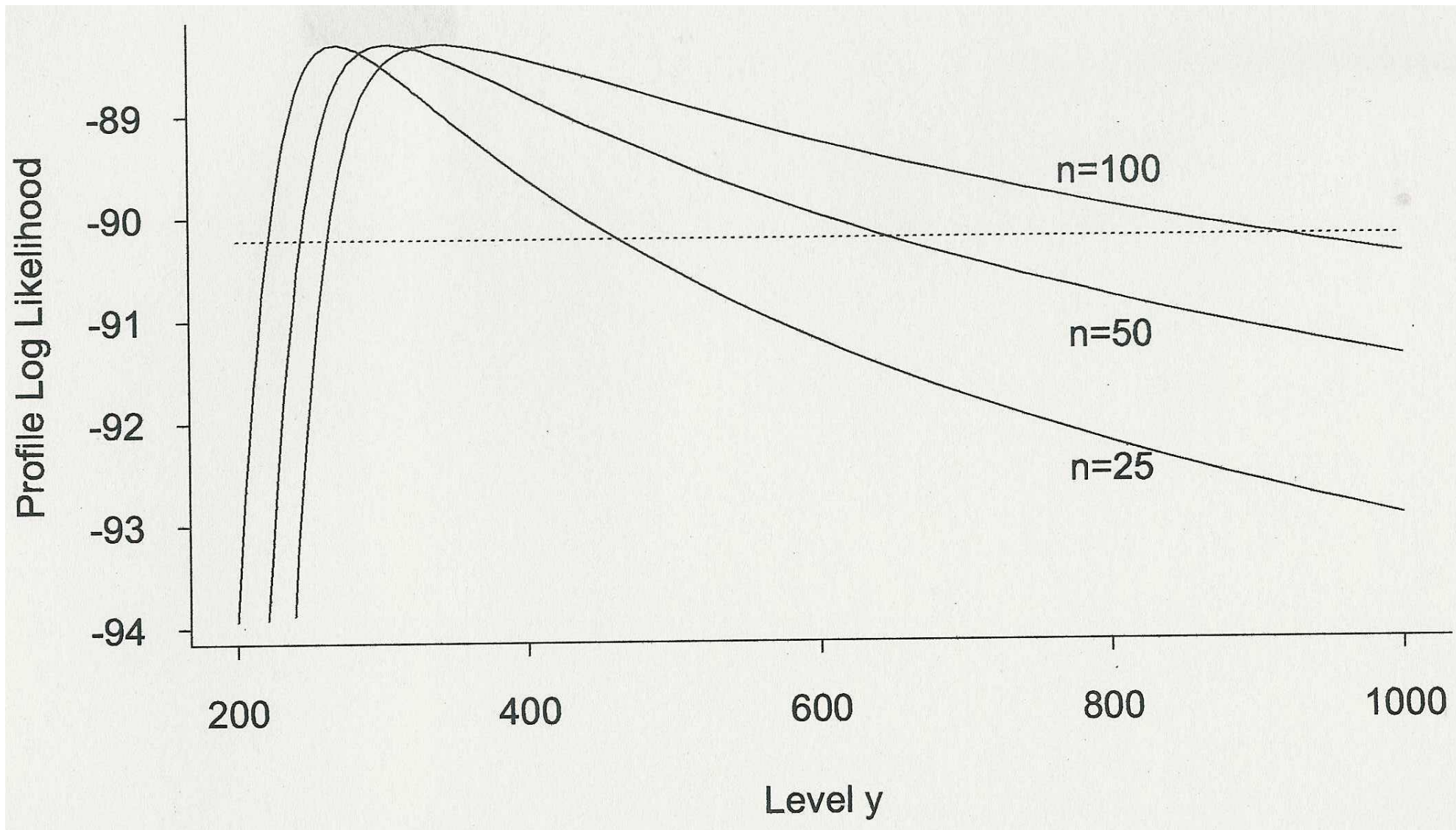
provided  $1 + \xi Y_i/\sigma > 0$  for all  $i$ .

Usual asymptotics valid if  $\xi > -\frac{1}{2}$  (Smith 1985)

## Profile Likelihoods for Quantiles

Suppose we are interested in the  $N$ -year return level  $y_N$ , i.e. the  $(1 - 1/N)$ -quantile of the annual maximum distribution. We can construct a *profile likelihood* by reparameterizing the GEV so that  $y_N$  is one of the three parameters, and maximizing with respect to the other two. Likelihood ratio asymptotics can then be used to construct a confidence interval for  $y_N$ .

Example from the Nidd data:



Profile log-likelihoods for extreme quantiles based on Nidd data



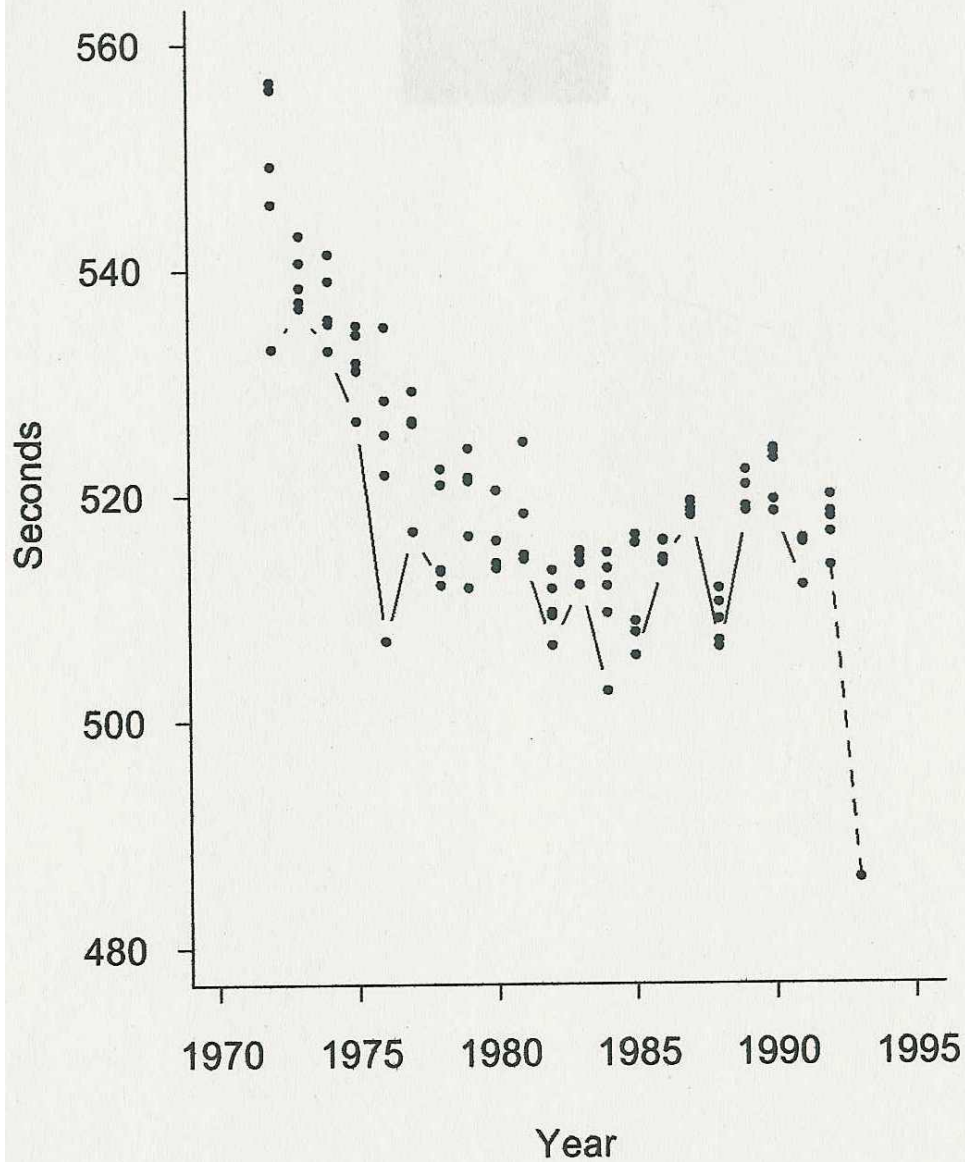
## Bayesian approaches

An alternative approach to extreme value inference is Bayesian, using vague priors for the GEV parameters and MCMC samples for the computations. Bayesian methods are particularly useful for *predictive inference*, e.g. if  $Z$  is some as yet unobserved random variable whose distribution depends on  $\mu, \psi$  and  $\xi$ , estimate  $\Pr\{Z > z\}$  by

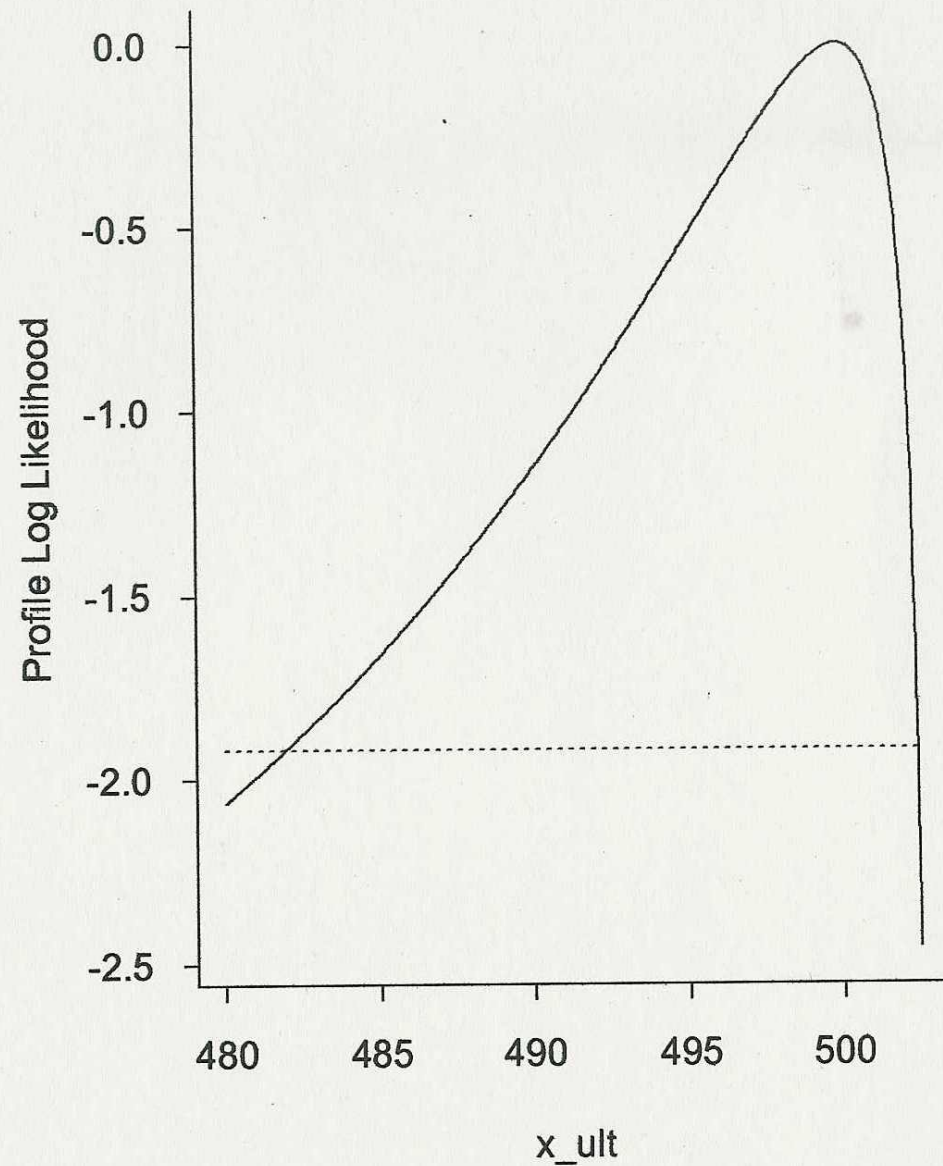
$$\int \Pr\{Z > z; \mu, \psi, \xi\} \pi(\mu, \psi, \xi | Y) d\mu d\psi d\xi$$

where  $\pi(\dots|Y)$  denotes the posterior density given past data  $Y$

(a)



(b)



Plots of women's 3000 meter records, and profile log-likelihood for ultimate best value based on pre-1993 data.

*Example.* The left figure shows the five best running times by different athletes in the women's 3000 metre track event for each year from 1972 to 1992. Also shown on the plot is Wang Junxia's world record from 1993. Many questions were raised about possible illegal drug use.

We approach this by asking how implausible Wang's performance was, given all data up to 1992.

Robinson and Tawn (1995) used the  $r$  largest order statistics method (with  $r = 5$ , translated to smallest order statistics) to estimate an extreme value distribution, and hence computed a profile likelihood for  $x_{\text{ult}}$ , the lower endpoint of the distribution, based on data up to 1992 (right plot of previous figure)

*Alternative Bayesian calculation:*

(Smith 1997)

Compute the (Bayesian) predictive probability that the 1993 performance is equal or better to Wang's, given the data up to 1992, and conditional on the event that there is a new world record.

The answer is approximately 0.0006.

# DIAGNOSTICS

Gumbel plots

QQ plots of residuals

Mean excess plot

Z- and W-statistic plots

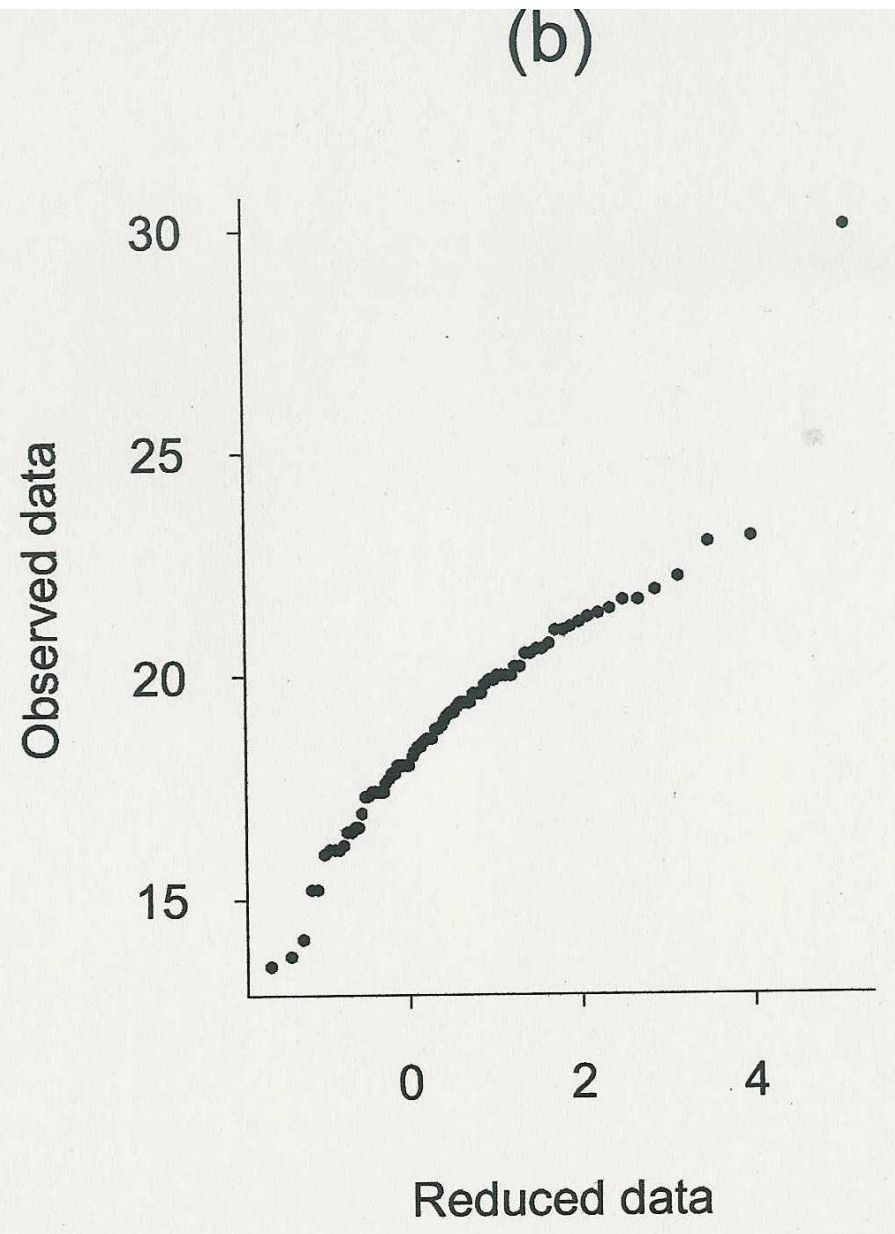
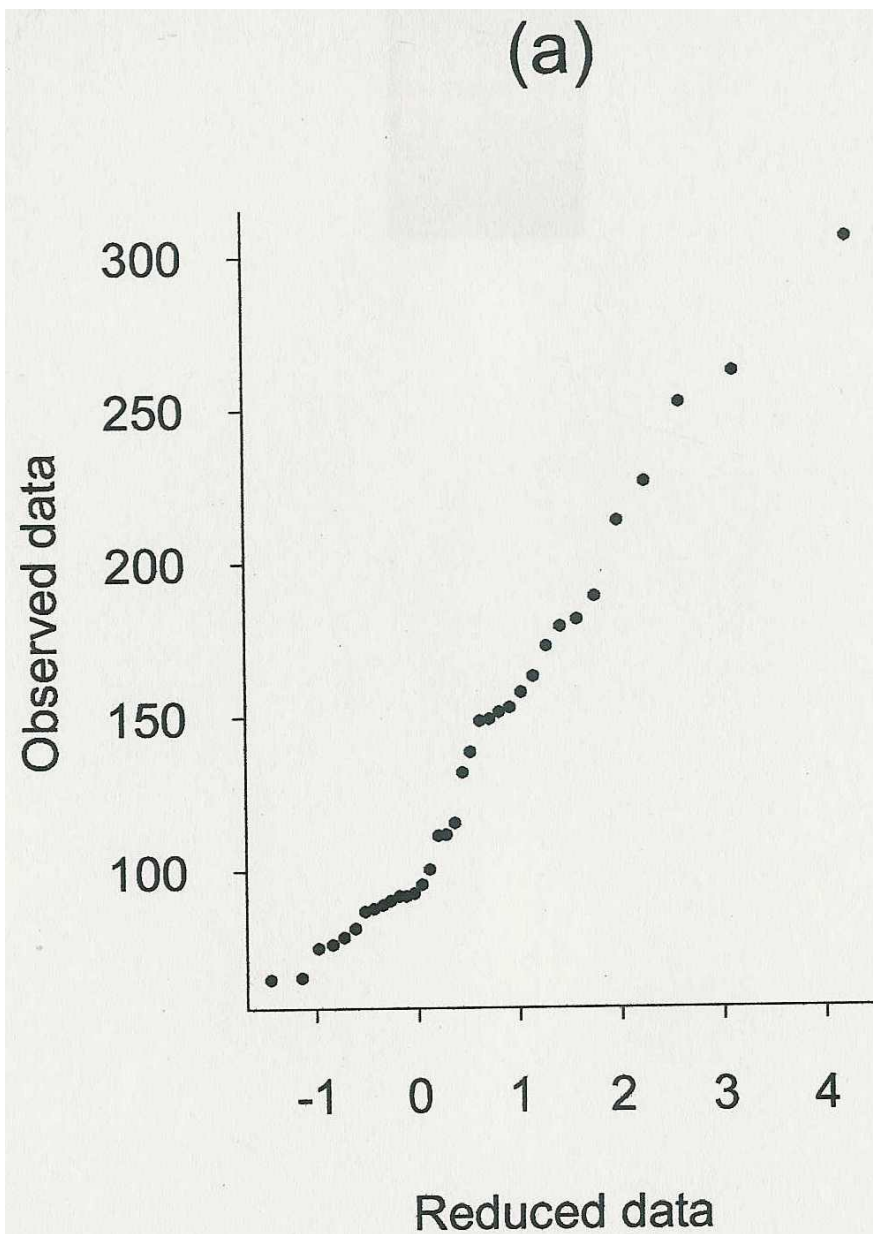
## *Gumbel plots*

