

# Predictive confidence distributions for time series



Gudmund Hermansen (with Nils Lid Hjort)

University of Oslo

May 13, 2015

# Motivation

- In several applications the main goal of fitting a statistical model is to predict future outcomes.
- The construction of confidence intervals, or credibility intervals, are therefore of great interest, i.e. confidence distributions.
- Here we will discuss a framework for confidence distributions for a future unobserved variable  $Y_0$  that we wish to predict.
- In addition, we discuss how to make confidence distributions for
  - the probability of exceeding a threshold and
  - the time point  $t$  corresponding to the event that  $E Y_t$  crosses  $y_0$ .
- The focus here is on various time series models.
- It is preliminary work and an extended discussion of ideas and methods presented in Schweder & Hjort (2015, Chapter 12)

# Confidence distributions for future observations

- Following the framework of Schweder & Hjort (2015) the CD

$$C_{\text{pred}}(y_0, Y),$$

for variable  $Y_0$  that we wish to predict based on data  $Y$ , has to satisfy two conditions:

- (i)  $C_{\text{pred}}(y_0, \mathbf{y}_{\text{obs}})$  is a distribution function in  $y_0$  for each  $\mathbf{y}_{\text{obs}}$  and
  - (ii)  $C_{\text{pred}}(Y_0, Y) \sim U[0, 1]$ , when  $(Y_0, Y)$  has the joint distribution of the true model.
- Note that this is different from estimating  $E Y_0$  or  $E Y_0 | \mathbf{y}_{\text{obs}}$ .
  - We will (try to) motivate an approximate pivot for linear predictors and apply the general CD framework in some applications.

# Prototype illustration with independent realisations

- Before we introduce dependency, consider the prototype case of predicting  $Y_{n+1}$  given a sequence of i.i.d. realisations.
- Suppose we wish to predict  $Y_{n+1}$  given a series of  $n$  i.i.d. realisations, where  $Y_1 \sim N(\mu, \sigma^2)$  (with  $\mu$  and  $\sigma$  unknown).
- Given  $Y = Y_1, \dots, Y_n$  we know that that

$$R = \frac{Y_{n+1} - \bar{Y}_n}{\hat{\sigma}} \sim H_n$$

is a pivot (its distribution is independent of  $\mu$  and  $\sigma$ ).

- And hence that

$$C_{\text{pred}}(y_{n+1}, \mathbf{y}_{\text{obs}}) = H_n \left( \frac{y_{n+1} - \bar{y}_{\text{obs}}}{\hat{\sigma}_{\text{obs}}} \right)$$

satisfies the two criteria above and is a valid CD construction.

# Prototype illustration with dependent realisations

- Let

$$\begin{pmatrix} Y_0 \\ Y_1 \end{pmatrix} \sim N_2 \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right).$$

- Then  $Y_1 | Y_0 = y_0 \sim N(\rho y_0, 1 - \rho^2)$  and if  $\rho$  and  $\sigma$  are known

$$\frac{Y_1 - \rho y_0}{\sqrt{1 - \rho^2}} | y_0 \sim N(0, 1),$$

but this is true for all  $y_0$ , meaning that it is also true unconditional.

- This motivates

$$C_{\text{pred}}(y_1) = \Phi \left( \frac{y_1 - \rho y_0}{\sqrt{1 - \rho^2}} \right),$$

where  $\Phi$  is the cdf. of a standard normal, but is it more ‘natural’ to consider

$$C_{\text{pred}}(y_1) = \Phi \left( \frac{y_1 - \rho y_0}{\sqrt{1 - \rho^2}} \right) | y_0?$$

- It is less clear what happens in more general models, and in cases where parameters need to be estimated.

# A family of time series models

- Let

$$Y_t = m(x_t, \beta) + \epsilon_t, \quad \text{for all } t \geq n \text{ and } n \geq 1,$$

where:

- $m$  is known, but depending on an unknown parameter  $\beta$
  - $x_t$  is a vector-valued sequence of known covariates
  - $\epsilon_t$  is a stationary time series process with mean zero and dependency that is known up to an unknown parameter  $\theta$ .
- Here  $m$  is typically a polynomial of low order, e.g.

$$Y_t = \beta_0 + \beta_1 x_{t,1} + \cdots + \beta_p x_{t,p} + \epsilon_t,$$

and we will also often assume that  $\epsilon_t$  is an autoregressive process, i.e.

$$\epsilon_t = \rho_1 \epsilon_{t-1} + \cdots + \rho_q \epsilon_{t-q} + \sigma \delta_t,$$

where  $\delta_t$  are independent  $N(0, 1)$ .

- The framework is easily extended to other more general types of model constructions.

# How to make good predictions?

- If  $\epsilon_t$  above is Gaussian, the ‘answer’ is conditional expectation, i.e.

$$\hat{y}_0 = E[Y_0 | \mathbf{y}_{\text{obs}}, X, \hat{\beta}, \hat{\theta}],$$

where the unknown parameters are fitted by maximum likelihood.

- If  $\epsilon_t$  is an autoregressive process, it is quite common to use so-called ‘conditional quasi maximum likelihood’ (which simplifies the estimation to ordinary least squares (OLS)).

**Example:** If  $\epsilon_t \sim \text{AR}(q)$  with Gaussian errors, then the likelihood  $\ell_n(\beta, \theta | \mathbf{y}_{\text{obs}}) | y_1, \dots, y_q$  is structurally equivalent to that of a linear regression model.

- There is a rich theory that aims at linear predictors, say

$$\hat{y}_0 = \hat{y}_{n+h} = \alpha^t \mathbf{x}_{n+h} + \gamma^t \mathbf{y}_{\text{obs}}, \quad \text{for } h \geq 1,$$

which minimises the mean squared error  $E(Y_0 - \hat{Y}_0)^2$ , see e.g. Grenander & Rosenblatt (1957).

- We will argue that linear predictors fitted by minimising the sum of squared prediction errors have some desired qualities.

## Predicting with the true model

- As simple starting point suppose  $Y_t$  is an autoregressive model of order one, then

$$Y_t = \epsilon_t = \rho\epsilon_{t-1} + \sigma\delta_t = \rho Y_{t-1} + \sigma\delta_t, \quad \text{for } t \leq n,$$

where  $\delta_t$  are independent  $N(0, 1)$ .

- Given the past  $Y_1, \dots, Y_n$  our best guess for  $Y_0 = Y_{n+1}$  with  $\rho$  known is  $\hat{Y}_{n+1} = \rho Y_n$ .
- Moreover, if  $\sigma$  is also known

$$R_n = \frac{Y_{n+1} - \hat{Y}_{n+1}}{\sigma} = \frac{\rho Y_n + \sigma\epsilon_t - \rho Y_n}{\sigma} = \epsilon_t \sim N(0, 1),$$

which is a pivot.

- Resulting in the simple confidence distribution construction:

$$C(y_0, \mathbf{y}_{\text{obs}}) = \Phi\left(\frac{y_0 - \rho y_n}{\sigma}\right).$$

## What if some parameters are unknown?

- If  $\rho$  is known and  $\sigma$  has been estimated, then

$$R_n = \frac{Y_{n+1} - \hat{Y}_{n+1}}{\hat{\sigma}} \sim t_{n-2},$$

since

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{t=2}^n (Y_t - \hat{Y}_{t-1})^2 \sim \chi_{n-2}^2$$

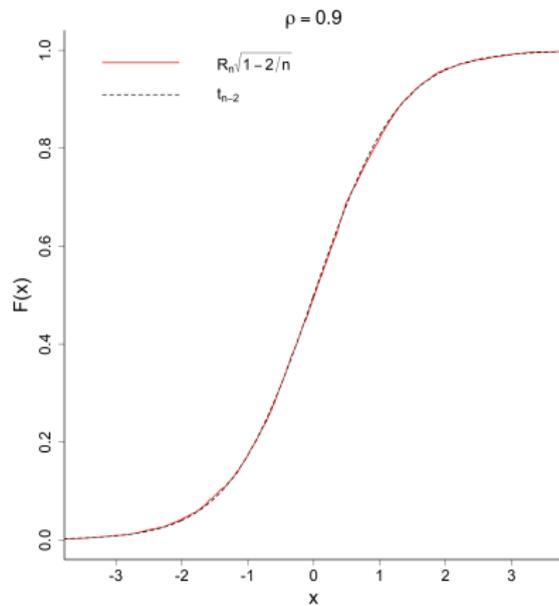
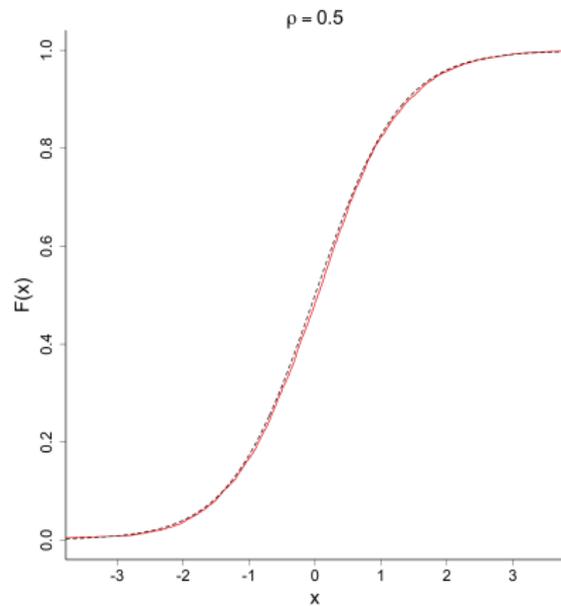
- If  $\rho$  also has to be estimated, the hope is that if  $\hat{\rho}$  is close to  $\rho$

$$R_n = \frac{Y_{n+1} - \hat{\rho}Y_n}{\hat{\sigma}}$$

is an approximate pivot; at least for large  $n$ .