Used as a diagnostic for Gumbel distribution with annual maxima data. Order data as  $Y_{1:N} \leq \dots \leq Y_{N:N}$ , then plot  $Y_{i:N}$  against *reduced value*  $x_{i:N}$ ,

$$x_{i:N} = -\log(-\log p_{i:N}),$$

$p_{i:N}$  being the  $i$ 'th *plotting position*, usually taken to be  $(i - \frac{1}{2})/N$ .

A straight line is ideal. Curvature may indicate Fréchet or Weibull form. Also look for outliers.



Gumbel plots. (a) Annual maxima for River Nidd flow series. (b) Annual maximum temperatures in Ivigtut, Iceland.

## QQ plots of residuals

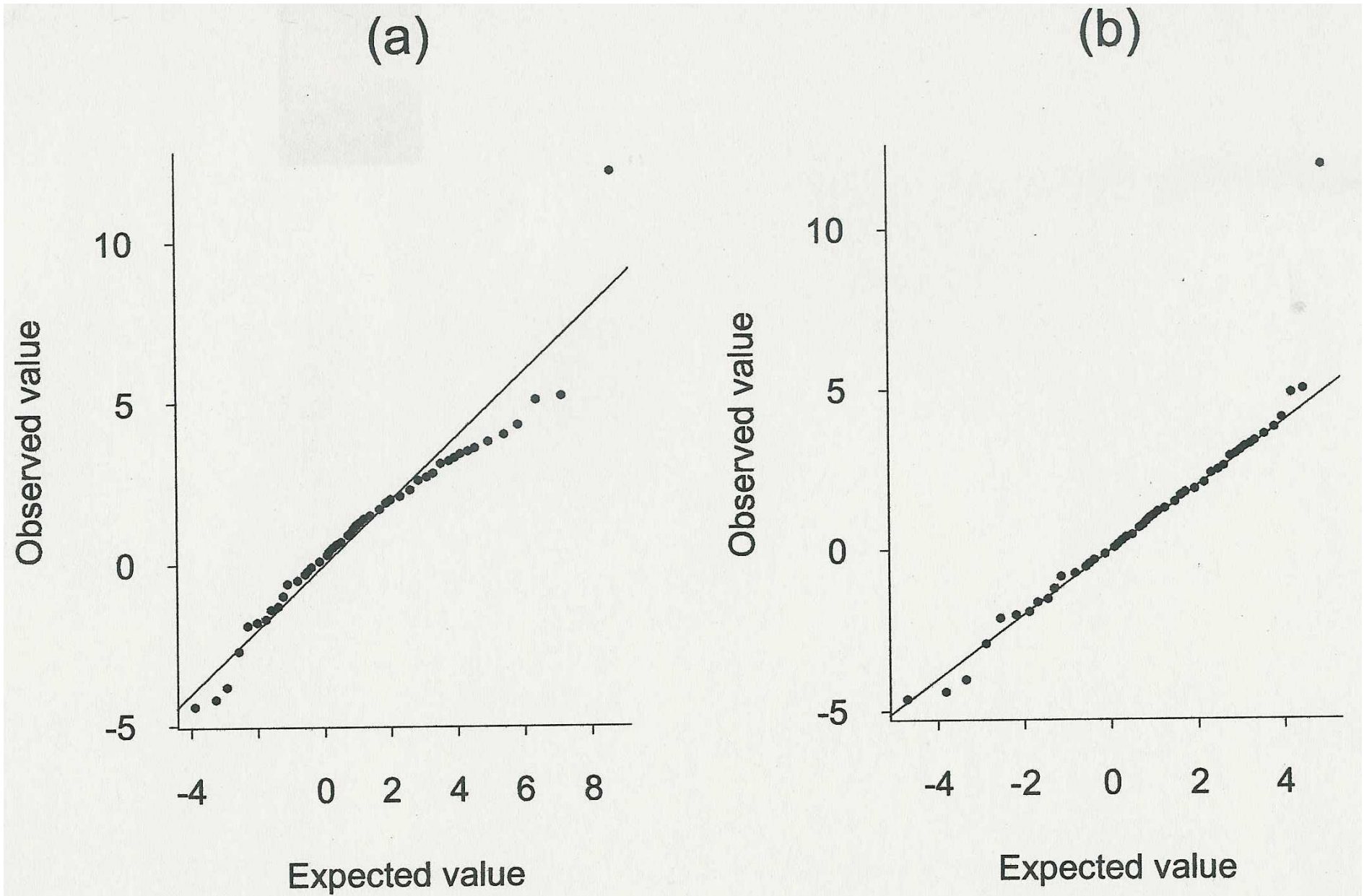
A second type of probability plot is drawn *after* fitting the model. Suppose  $Y_1, \dots, Y_N$  are IID observations whose common distribution function is  $G(y; \theta)$  depending on parameter vector  $\theta$ . Suppose  $\theta$  has been estimated by  $\hat{\theta}$ , and let  $G^{-1}(p; \theta)$  denote the inverse distribution function of  $G$ , written as a function of  $\theta$ . A QQ (quantile-quantile) plot consists of first ordering the observations  $Y_{1:N} \leq \dots \leq Y_{N:N}$ , and then plotting  $Y_{i:N}$  against the reduced value

$$x_{i:N} = G^{-1}(p_{i:N}; \hat{\theta}),$$

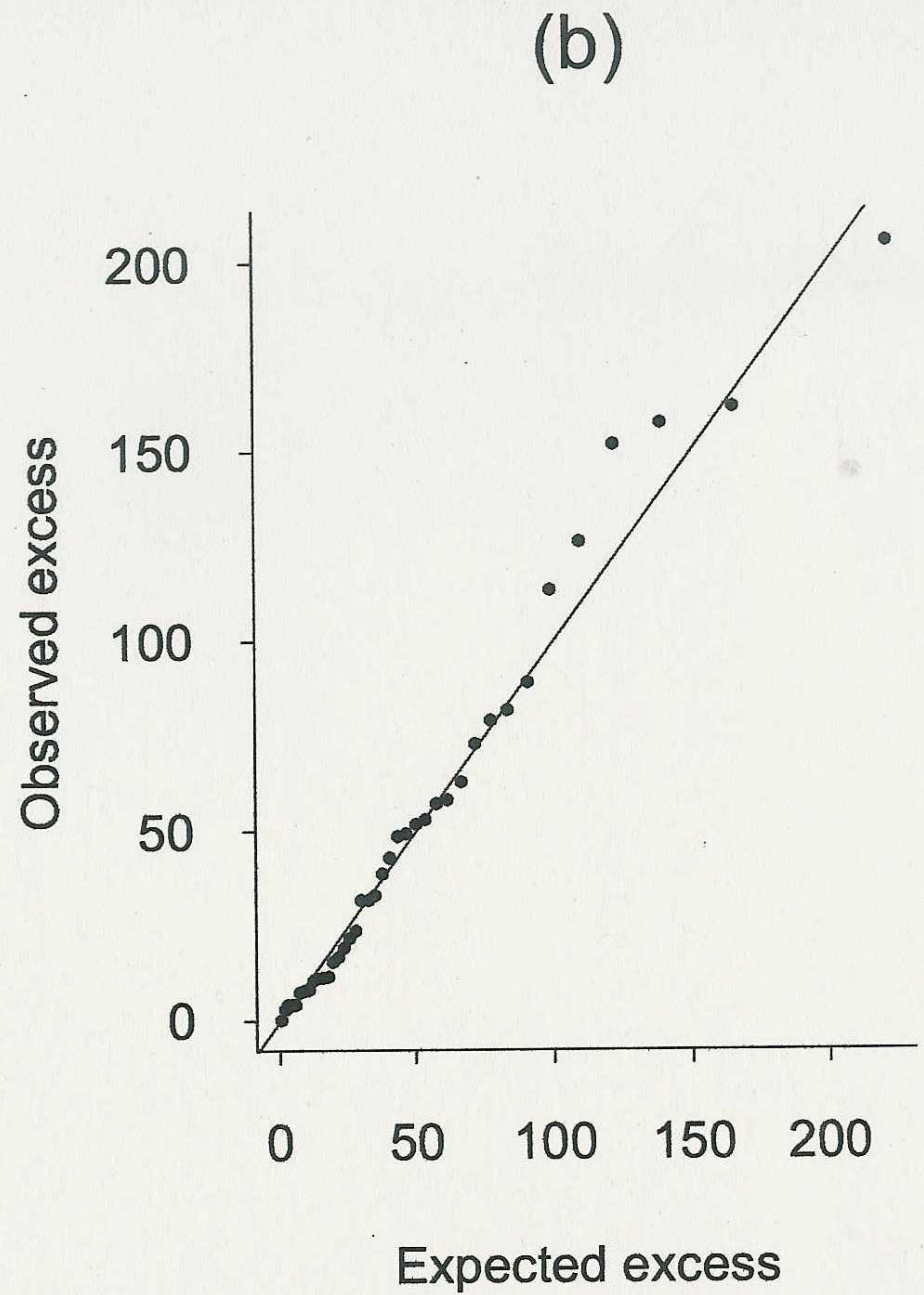
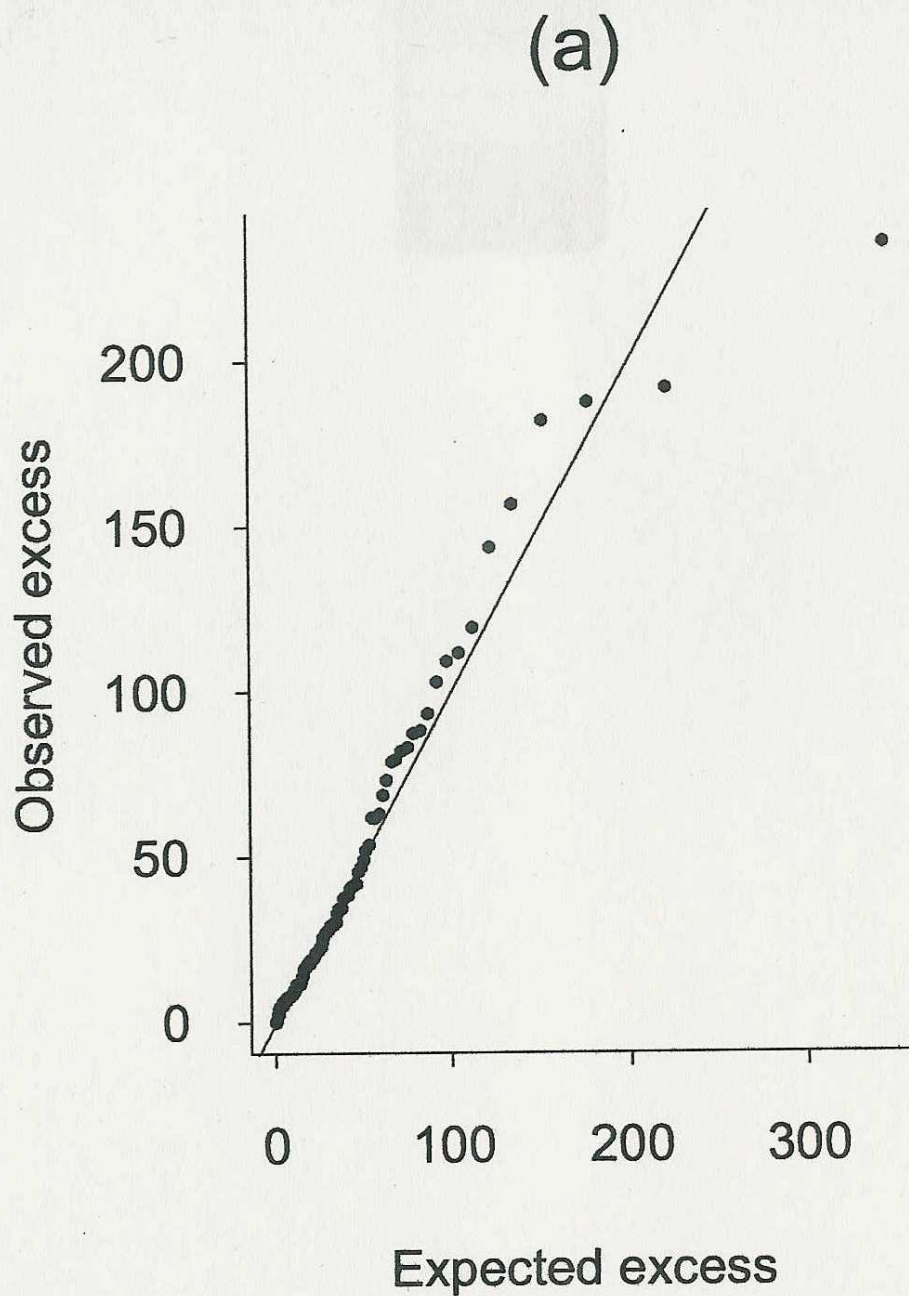
where  $p_{i:N}$  may be taken as  $(i - \frac{1}{2})/N$ . If the model is a good fit, the plot should be roughly a straight line of unit slope through the origin.

Examples...





GEV model to Ivigtut data, (a) without adjustment, (b) excluding largest value from model fit but including it in the plot.