- In simulation studies, we observe that  $R_n \sqrt{1 - 2/n}$  is close to a  $t_{n-2}$ -distribution, for all  $\rho$  and  $\sigma$ .

# 'Proof of concept' with $n = 10$ samples



## A more general modelling framework

- In practice, we (currently) obtain the best simulation results (fit and coverage) with

$$R_n = \frac{Y_{n+1} - \widehat{Y}_{n+1}}{s},$$

where

$$s^2 = \frac{1}{n-2} \sum_{t=2}^n (r_t - \bar{r}_n)^2$$

and where  $r_t = Y_t - \widehat{Y}_t$  are the residuals.

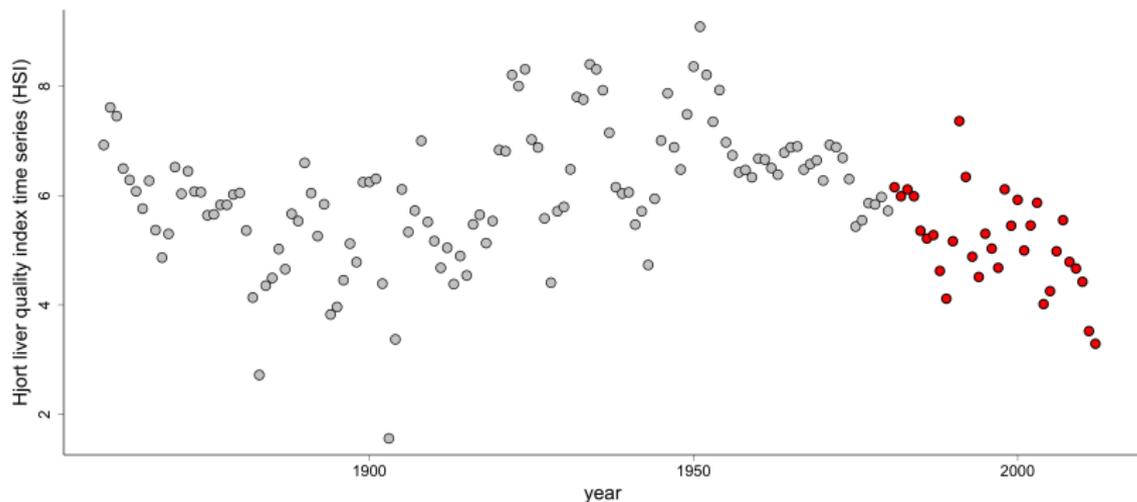
- A similar results seems to be true (in simulations) for general linear predictors of the type

$$\widehat{y}_0 = \widehat{y}_{n+h} = \widehat{\alpha}^t \mathbf{x}_{n+h} + \widehat{\gamma}^t \mathbf{y}_{\text{obs}}, \quad \text{for } h \geq 1,$$

where  $\widehat{\alpha}$  and  $\widehat{\gamma}$  are fitted by minimising the empirical sum of squared prediction errors.

- The approximation seems to be quite robust, also in very small samples, and to work well under model misspecification.

# Illustration: The liver quality index

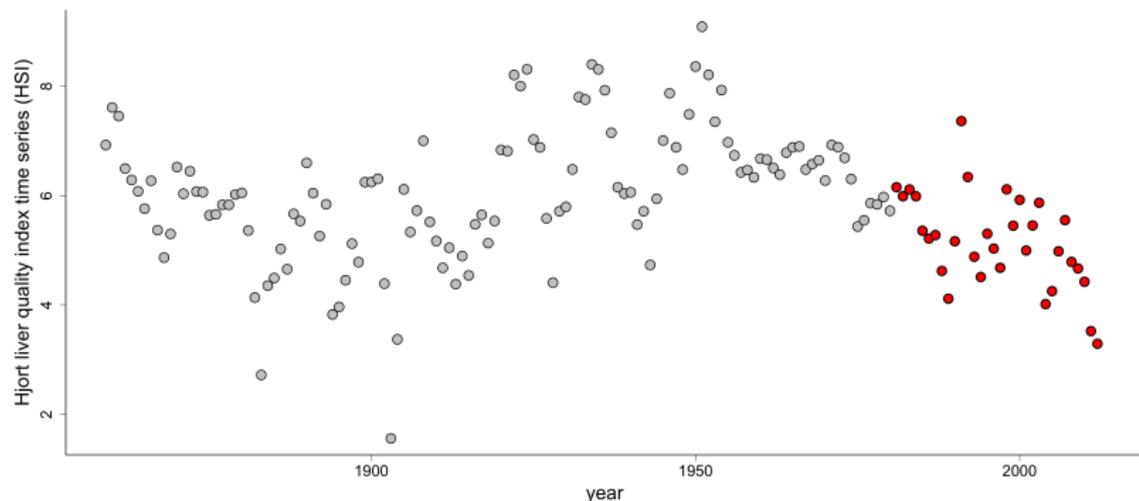


- The liver quality index (HSI) for a individual fish is defined as

$$\text{HSI}_{\text{fish}} = 100 \times \frac{\text{weight of liver}}{\text{weight of fish}}.$$

- The HSI index is viewed as a measure of the quality of life and is e.g. related to reproduction.

# Illustration: The liver quality index

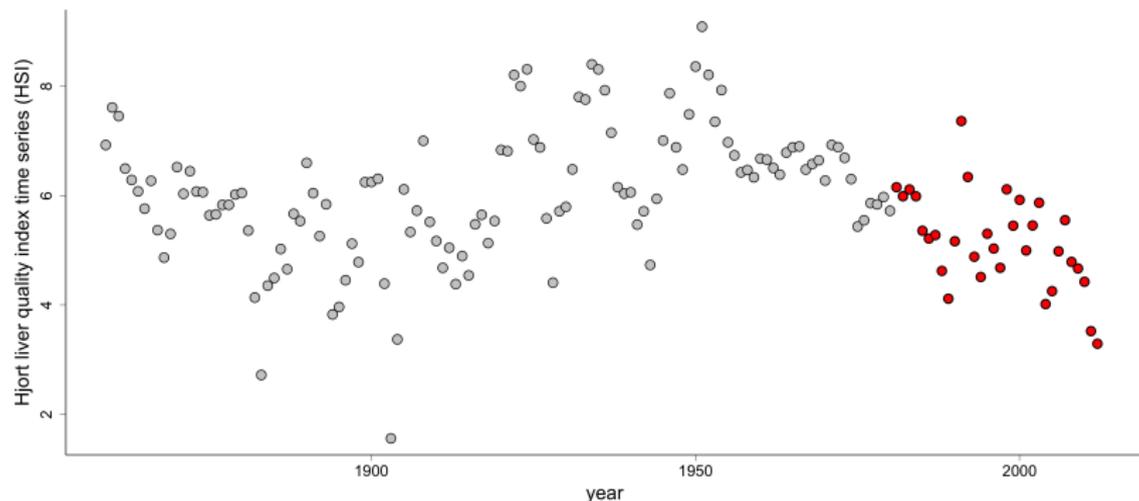


- For a population of fish in a given year  $i$ , the liver quality index is commonly estimated by

$$\widehat{\text{HSI}}_{\text{bulk},i} = 100 \times \frac{\text{total amount of liver}}{\text{total amount of fish}},$$

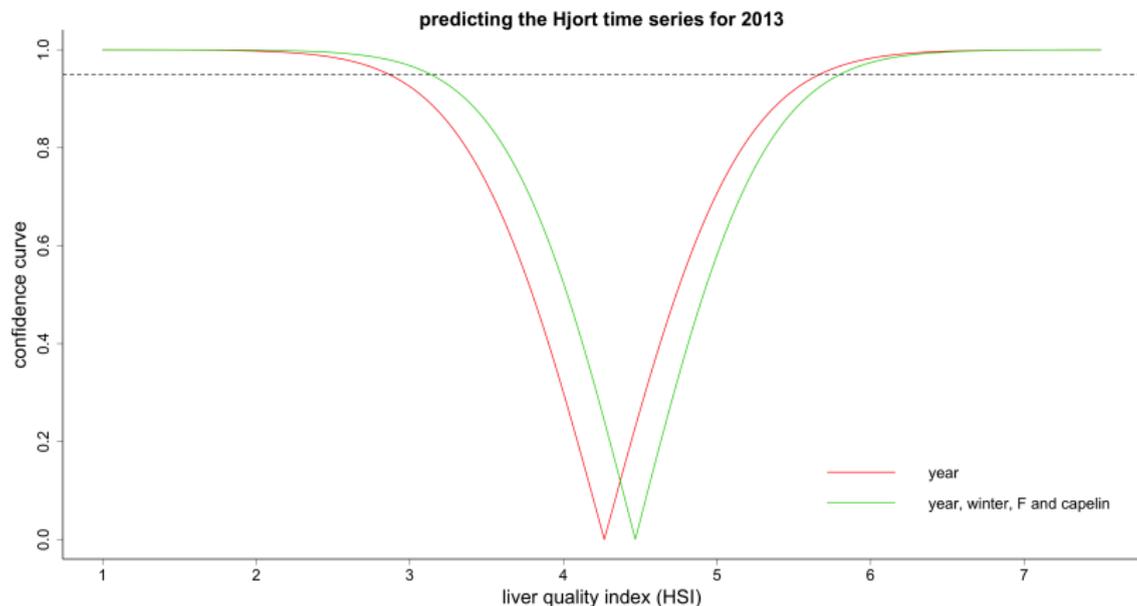
for fish caught in year  $i$ .