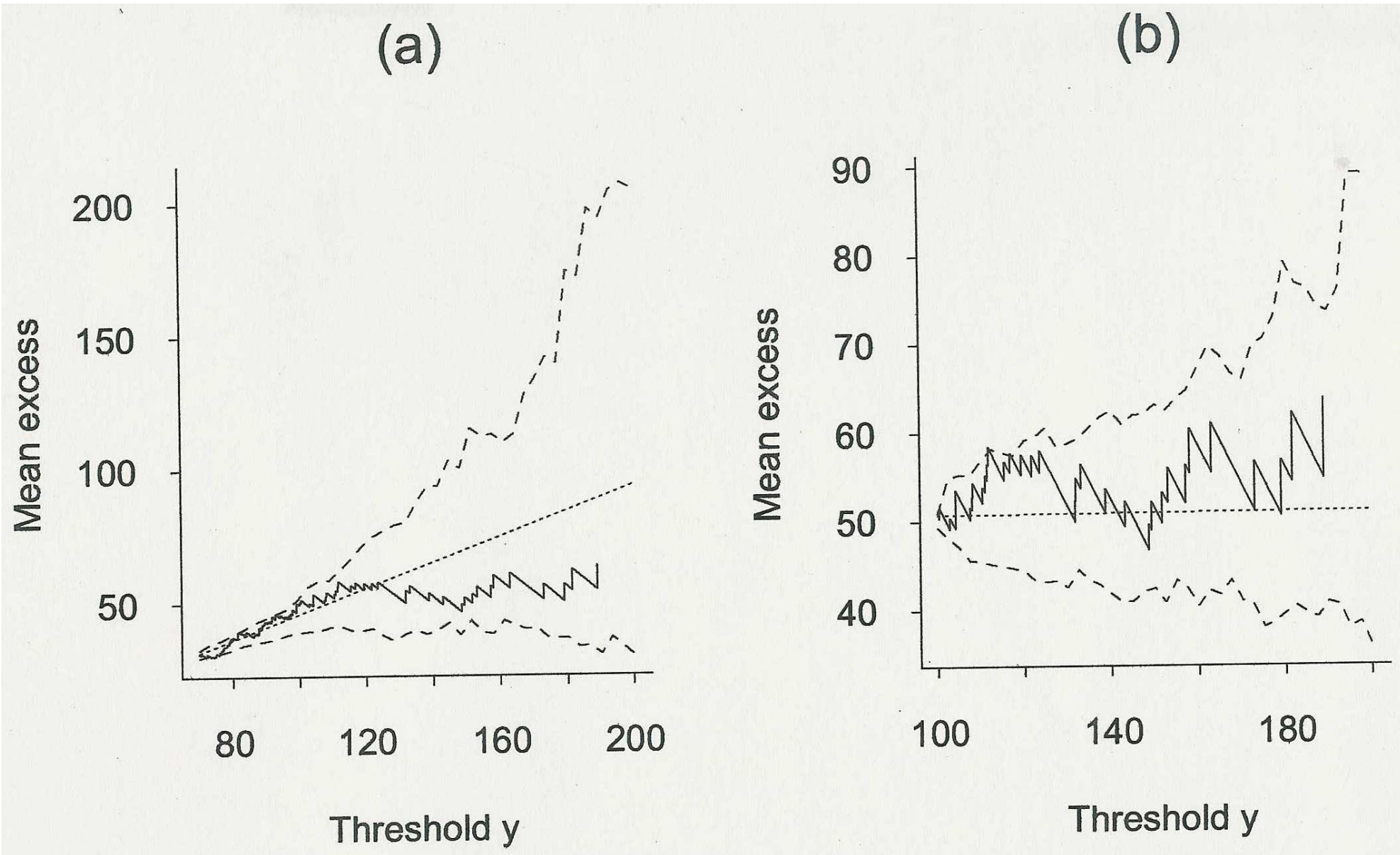
QQ plots for GPD, Nidd data. (a)  $u = 70$ . (b)  $u = 100$ .

## *Mean excess plot*

Idea: for a sequence of values of  $w$ , plot the mean excess over  $w$  against  $w$  itself. If the GPD is a good fit, the plot should be approximately a straight line.

In practice, the actual plot is very jagged and therefore its “straightness” is difficult to assess. However, a Monte Carlo technique, *assuming* the GPD is valid throughout the range of the plot, can be used to assess this.

Examples...



Mean excess over threshold plots for Nidd data, with Monte Carlo confidence bands, relative to threshold 70 (a) and 100 (b).

## *Z- and W-statistic plots*

Consider nonstationary model with  $\mu_t, \psi_t, \xi_t$  dependent on  $t$ .

$Z$  statistic based on intervals between exceedances  $T_k$ :

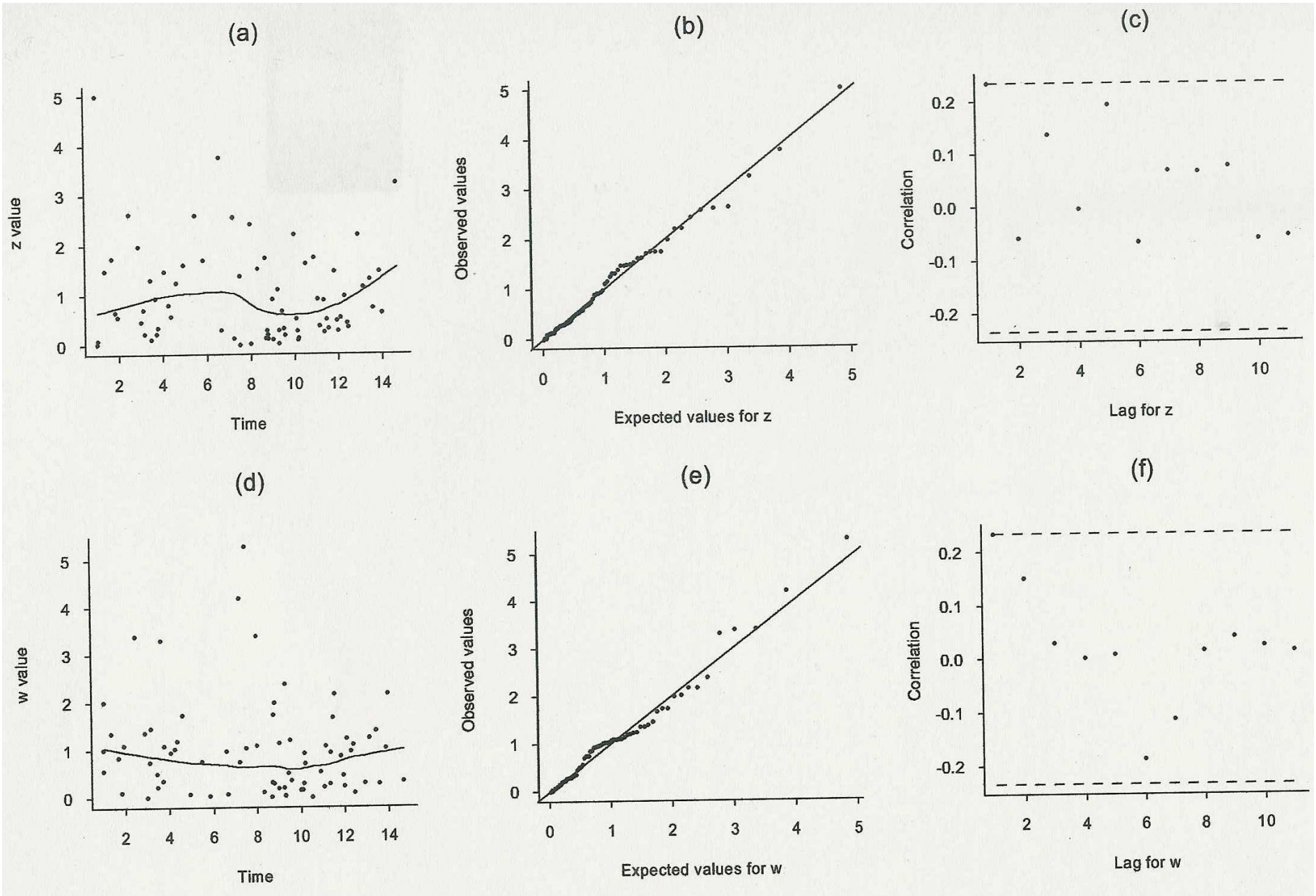
$$Z_k = \int_{T_{k-1}}^{T_k} \lambda_u(s) ds,$$
$$\lambda_u(s) = \{1 + \xi_s(u - \mu_s)/\psi_s\}^{-1/\xi_s}.$$

$W$  statistic based on excess values: if  $Y_k$  is excess over threshold at time  $T_k$ ,

$$W_k = \frac{1}{\xi_{T_k}} \log \left\{ 1 + \frac{\xi_{T_k} Y_k}{\psi_{T_k} + \xi_{T_k} (u - \mu_{T_k})} \right\}.$$

Idea: if the model is exact, both  $Z_k$  and  $W_k$  and i.i.d. exponential with mean 1. Can test this with various plots.





Diagnostic plots based on  $Z$  and  $W$  statistics for oil company insurance data ( $u = 5$ )

### III. INSURANCE EXTREMES I

We return to the oil company data set discussed in section I. Prior to any of the analysis, some examination was made of clustering phenomena, but this only reduced the original 425 claims to 393 “independent” claims (Smith & Goodman 2000)

GPD fits to various thresholds:

$u$	$N_u$	Mean Excess	$\sigma$	$\xi$
0.5	393	7.11	1.02	1.01
2.5	132	17.89	3.47	0.91
5	73	28.9	6.26	0.89
10	42	44.05	10.51	0.84
15	31	53.60	5.68	1.44
20	17	91.21	19.92	1.10
25	13	113.7	74.46	0.93
50	6	37.97	150.8	0.29

Point process approach:

$u$	$N_u$	$\mu$	$\log \psi$	$\xi$
0.5	393	26.5 (4.4)	3.30 (0.24)	1.00 (0.09)
2.5	132	26.3 (5.2)	3.22 (0.31)	0.91 (0.16)
5	73	26.8 (5.5)	3.25 (0.31)	0.89 (0.21)
10	42	27.2 (5.7)	3.22 (0.32)	0.84 (0.25)
15	31	22.3 (3.9)	2.79 (0.46)	1.44 (0.45)
20	17	22.7 (5.7)	3.13 (0.56)	1.10 (0.53)
25	13	20.5 (8.6)	3.39 (0.66)	0.93 (0.56)

Standard errors are in parentheses



## *Predictive Distributions of Future Losses*

What is the probability distribution of future losses over a specific time period, say 1 year?

Let  $Y$  be future total loss. Distribution function  $G(y; \mu, \psi, \xi)$  — in practice this must itself be simulated.

Traditional frequentist approach:

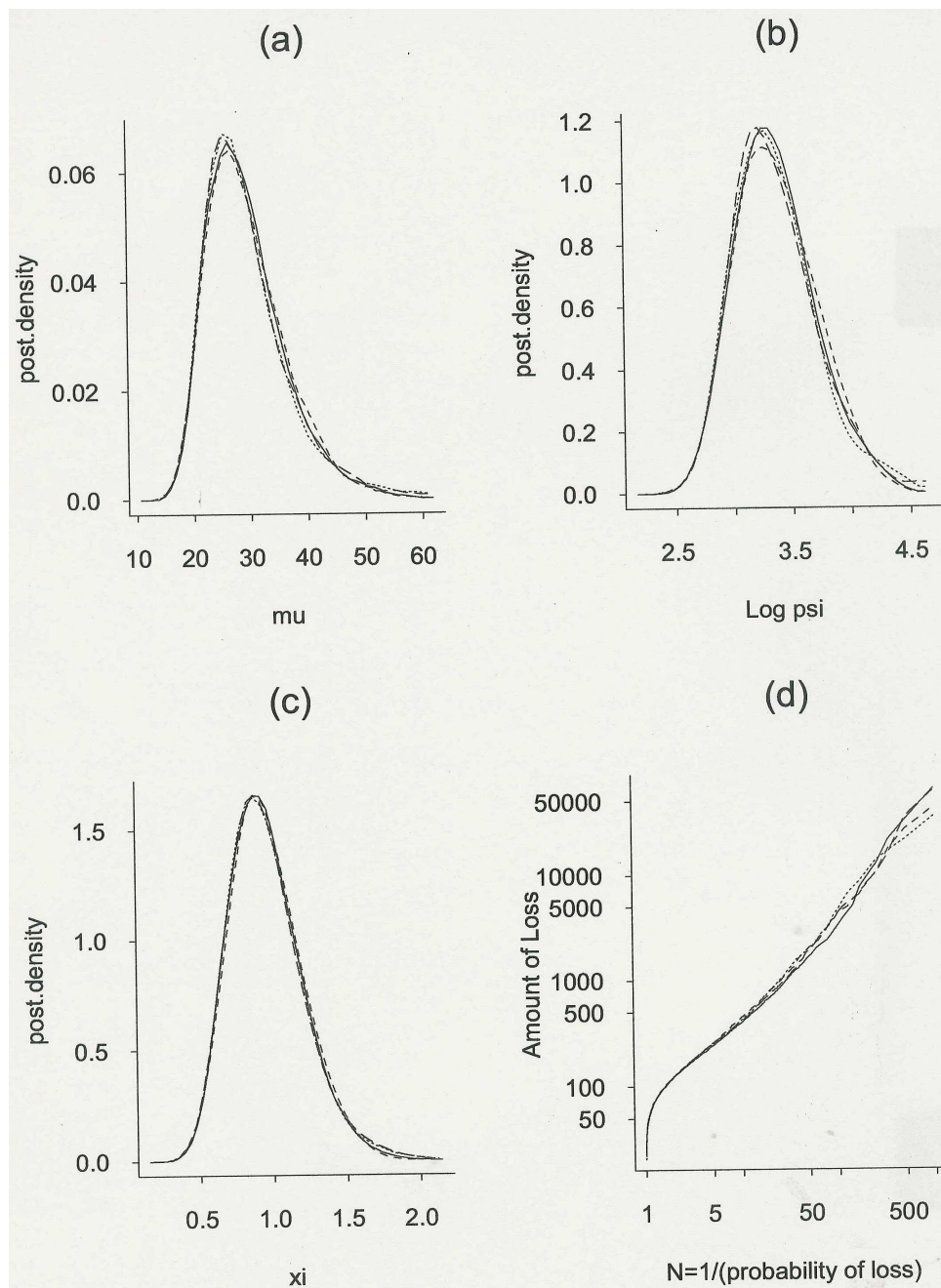
$$\hat{G}(y) = G(y; \hat{\mu}, \hat{\psi}, \hat{\xi})$$

where  $\hat{\mu}$ ,  $\hat{\psi}$ ,  $\hat{\xi}$  are MLEs.

Bayesian:

$$\tilde{G}(y) = \int G(y; \mu, \psi, \xi) d\pi(\mu, \psi, \xi | \mathbf{X})$$

where  $\pi(\cdot | \mathbf{X})$  denotes posterior density given data  $\mathbf{X}$ .



Estimated posterior densities for the three parameters, and for the predictive distribution function. Four independent Monte Carlo runs are shown for each plot.

## Hierarchical models for claim type and year effects

Further features of the data:

1. When separate GPDs are fitted to each of the 6 main types, there are clear differences among the parameters.
2. The rate of high-threshold crossings does not appear uniform, but peaks around years 10–12.

### *A Hierarchical Model:*

Level I. Parameters  $m_\mu$ ,  $m_\psi$ ,  $m_\xi$ ,  $s_\mu^2$ ,  $s_\psi^2$ ,  $s_\xi^2$  are generated from a prior distribution.

Level II. Conditional on the parameters in Level I, parameters  $\mu_1, \dots, \mu_J$  (where  $J$  is the number of types) are independently drawn from  $N(m_\mu, s_\mu^2)$ , the normal distribution with mean  $m_\mu$ , variance  $s_\mu^2$ . Similarly,  $\log \psi_1, \dots, \log \psi_J$  are drawn independently from  $N(m_\psi, s_\psi^2)$ ,  $\xi_1, \dots, \xi_J$  are drawn independently from  $N(m_\xi, s_\xi^2)$ .

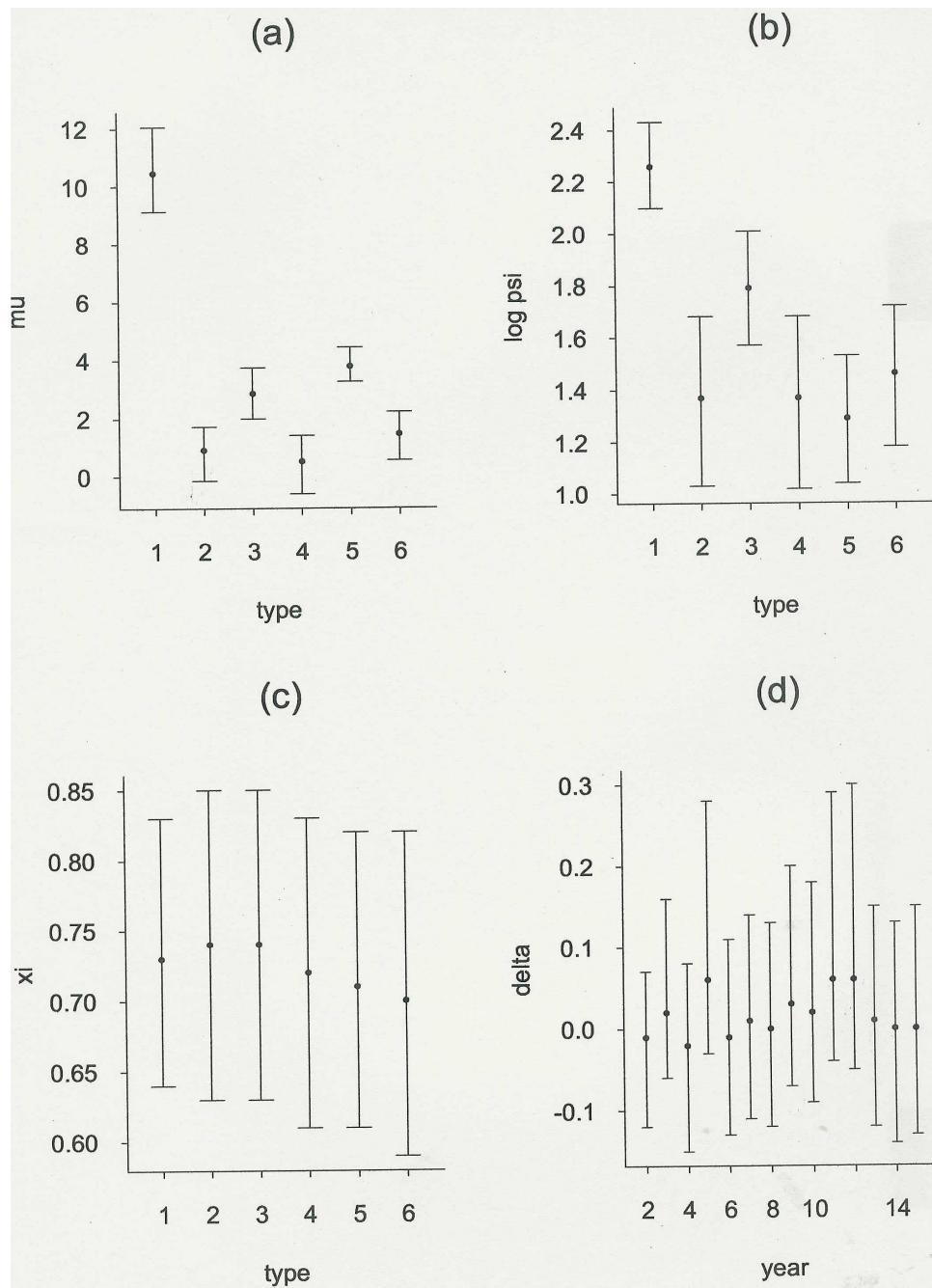
Level III. Conditional on Level II, for each  $j \in \{1, \dots, J\}$ , the point process of exceedances of type  $j$  is generated from the Poisson process with parameters  $\mu_j$ ,  $\psi_j$ ,  $\xi_j$ .