# Illustration: The liver quality index



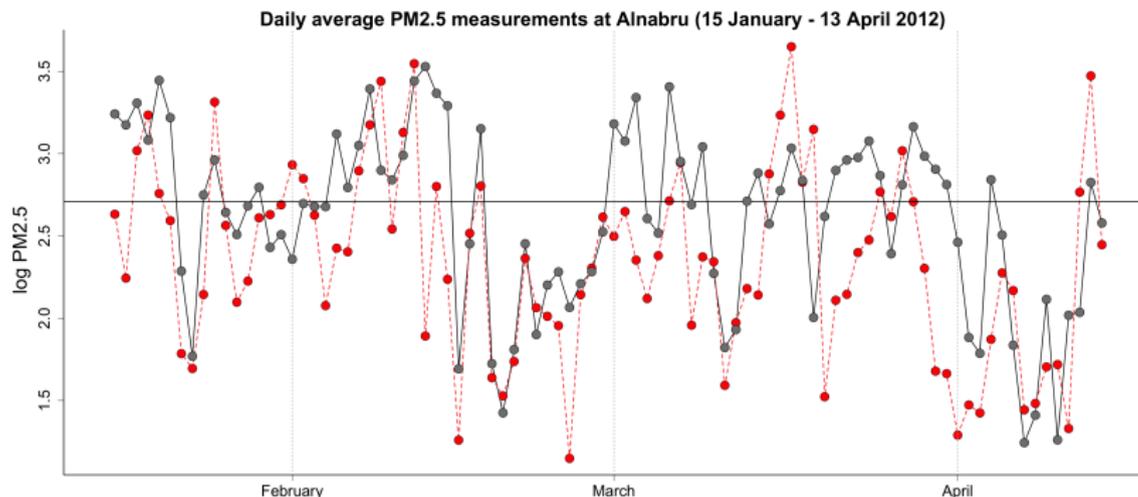
- It is important to understand the dynamic of the HSI index, e.g. how external factors, like changes in **sea temperature, death rates and food supply** affects the series.
- As an illustration we will compare two confidence curves for the ‘future’ HSI index of 2013 (i.e. the next observation).

# Confidence distribution for the next observation



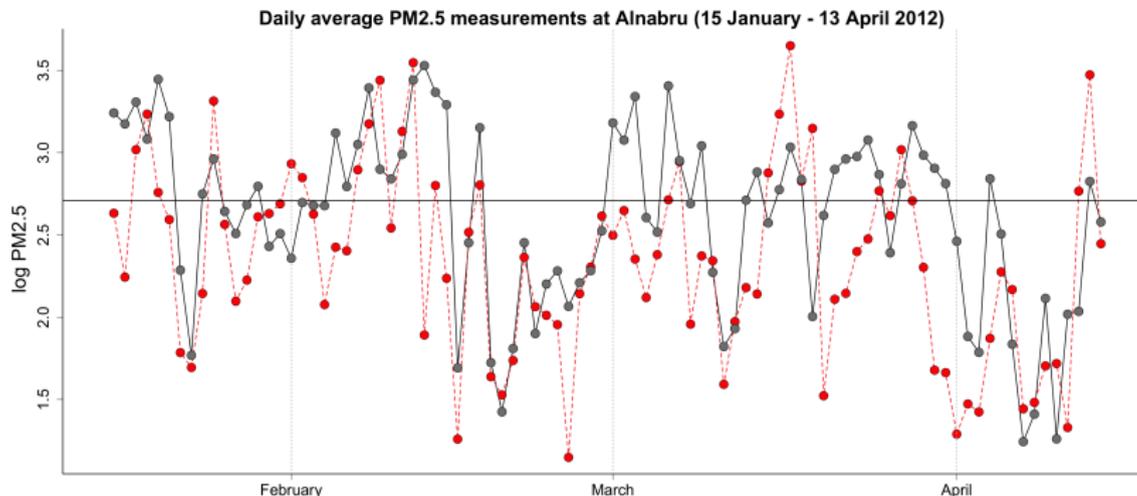
- The confidence curve are constructed using the approximatin introduced above.
- The idea is to see how adding additional covariates affects the prediction with respect to the confidence curve.

# Illustration: Particle pollution (PM<sub>2.5</sub>)



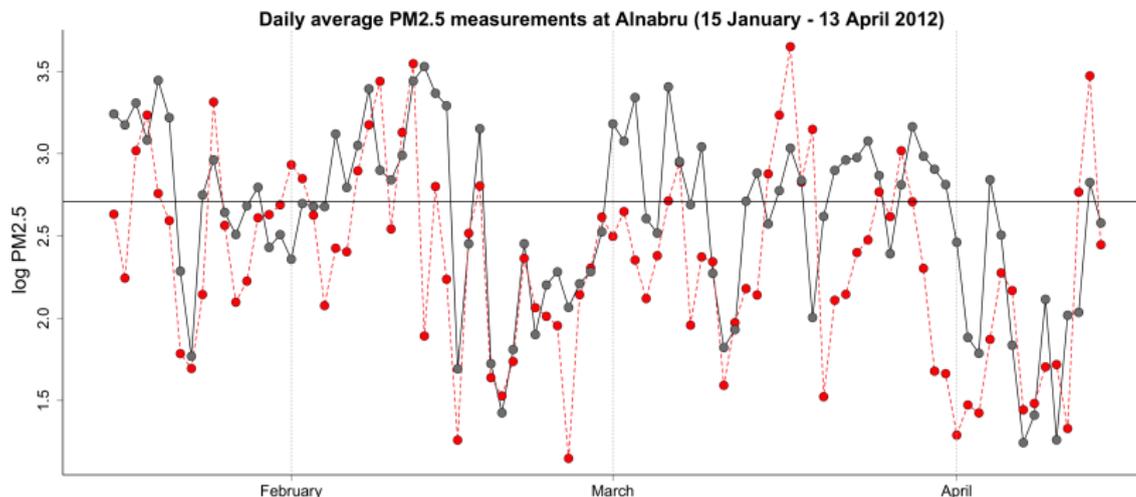
- Daily average measurements of fine particulate matter (PM<sub>2.5</sub>), a type of air pollutant of tiny particles, or droplets, in the air.
- Exposure can cause short-term effects, like eye, nose, throat and lung irritation, long-term exposure can affect lung function and has also been shown to be related to asthma and heart disease.

# Illustration: Particle pollution (PM<sub>2.5</sub>)



- For health reasons, government regulations typically restrict the daily average emission to 25-35  $\mu\text{g}/\text{m}^3$ .
- New directives in Norway suggest that the level may be set as low as 15  $\mu\text{g}/\text{m}^3$  with the possibility of 8 exceedances per year.

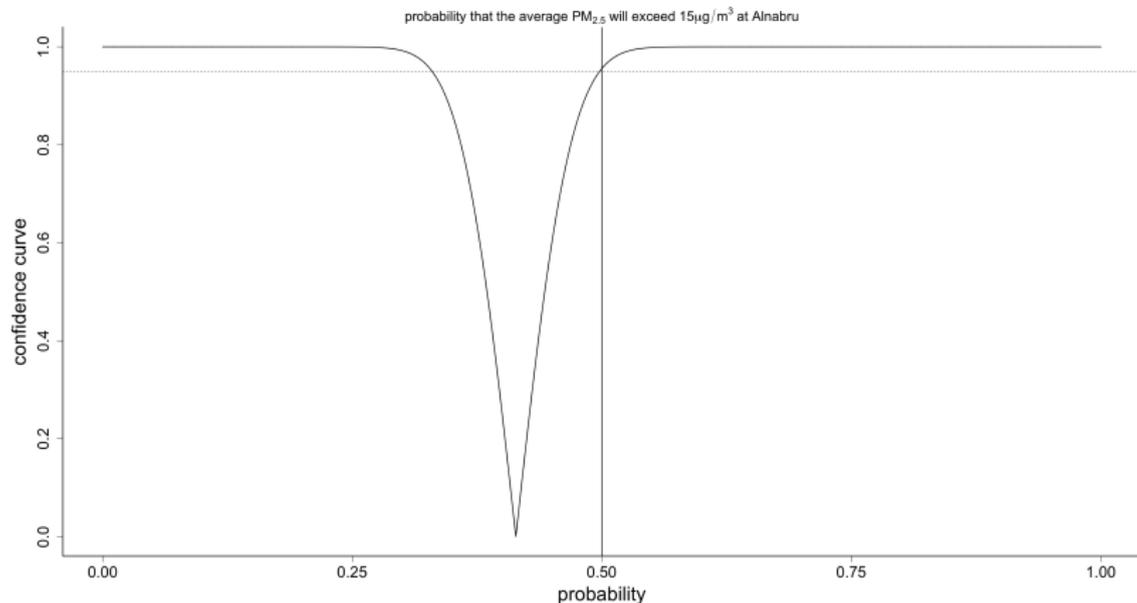
# Illustration: Particle pollution (PM<sub>2.5</sub>)