This model may be further extended to include a year effect, as follows. Suppose the extreme value parameters for type  $j$  in year  $k$  are not  $\mu_j, \psi_j, \xi_j$  but  $\mu_j + \delta_k, \psi_j, \xi_j$ . We fix  $\delta_1 = 0$  to ensure identifiability, and let  $\{\delta_k, k > 1\}$  follow an AR(1) process:

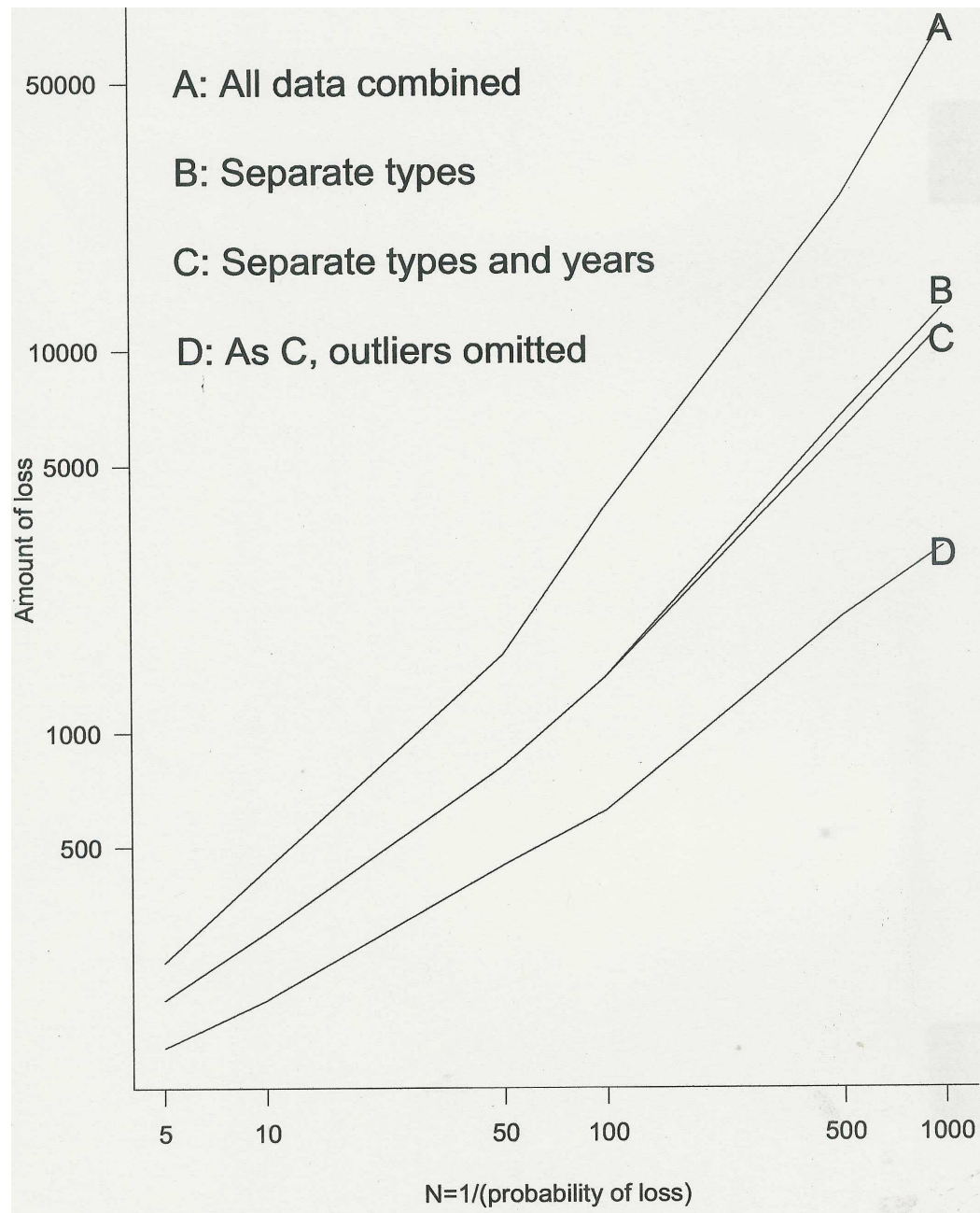
$$\delta_k = \rho\delta_{k-1} + \eta_k, \quad \eta_k \sim N(0, s_\eta^2)$$

with a vague prior on  $(\rho, s_\eta^2)$ .

We show boxplots for each of  $\mu_j, \log \psi_j, \xi_j, j = 1, \dots, 6$  and for  $\delta_k, k = 2, 15$ .



Posterior means and quartiles for  $\mu_j$ ,  $\log \psi_j$ ,  $\xi_j$  ( $j = 1, \dots, 6$ ) and for  $\delta_k$  ( $k = 2, \dots, 15$ ).



Computations of posterior predictive distribution functions (plotted on a log-log scale) corresponding to the homogenous model (curve A) and three different versions of the hierarchical model.



## IV. INSURANCE EXTREMES II

This example is based on a data set constructed by the U.K. insurance company Benfield-Greig, consisting of 57 historical storm events and losses calculated “as if” they occurred in 1998. Here,  $\mu_t$  is modelled as a function of time  $t$  through various co-variates —

Seasonal effects (dominant annual cycle)

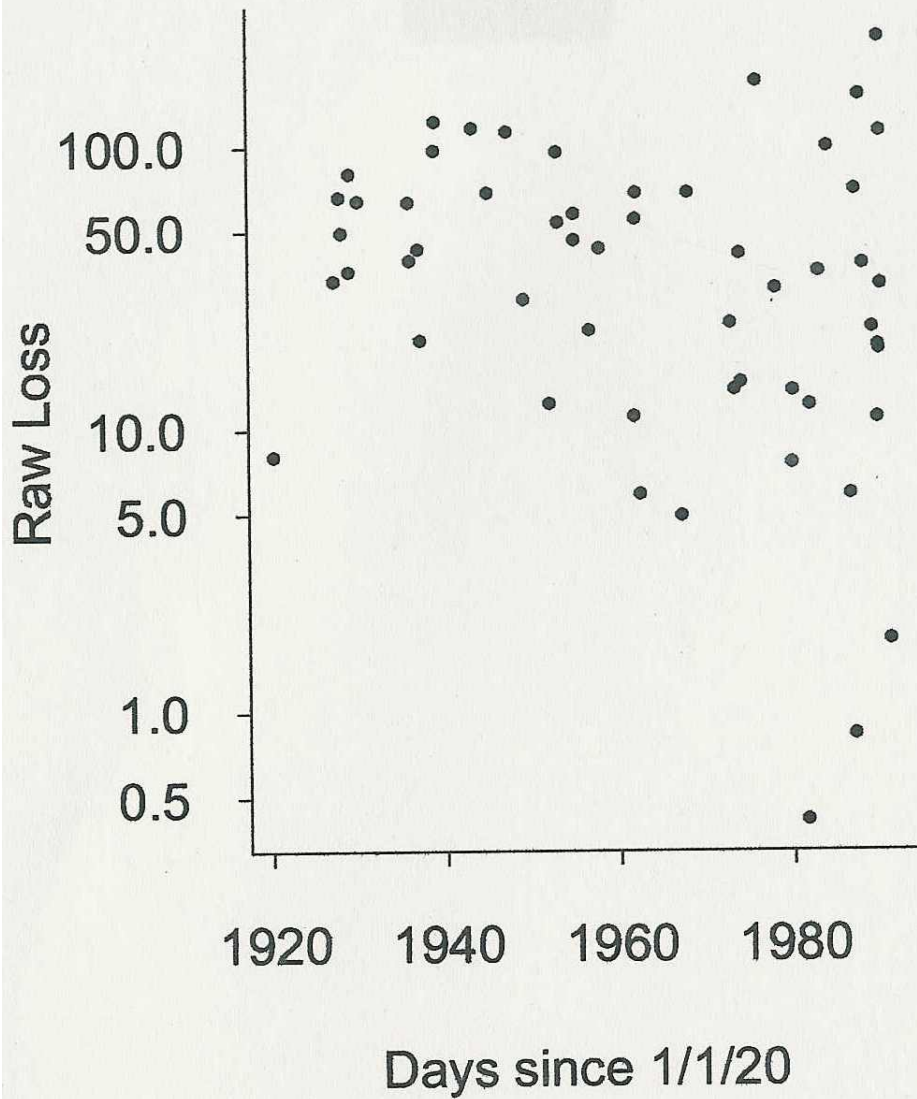
Polynomial terms in  $t$

Nonparametric trends

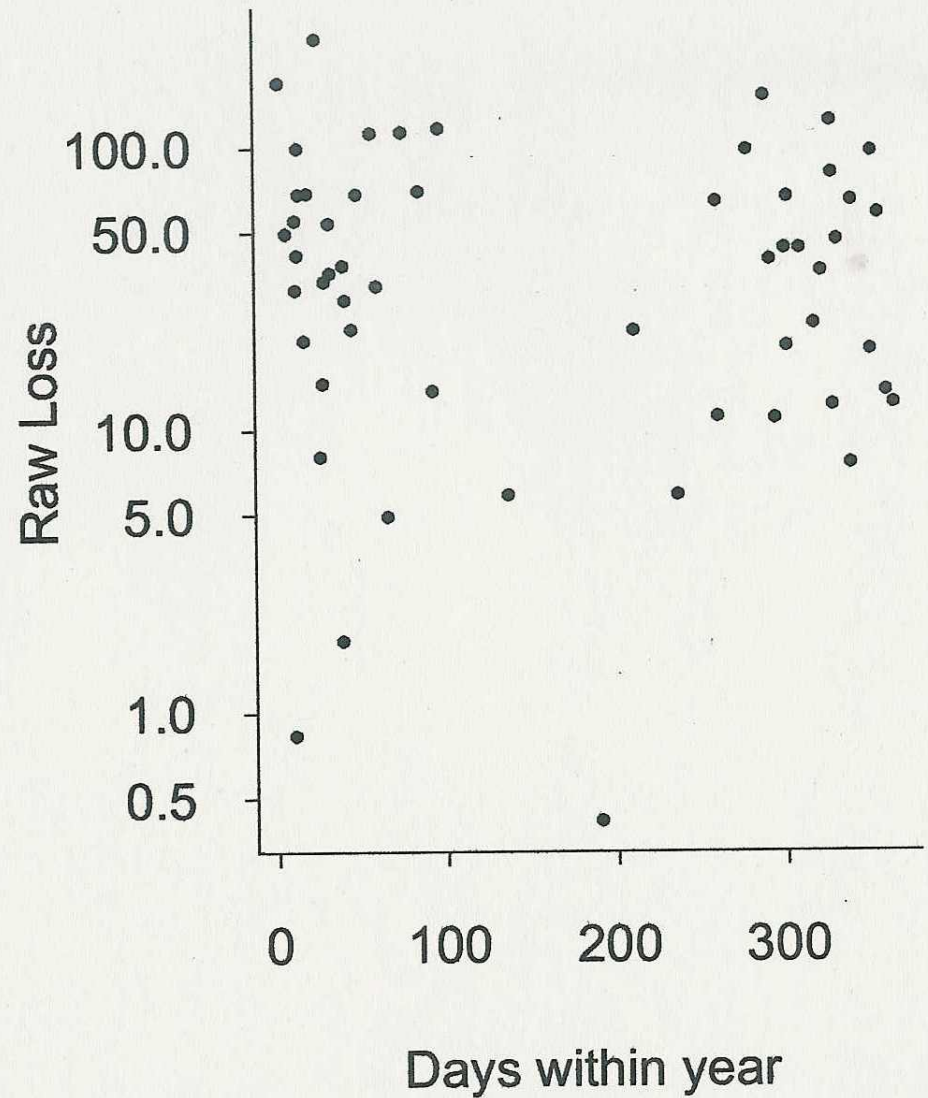
Dependence on oscillation indices

Other models in which  $\psi_t$  and  $\xi_t$  depend on  $t$  were also tried but do not produce significant differences from the constant case.

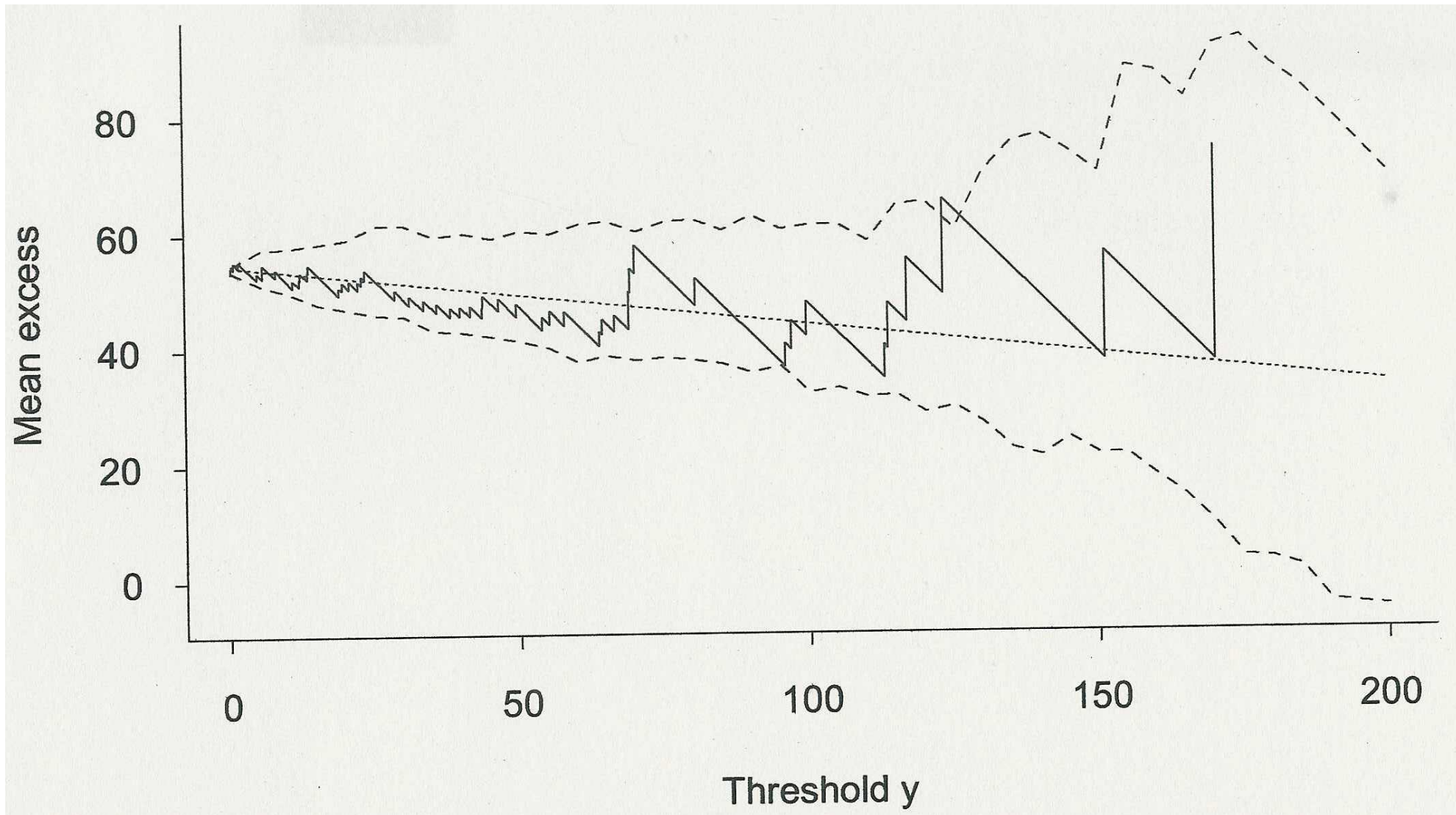
(a)



(b)



Plots of estimated storm losses against (a) time measured in years, (b) day within the year.