- One of the main objective is to use forecast to detect periods, or days with potential high emissions, enabling the legislature to take appropriate actions in advanced.
- As an application we will make the confidence curve for

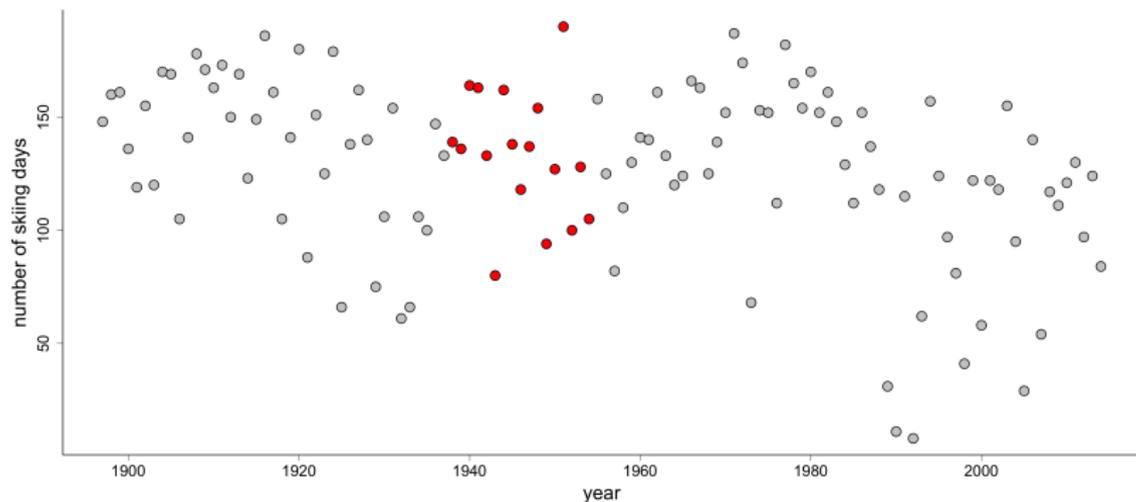
$$\psi = \psi(\theta) = \Pr\{Y_{n+1} \geq \log(15) \mid \mathbf{x}_{\text{forecast}}, \mathbf{y}_{\text{obs}}, \theta\}.$$

## Illustration: Particle pollution ( $PM_{2.5}$ )



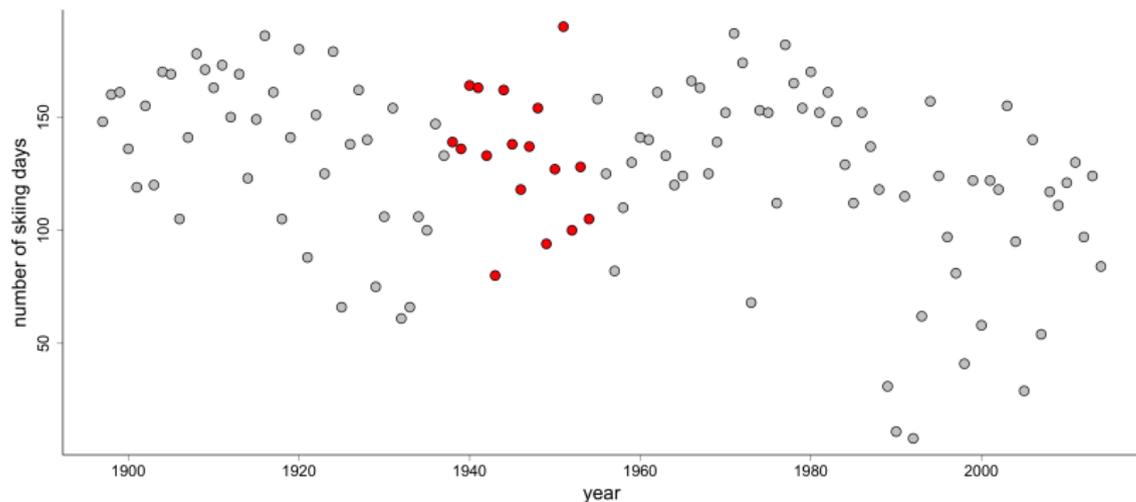
- The confidence curve is constructed using the delta method and the asymptotic normality of the estimated parameters.
- In simulations the large-sample approximate seems robust and the profile log-likelihood approach is often unstable.