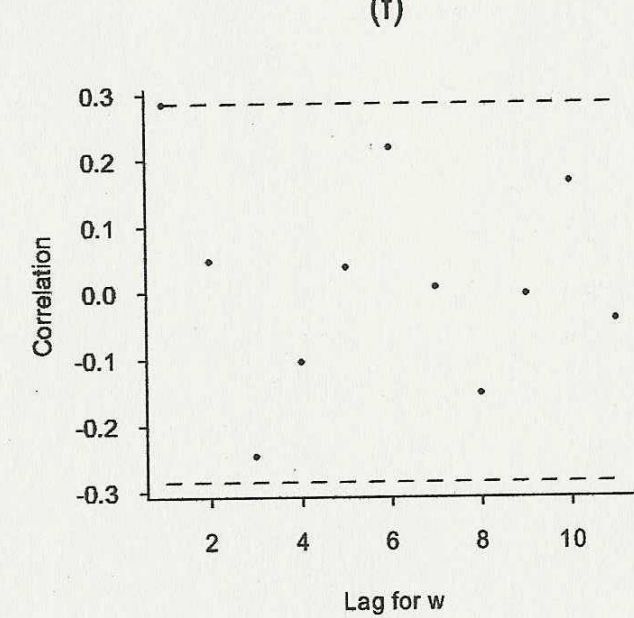
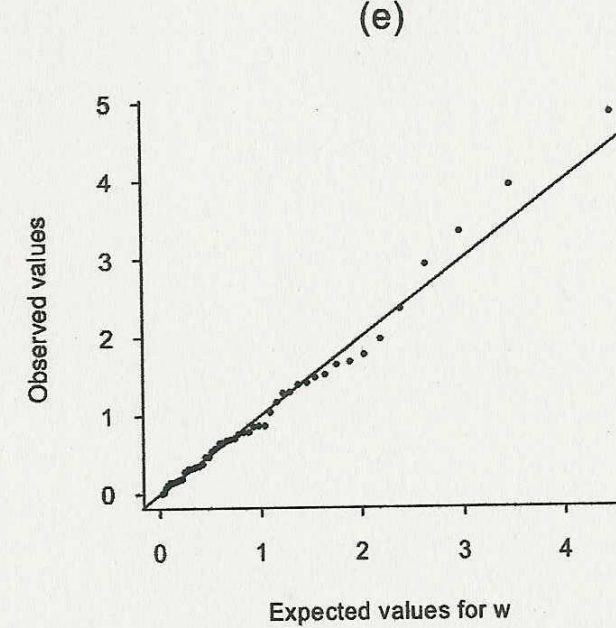
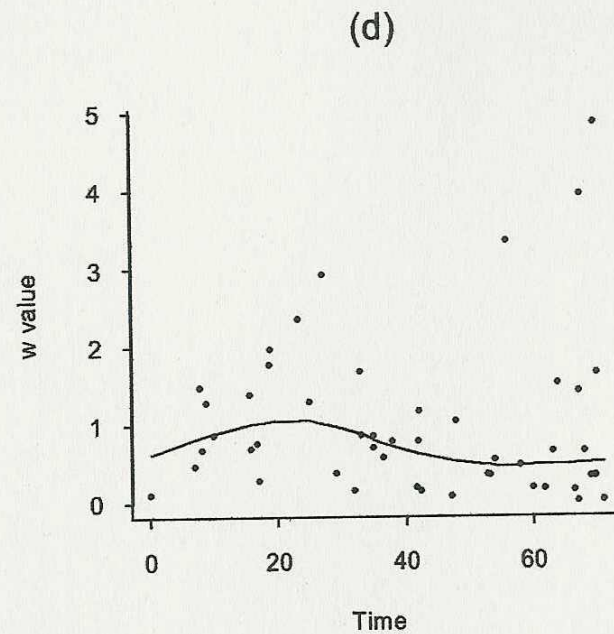
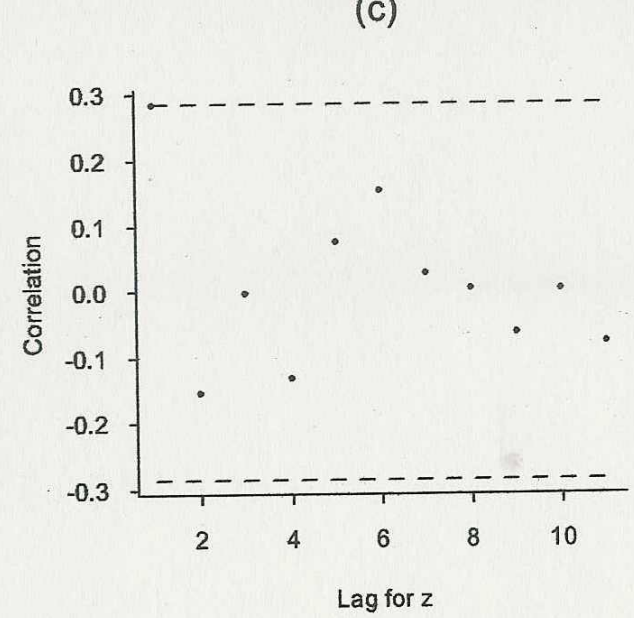
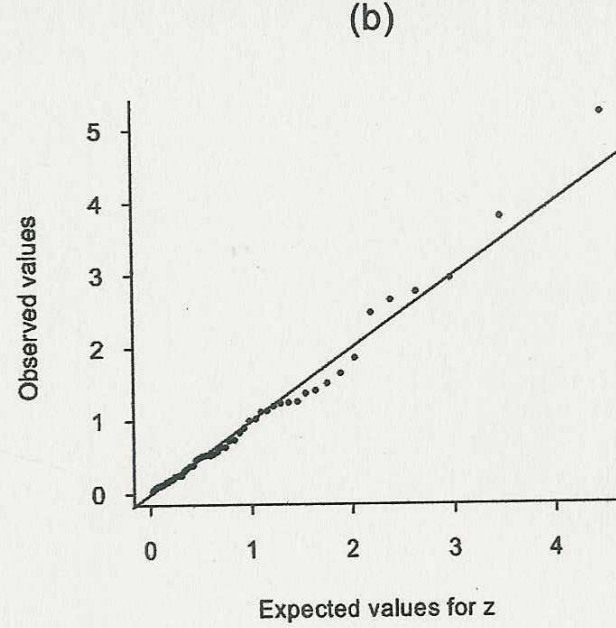
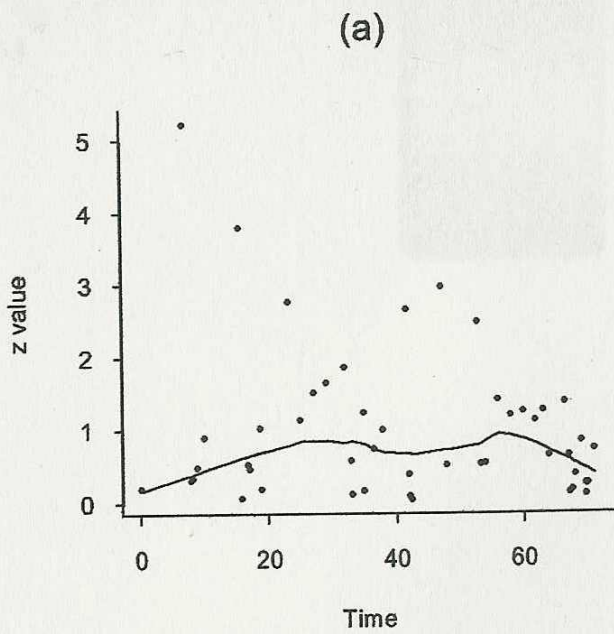


Mean excess plot with Monte Carlo confidence bands.

## Comparison of models

Model	Number of parameters	NLLH	AIC
Simple GPD	3	312.5	631.0
Seasonal	5	294.8	599.6
Seasonal+cubic	8	289.3	594.6
Seasonal+spline	10	289.6	599.2
Seasonal+SOI	6	294.5	601.0
Seasonal+NAO	6	289.0	590.0
Sea+NAO+cub	9	284.4	586.8
Sea+NAO+spl	11	284.6	591.2

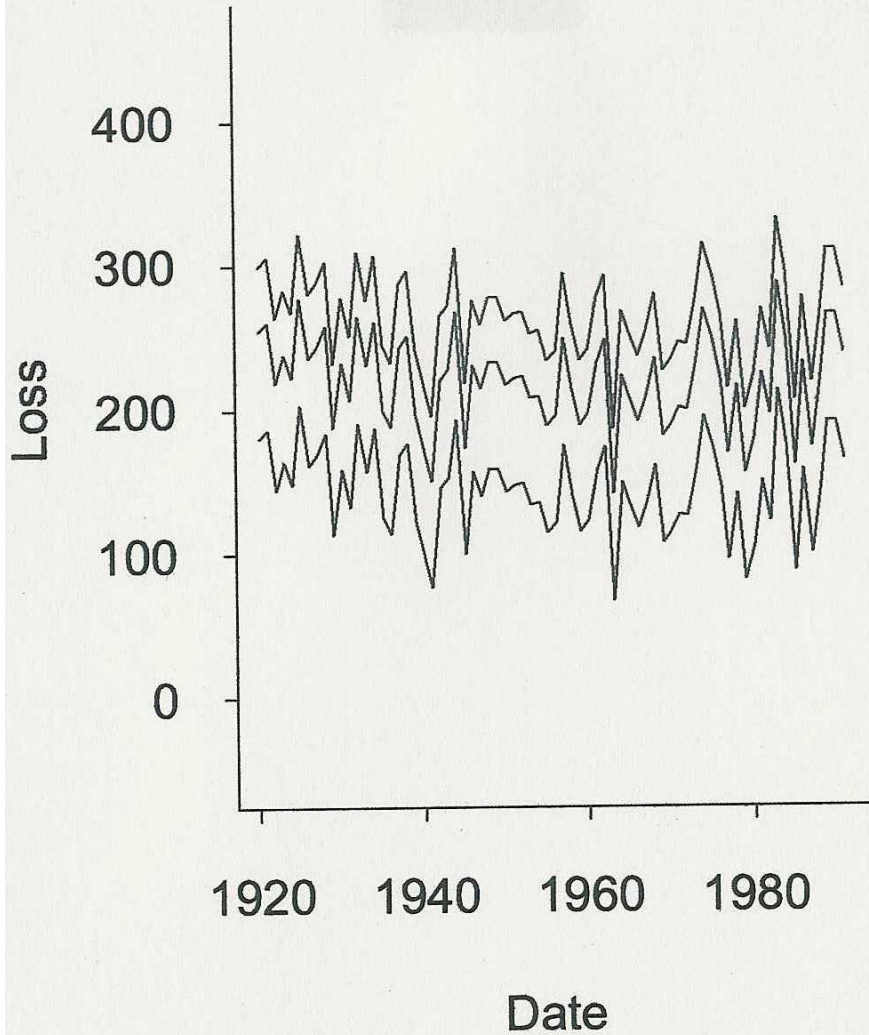
$$(AIC = 2NLLH + 2p)$$



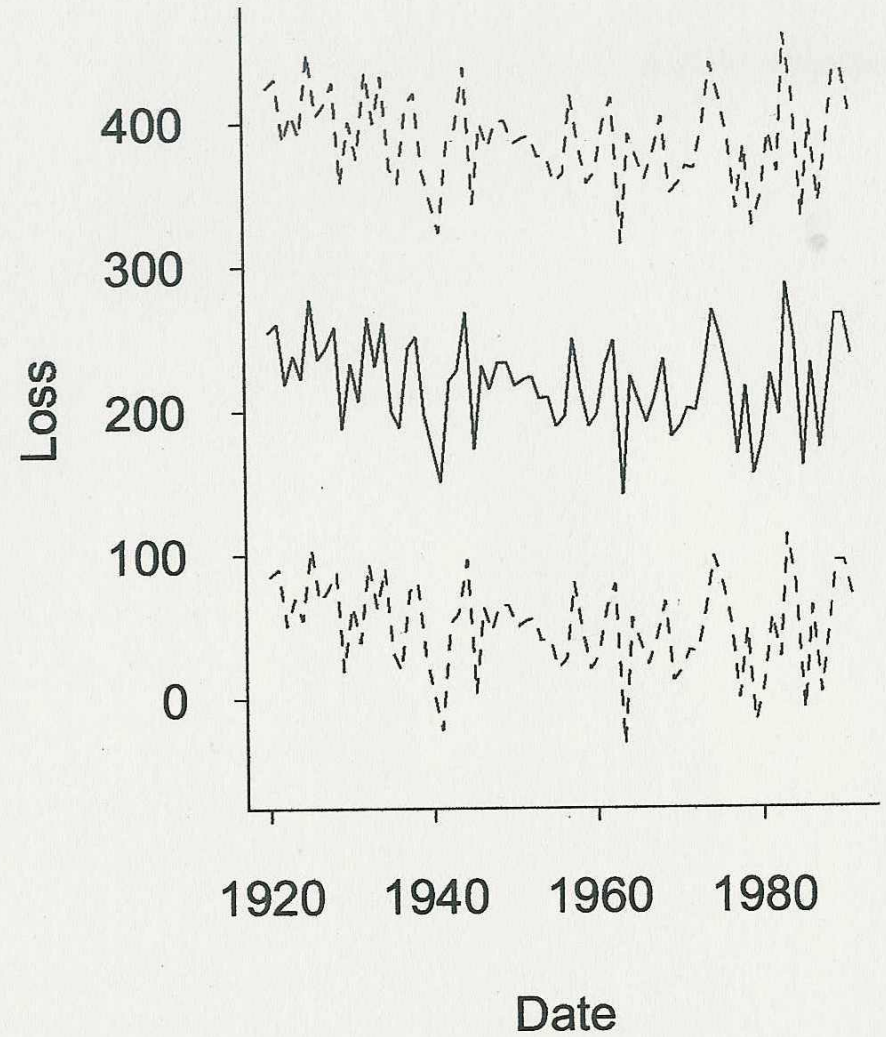
Diagnostic plots based on the  $Z$  (plots a,b,c) and  $W$  (plots d,e,f) statistics, seasonal model with NAO.



(a)

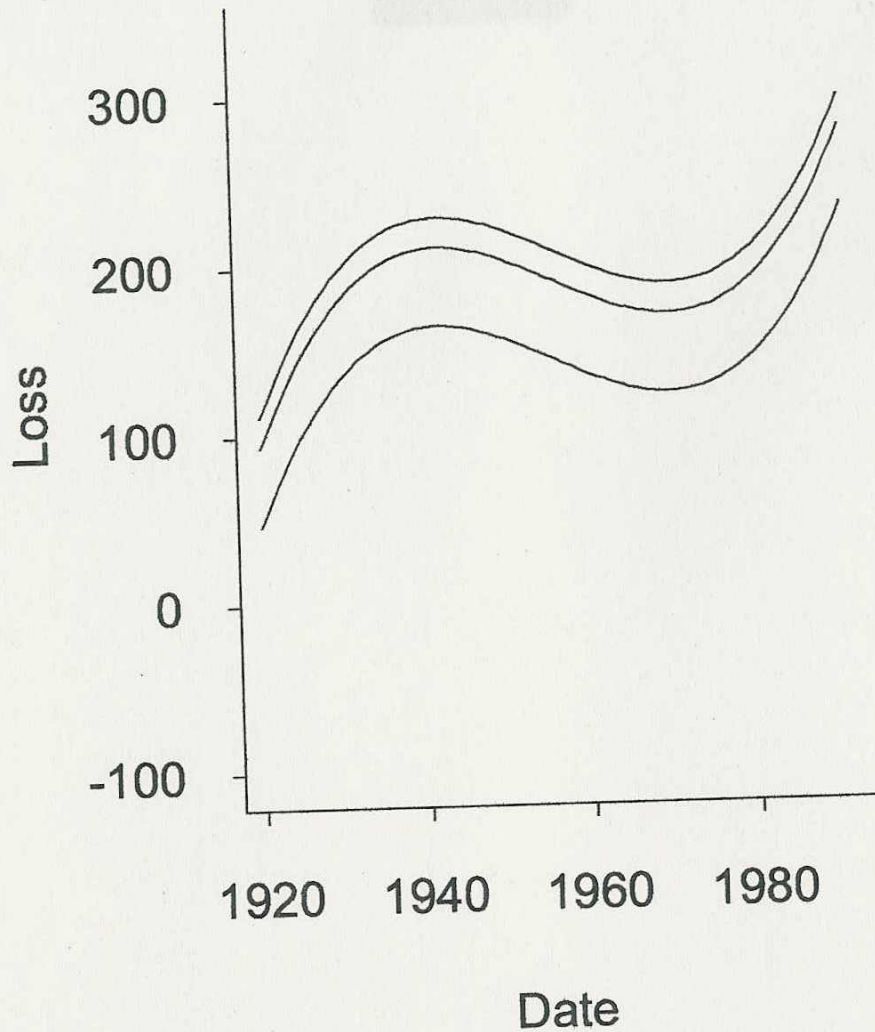


(b)

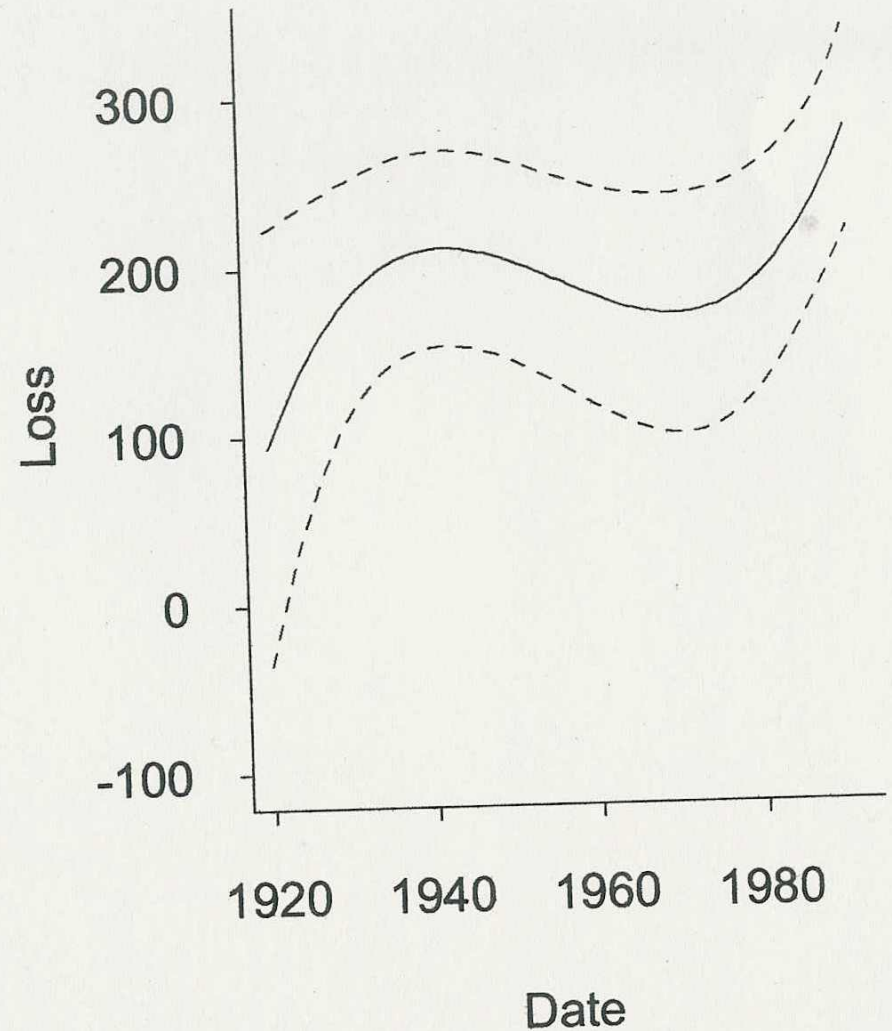


(a) Estimates of 10-year (bottom), 100-year (middle) and 1000-year (top) return levels based on the fitted model for January, assuming long-term trend based on NAO. (b) 100-year return level with confidence limits.

(a)

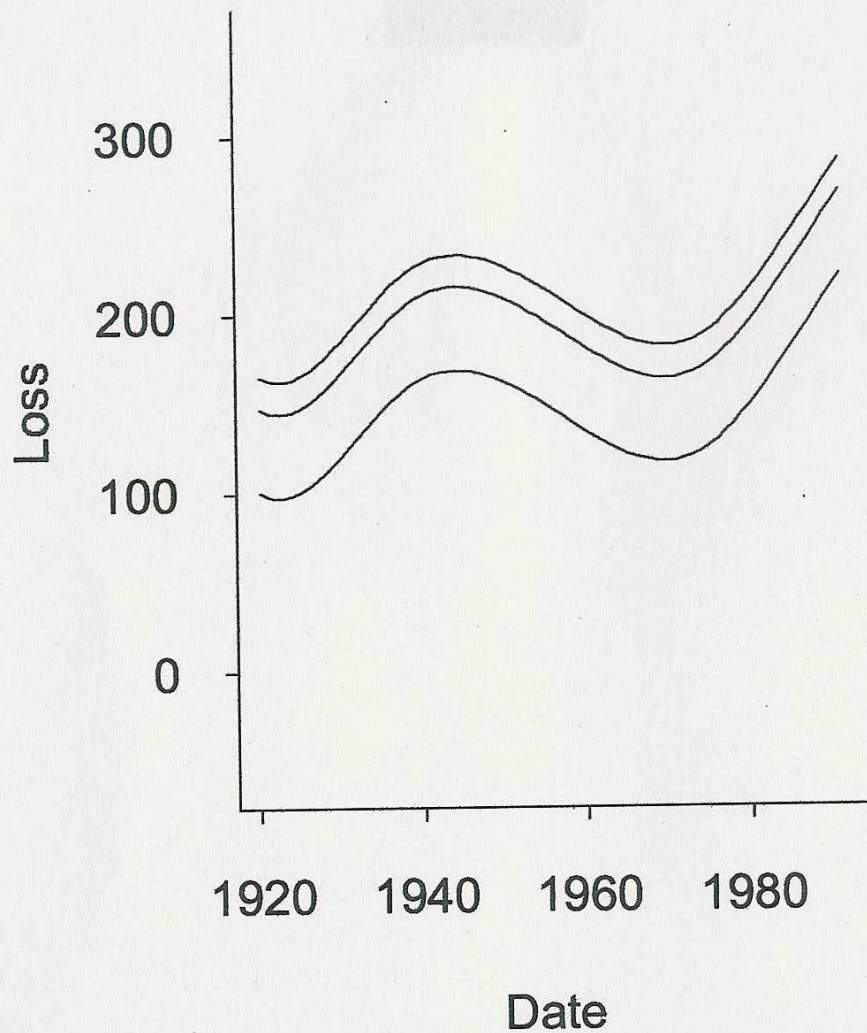


(b)

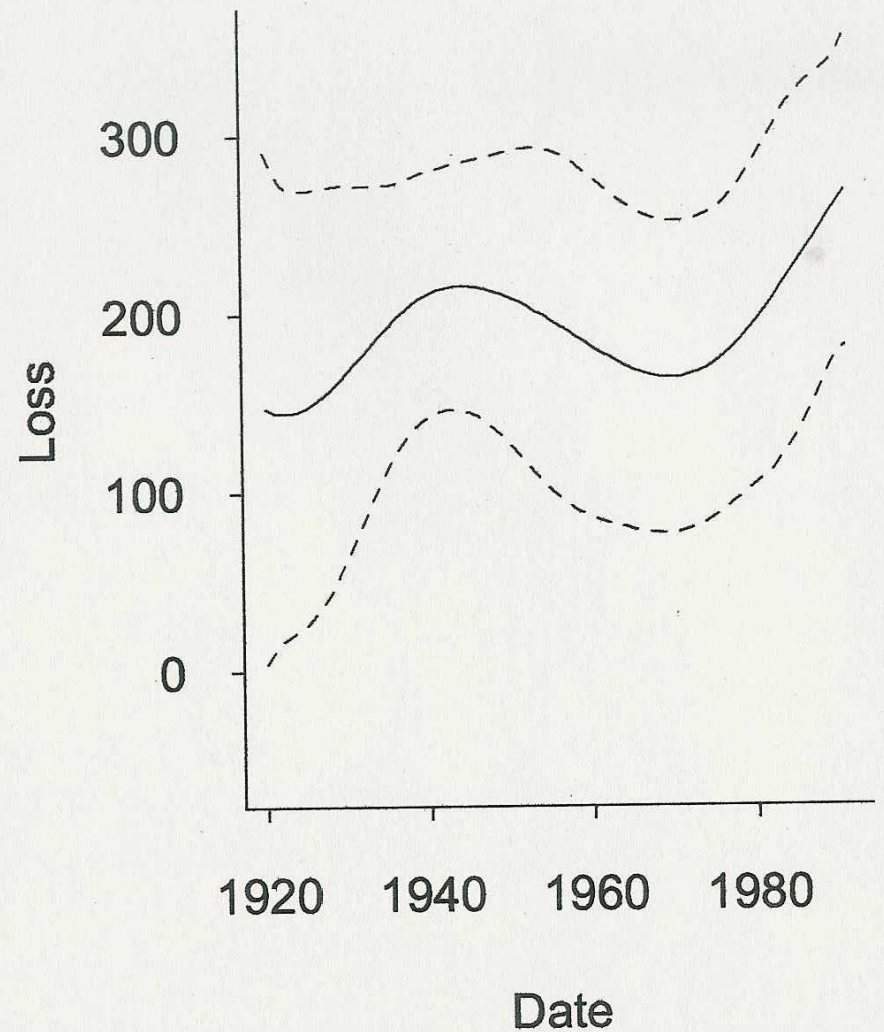


(a) Estimates of 10-year (bottom), 100-year (middle) and 1000-year (top) return levels based on the fitted model for January, assuming long-term trend based on a cubic polynomial. (b) 100-year return level with confidence limits.

(a)



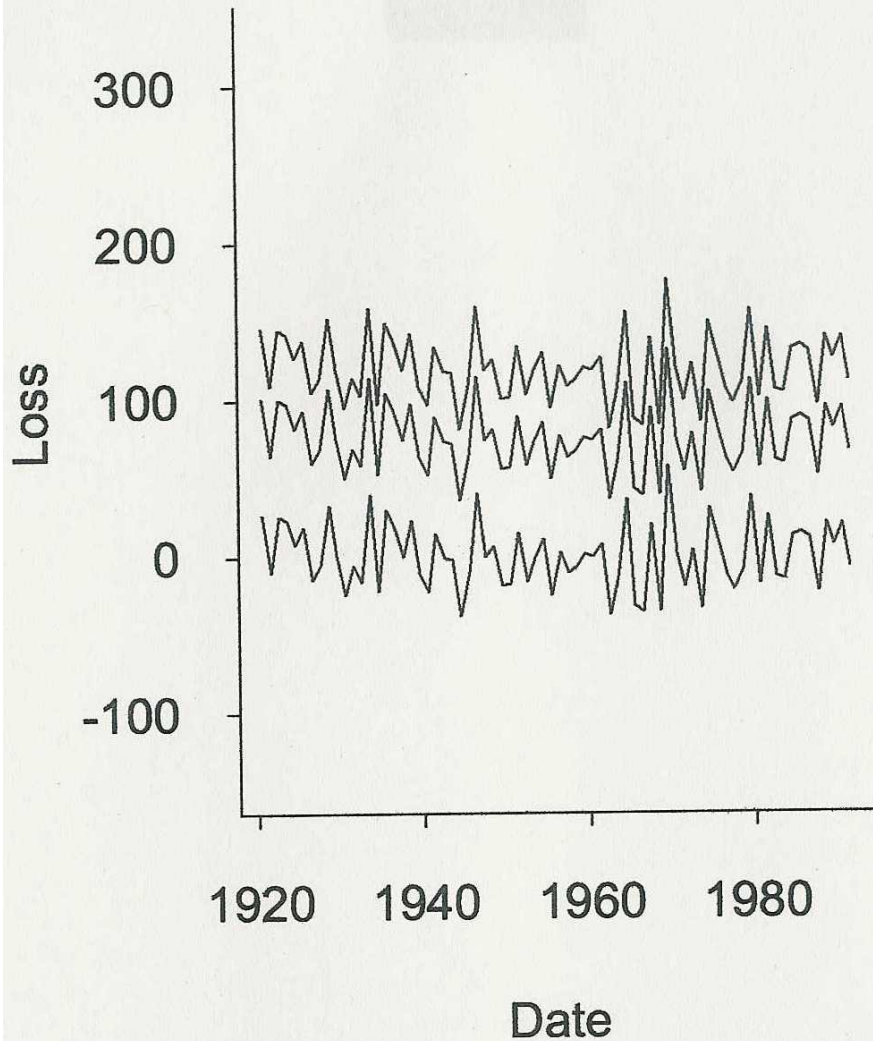
(b)



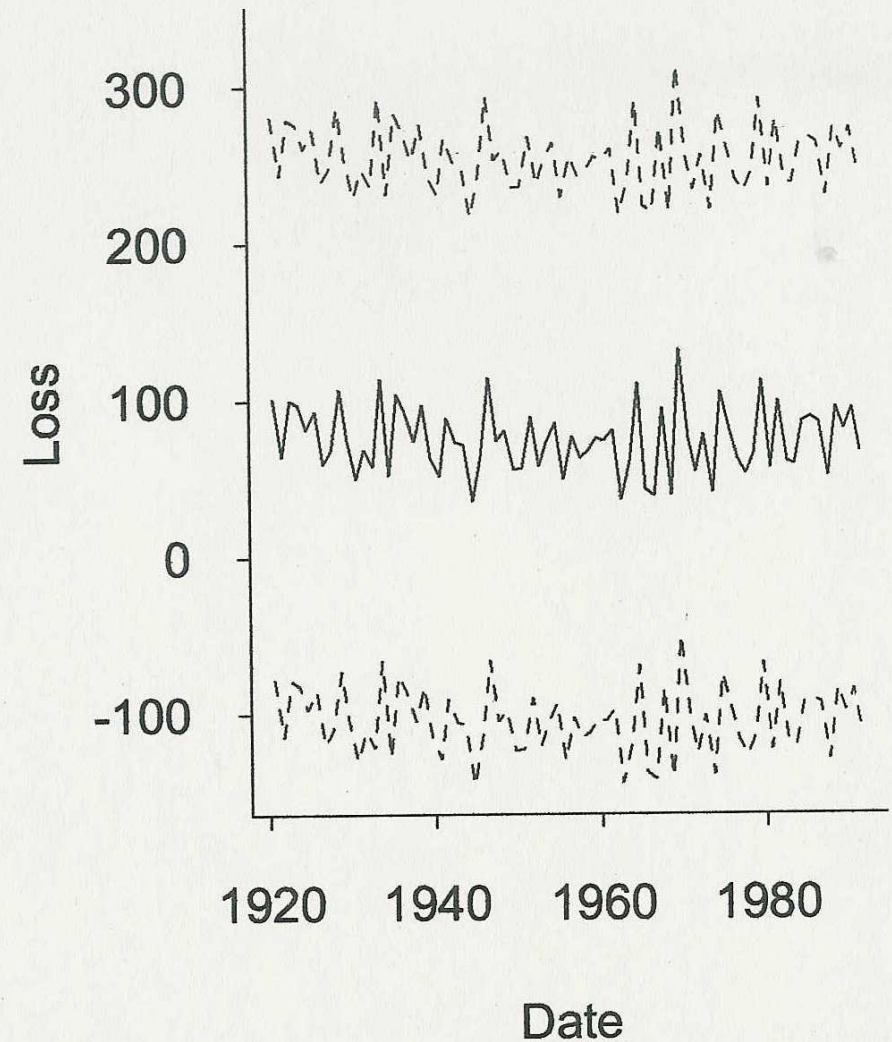
(a) Estimates of 10-year (bottom), 100-year (middle) and 1000-year (top) return levels based on the fitted model for January, assuming long-term trend based on a cubic spline with 5 knots.  
(b) 100-year return level with confidence limits.



(a)



(b)



(a) Estimates of 10-year (bottom), 100-year (middle) and 1000-year (top) return levels based on the fitted model for July, assuming long-term trend based on NAO. (b) 100-year return level with confidence limits.

## How extreme was the 1990 Event?

Model	Return Period (Years)
Pareto	18.7
Exponential	487
Gen. Pareto	333
With seasons	348
High NAO	187
Low NAO	1106
Random NAO	432
Current industry estimate?	<50?

## Note:

One reason I wanted to bring up this example is because a dataset that is constructed along very similar lines, but for U.S. hurricane damage, has recently been published in

R.A. Pielke Jr. *et al.* (2007), Normalized Hurricane Damages in the United States: 1900-2005. *Natural Hazards Review*, (accepted).

(updating an earlier paper by Pielke and Landsea, 1998). They claim, in essence, that there is no evidence of any trend in normalized hurricane damage.

The data are available from Roger Pielke's webpage and it would be interesting to apply the same method as we have described here!

# V. TREND IN PRECIPITATION EXTREMES

(joint work with Amy Grady and Gabi Hegerl)

During the past decade, there has been extensive research by climatologists documenting increases in the levels of extreme precipitation, but in observational and model-generated data.