## Illustration: Number of skiing days



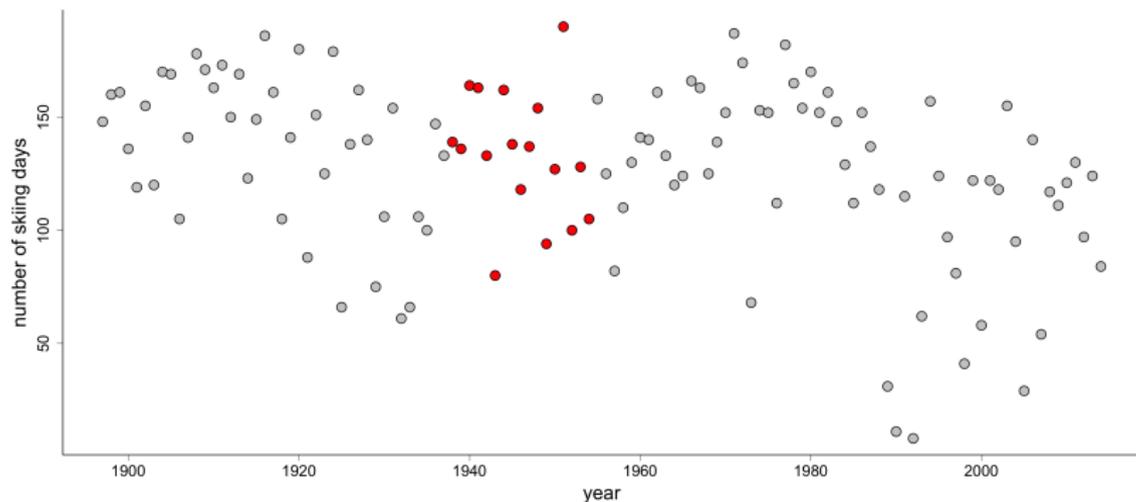
- The number of skiing days in a winter season is defined as the number of days with at least 25 cm snow.
- Meteorologist Gustav Bjørnbæk introduced the definition as the least amount of snow needed to avoid serious injury (in case of a fall).

## Illustration: Number of skiing days



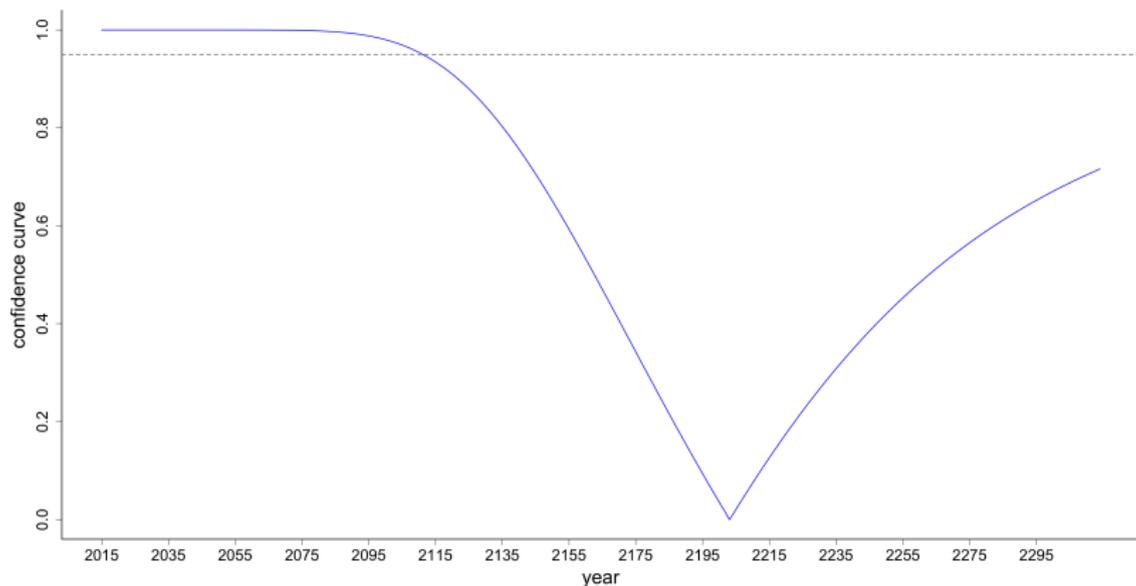
- Besides being of great interest to skiing enthusiasts, the number of skiing days can be seen as an indicator of how cold a winter is and is also an indication of the general temperature over a given period of time.
- Since 1896 the average decline is about  $-0.33$  skiing days per year.

## Illustration: Number of skiing days



- Since the mid 1950s the decline has a estimated slope of  $-0.9$ .
- The goal is to study, and make a confidence curve for, the time point  $t$  the expected value  $E Y_t$  crosses  $y_0 = 30$ .

## Illustration: Number of skiing days



- The confidence curve is fitted using the profile log-likelihood approach and converting the deviance function.