With a few exceptions (papers by Katz, Zwiers and co-authors) this literature have not made use of the extreme value distributions and related constructs

There are however a few papers by statisticians that have explored the possibility of using more advanced extreme value methods (e.g. Cooley, Naveau and Nychka, to appear *JASA*; Sang and Gelfand, submitted)

This discussion uses extreme value methodology to look for trends

## DATA SOURCES

- NCDC Rain Gauge Data (Groisman 2000)
  - Daily precipitation from 5873 stations
  - Select 1970–1999 as period of study
  - 90% data coverage provision — 4939 stations meet that
- NCAR-CCSM climate model runs
  - 20 × 41 grid cells of side 1.4°
  - 1970–1999 and 2070–2099 (A2 scenario)
- PRISM data
  - 1405 × 621 grid, side 4km
  - Elevations
  - Mean annual precipitation 1970–1997

## **EXTREME VALUES METHODOLOGY**

Based on “point process” extreme values methodology (cf. Smith 1989, Coles 2001, Smith 2003)

## Inhomogeneous case:

- Time-dependent threshold  $u_t$  and parameters  $\mu_t, \psi_t, \xi_t$

- Exceedance  $y > u_t$  at time  $t$  has probability

$$\frac{1}{\psi_t} \left( 1 + \xi_t \frac{y - \mu_t}{\psi_t} \right)_+^{-1/\xi_t - 1} \exp \left\{ - \left( 1 + \xi_t \frac{u_t - \mu_t}{\psi_t} \right)_+^{-1/\xi_t} \right\} dy dt$$

- Estimation by maximum likelihood

## Seasonal models without trends

General structure:

$$\begin{aligned}\mu_t &= \theta_{1,1} + \sum_{k=1}^{K_1} \left( \theta_{1,2k} \cos \frac{2\pi kt}{365.25} + \theta_{1,2k+1} \sin \frac{2\pi kt}{365.25} \right), \\ \log \psi_t &= \theta_{2,1} + \sum_{k=1}^{K_2} \left( \theta_{2,2k} \cos \frac{2\pi kt}{365.25} + \theta_{2,2k+1} \sin \frac{2\pi kt}{365.25} \right), \\ \xi_t &= \theta_{3,1} + \sum_{k=1}^{K_3} \left( \theta_{3,2k} \cos \frac{2\pi kt}{365.25} + \theta_{3,2k+1} \sin \frac{2\pi kt}{365.25} \right).\end{aligned}$$

Call this the  $(K_1, K_2, K_3)$  model.

*Note:* This is all for one station. The  $\theta$  parameters will differ at each station.



## Models with trend

Add to the above:

- Overall linear trend  $\theta_{j,2K+2}t$  added to any of  $\mu_t$  ( $j = 1$ ),  $\log \psi_t$  ( $j = 1$ ),  $\xi_t$  ( $j = 1$ ). Define  $K_j^*$  to be 1 if this term is included, o.w. 0.
- Interaction terms of form

$$t \cos \frac{2\pi kt}{365.25}, \quad t \sin \frac{2\pi kt}{365.25}, \quad k = 1, \dots, K_j^{**}.$$

Typical model denoted

$$(K_1, K_2, K_3) \times (K_1^*, K_2^*, K_3^*) \times (K_1^{**}, K_2^{**}, K_3^{**})$$

Eventually use  $(4, 2, 1) \times (1, 1, 0) \times (2, 2, 0)$  model (27 parameters for each station)

## SPATIAL SMOOTHING

Let  $Z_s$  be field of interest, indexed by  $s$  (typically the logarithm of the 25-year RV at site  $s$ , or a log of ratio of RVs. Taking logs improves fit of spatial model, to follow.)

Don't observe  $Z_s$  — estimate  $\hat{Z}_s$ . Assume

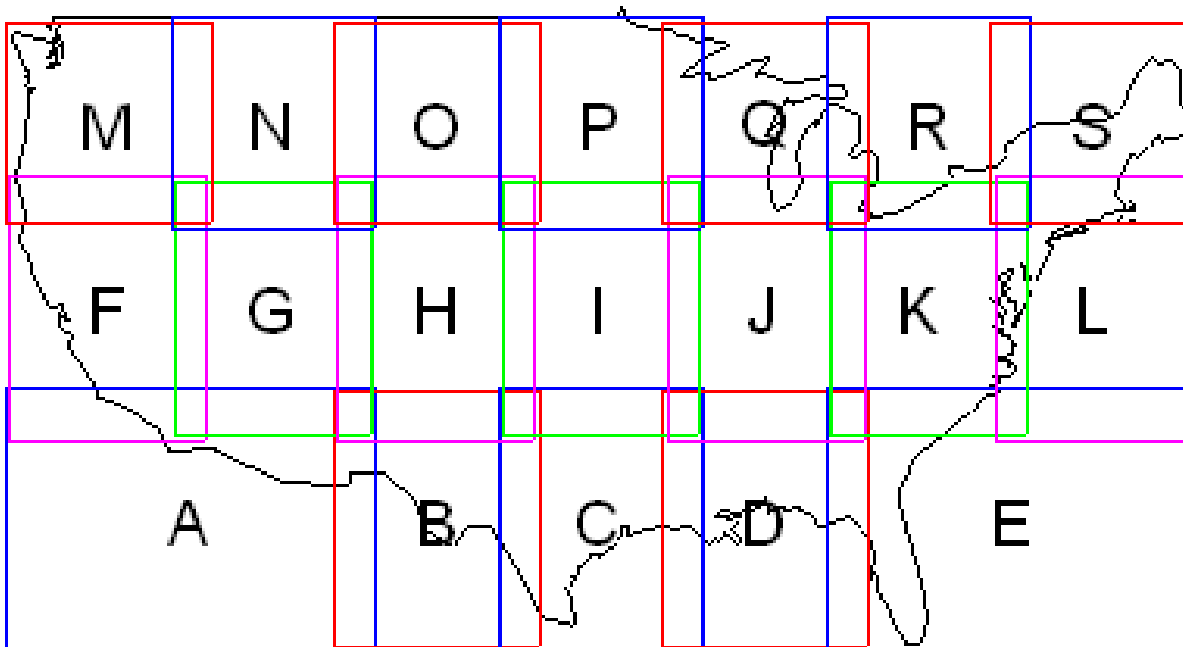
$$\begin{aligned}\hat{Z} | Z &\sim N[Z, W] \\ Z &\sim N[X\beta, V(\phi)] \\ \hat{Z} &\sim N[X\beta, V(\phi) + W].\end{aligned}$$

for known  $W$ ;  $X$  are covariates,  $\beta$  are unknown regression parameters and  $\phi$  are parameters of spatial covariance matrix  $V$ .

- $\phi$  by REML
- $\beta$  given  $\phi$  by GLS
- Predict  $Z$  at observed and unobserved sites by kriging

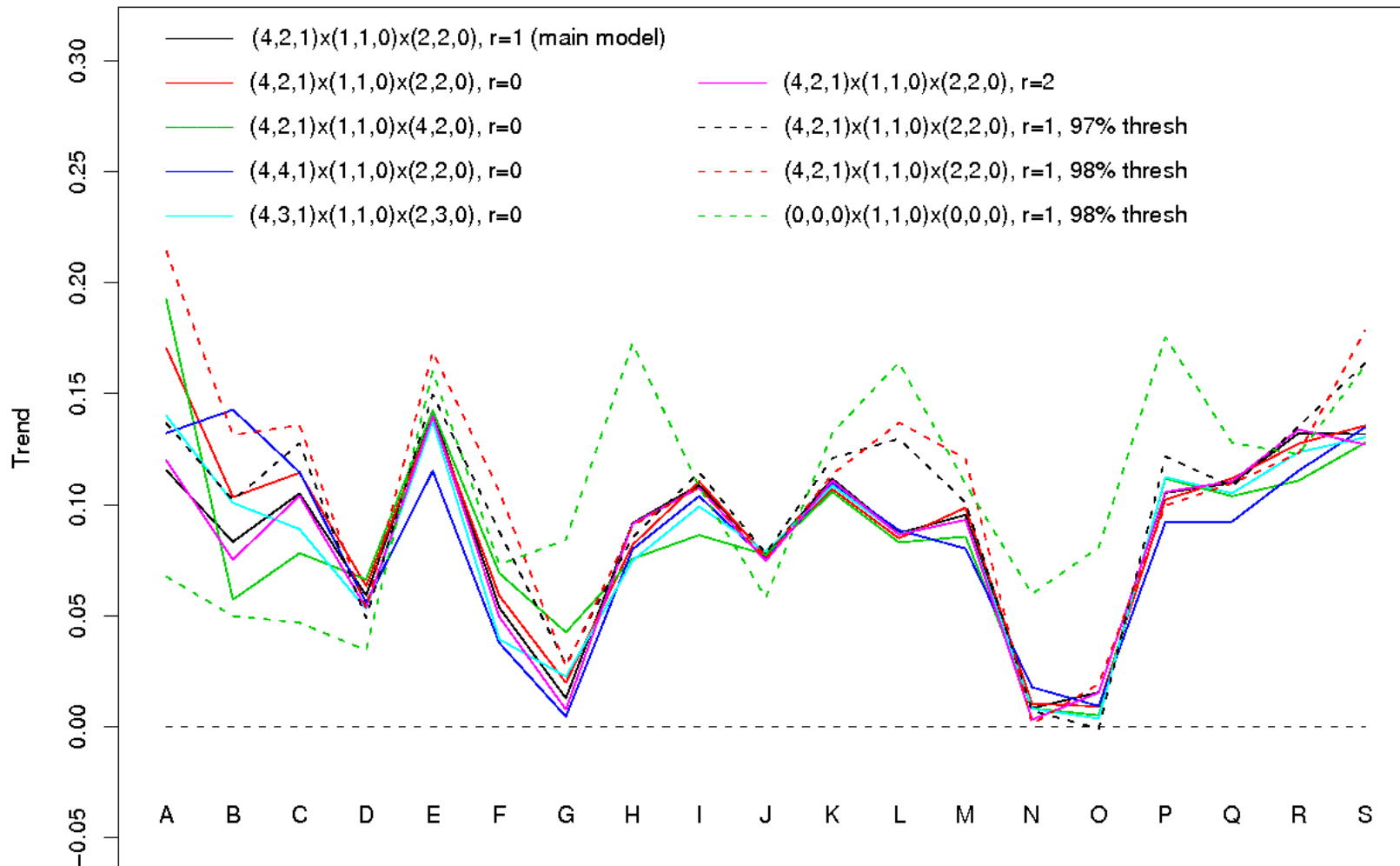
## *Spatial Heterogeneity*

- Divide US into 19 overlapping regions, most  $10^{\circ} \times 10^{\circ}$ 
  - Kriging within each region
  - Linear smoothing across region boundaries
  - Same for MSPEs
  - Also calculate regional averages, including MSPE

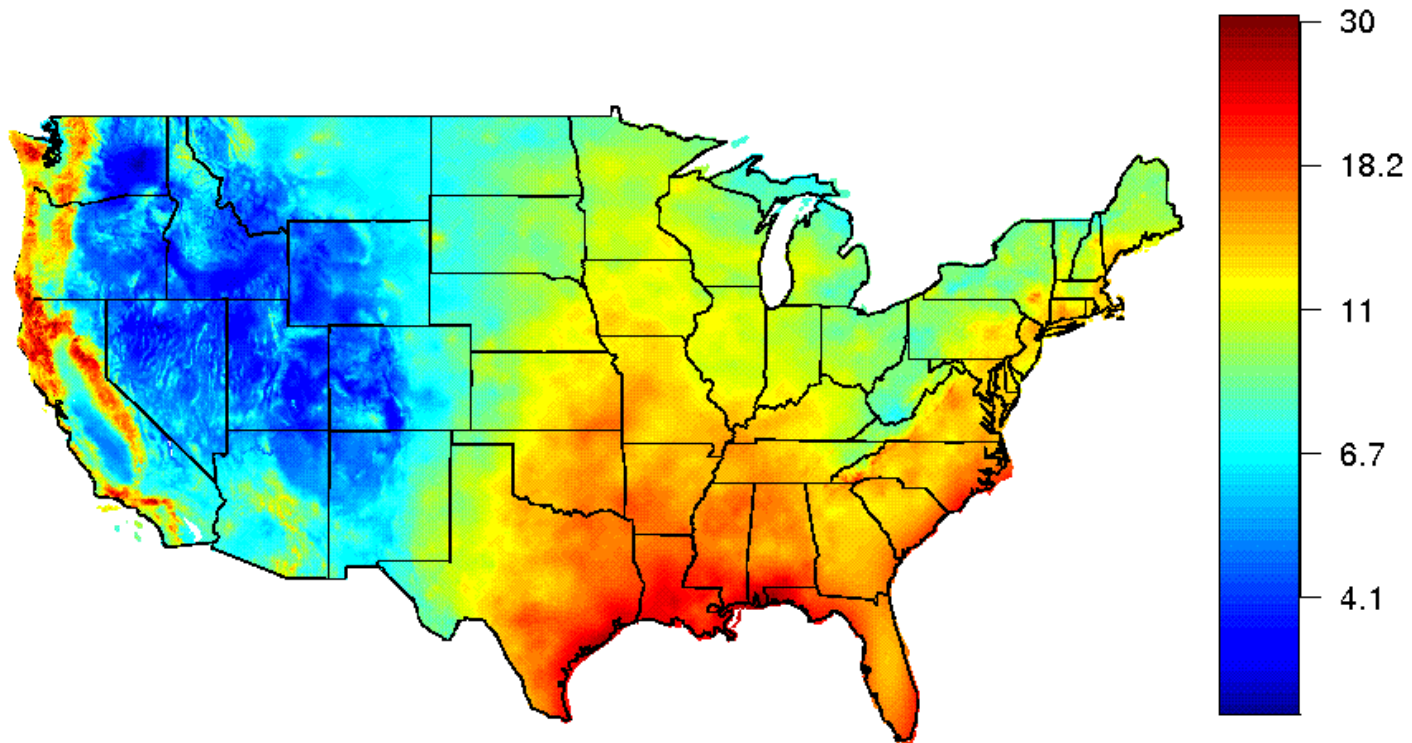


Continental USA divided into 19 regions

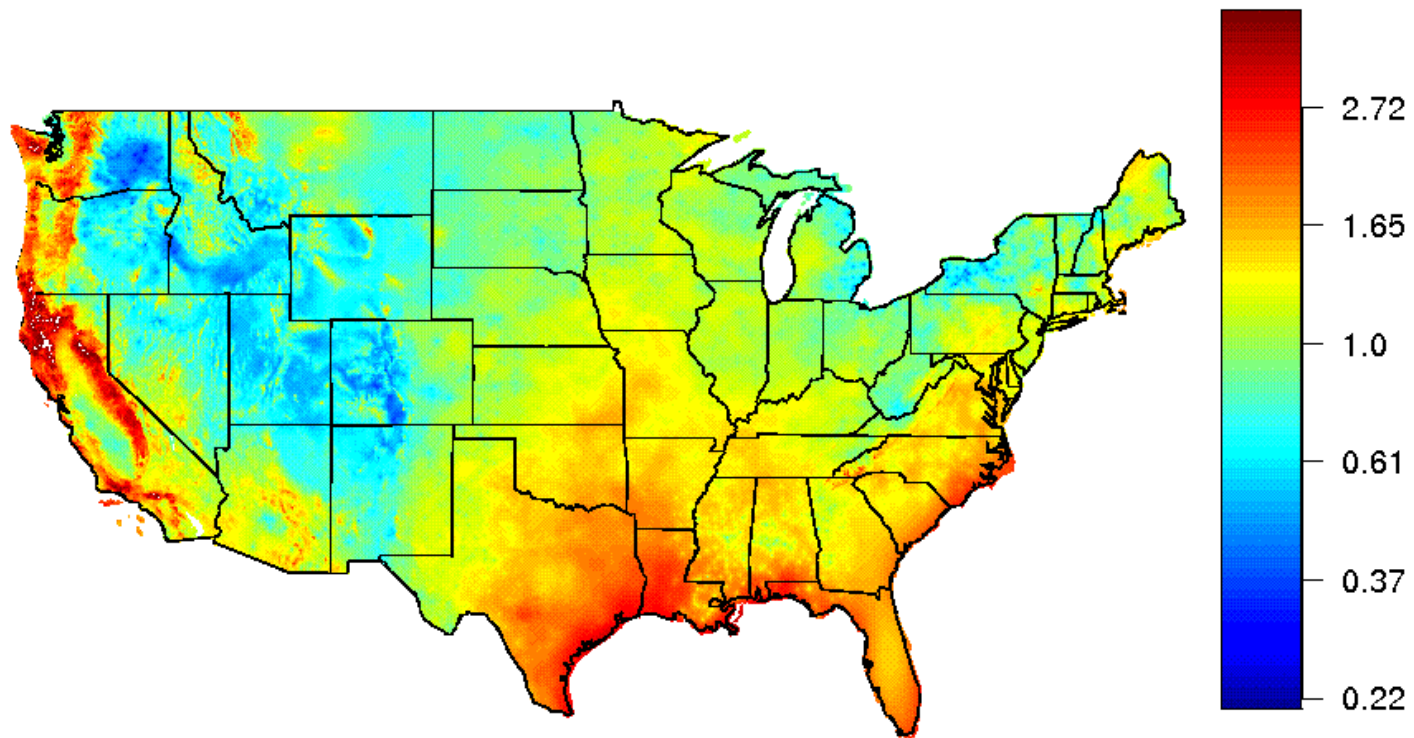
### REGIONAL AVERAGE TRENDS FOR 9 EV MODELS (GWA METHOD)



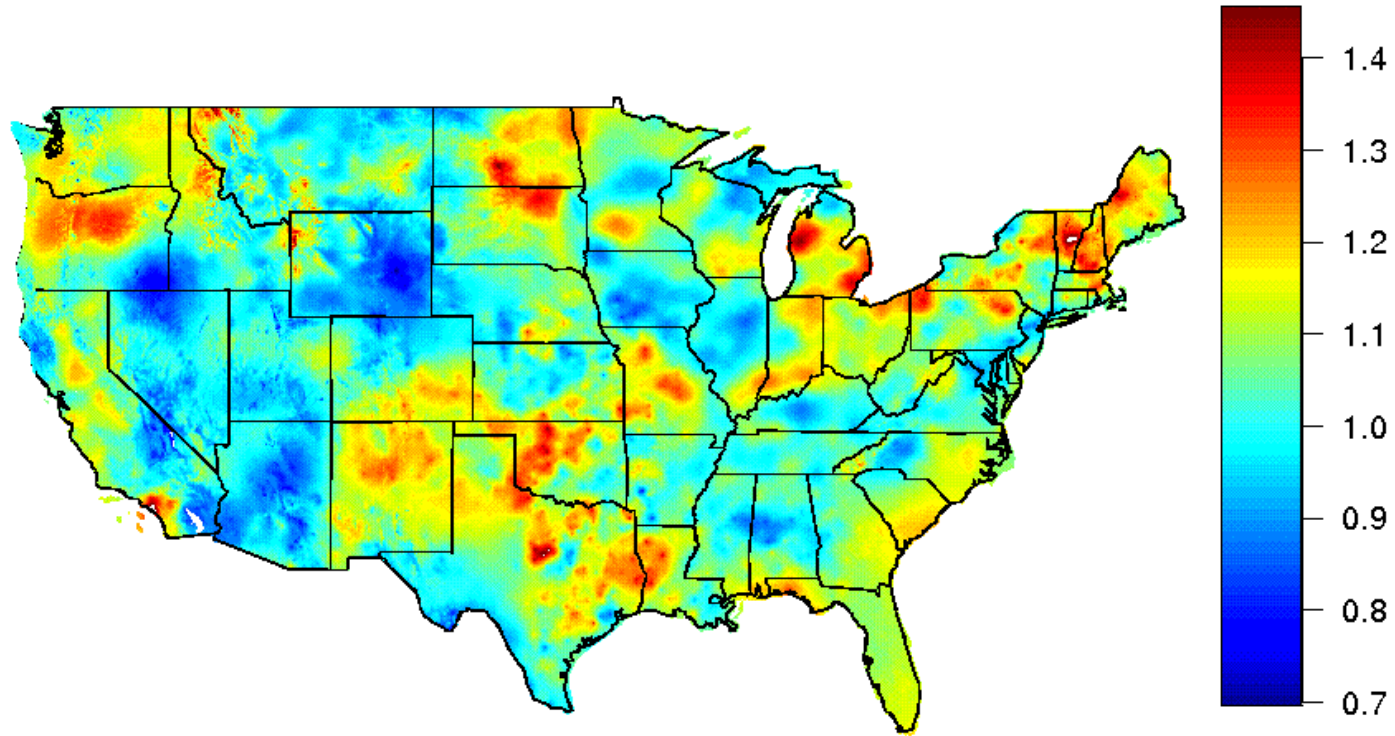
Trends across 19 regions (measured as change in log RV25) for 8 different seasonal models and one non-seasonal model with simple linear trends. Regional averaged trends by geometric weighted average approach.



Map of 25-year return values (cm.) for the years 1970–1999



Root mean square prediction errors for map of 25-year return values for 1970–1999



Ratios of return values in 1999 to those in 1970



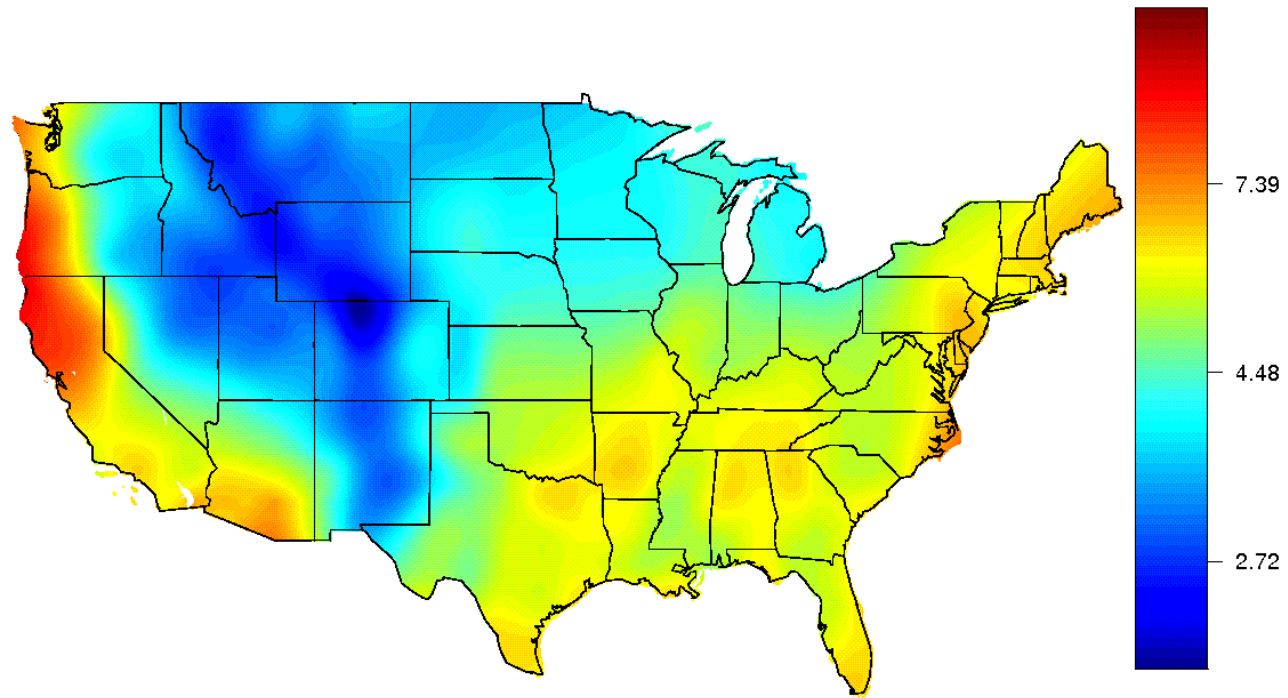
	$\Delta_1$	$S_1$	$\Delta_2$	$S_2$		$\Delta_1$	$S_1$	$\Delta_2$	$S_2$
A	-0.01	.03	0.05**	.05	K	0.08***	.01	0.09**	.03
B	0.07**	.03	0.08***	.04	L	0.07***	.02	0.07*	.04
C	0.11***	.01	0.10	.03	M	0.07***	.02	0.10**	.03
D	0.05***	.01	0.06	.05	N	0.02	.03	0.01	.03
E	0.13***	.02	0.14*	.05	O	0.01	.02	0.02	.03
F	0.00	.02	0.05*	.04	P	0.07***	.01	0.11***	.03
G	-0.01	.02	0.01	.03	Q	0.07***	.01	0.11***	.03
H	0.08***	.01	0.10***	.03	R	0.15***	.02	0.13***	.03
I	0.07***	.01	0.12***	.03	S	0.14***	.02	0.12*	.06
J	0.05***	.01	0.08**	.03					

$\Delta_1$ : Mean change in log 25-year return value (1970 to 1999) by kriging

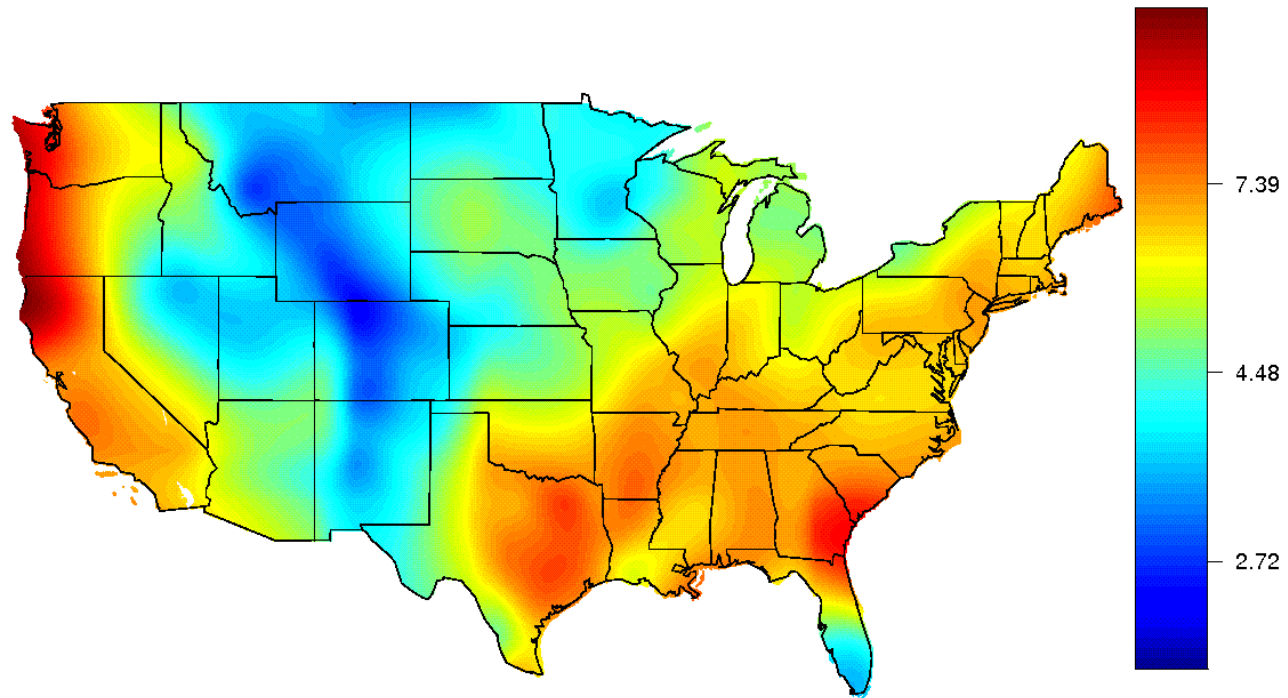
$S_1$ : Corresponding standard error (or RMSPE)

$\Delta_2$ ,  $S_2$ : same but using geometrically weighted average (GWA)

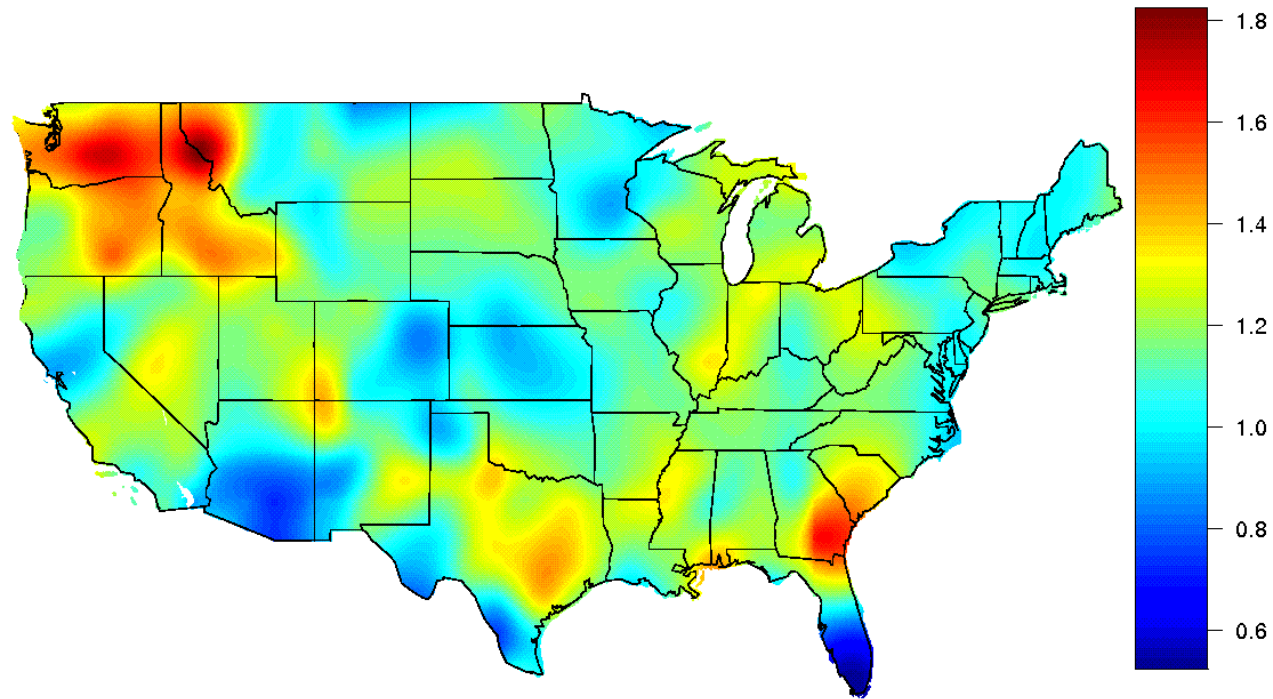
Stars indicate significance at 5%\*, 1%\*\* , 0.1%\*\*\*.



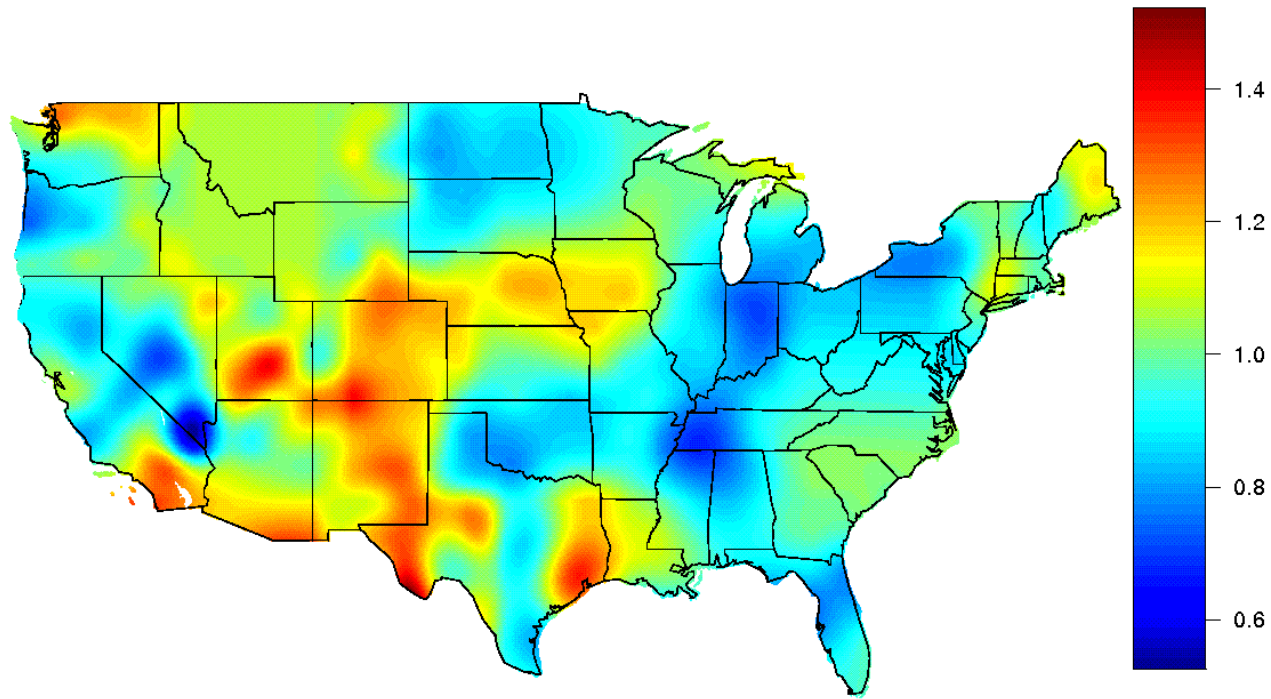
Return value map for CCSM data (cm.): 1970–1999



Return value map for CCSM data (cm.): 2070–2099



Estimated ratios of 25-year return values for 2070–2099 to those of 1970–1999, based on CCSM data, A2 scenario



Extreme value model with trend: ratio of 25-year return value in 1999 to 25-year return value in 1970, based on CCSM data

## CONCLUSIONS

1. Focus on  $N$ -year return values — strong historical tradition for this measure of extremes (we took  $N = 25$  here)
2. Seasonal variation of extreme value parameters is a critical feature of this analysis
3. Overall significant increase over 1970–1999 except for parts of western states — average increase across continental US is 7%
4. Projections to 2070–2099 show further strong increases but note caveat based on point 5
5. *But...* based on CCSM data there is a completely different spatial pattern and no overall increase — still leaves some doubt as to overall interpretation.

**THANK YOU FOR YOUR  
ATTENTION!**