

# Modeling and Analysis of Time Series Data

## Chapter 5: Parameter estimation and model identification for ARMA models

Edward L. Ionides


# Outline

- 1 Likelihood-based inference in the context of ARMA models
  - The maximum likelihood estimator
  - Fisher information
  - Profile likelihood confidence intervals
  - Bootstrap standard errors
- 2 Model selection for ARMA models
  - Likelihood ratio tests
  - Akaike's information criterion (AIC)
- 3 Fitting ARMA models in R
  - Examining the AR and MA roots
  - A simulation study
  - Assessing numerical correctness

## Background on likelihood-based inference

- For any data  $y_{1:N}^*$  and any probabilistic model  $f_{Y_{1:N}}(y_{1:N}; \theta)$  we define the likelihood function to be

$$\mathcal{L}(\theta) = f_{Y_{1:N}}(y_{1:N}^*; \theta).$$

- It is often convenient to work with the logarithm to base  $e$  of the likelihood, which we write as 

$$\ell(\theta) = \log \mathcal{L}(\theta).$$

- Using the likelihood function as a statistical tool is a very general technique, widely used since Fisher (1922) ([Wikipedia: Likelihood\\_function](#)).
- Time series analysis involves various situations where we can, with sufficient care, compute the likelihood function and take advantage of the general framework of likelihood-based inference.

- Computation of the likelihood function for ARMA models is not entirely straightforward.
- Computationally efficient algorithms exist, using a state space model representation of ARMA models that will be developed later in this course.
- For now, it is enough that software exists to evaluate and maximize the likelihood function for a Gaussian ARMA model. Our immediate task is to think about how to use that capability.

- Before evaluation of the ARMA likelihood became routine, it was popular to use a method of moments estimator called **Yule-Walker** estimation (Shumway and Stoffer, 2017, Section 3.5). This is nowadays mostly of historical interest.
- For massively long time series data and big ARMA models, it can be computationally infeasible to work with the likelihood function. However, we are going to focus on the common situation where we can (with due care) work with the likelihood.
- Likelihood-based inference (meaning statistical tools based on the likelihood function) provides tools for parameter estimation, standard errors, hypothesis tests and diagnosing model misspecification.
- Likelihood-based inference often (but not always) has favorable theoretical properties. Here, we are not especially concerned with the underlying theory of likelihood-based inference. On any practical problem, we can check the properties of a statistical procedure by simulation experiments.

# The maximum likelihood estimator (MLE)


- A maximum likelihood estimator (MLE) is

$$\hat{\theta}(y_{1:N}) = \arg \max_{\theta} f_{Y_{1:N}}(y_{1:N}; \theta),$$

where  $\arg \max_{\theta} g(\theta)$  means a value of argument  $\theta$  at which the maximum of the function  $g$  is attained, so  $g(\arg \max_{\theta} g(\theta)) = \max_{\theta} g(\theta)$ .

- If there are many values of  $\theta$  giving the same maximum value of the likelihood, then an MLE still exists but is not unique.
- The maximum likelihood estimate (also known as the MLE) is

*estimator  
evaluated  
at the data.*


$$\begin{aligned}\hat{\theta} &= \hat{\theta}(y_{1:N}^*) \\ &= \arg \max_{\theta} \mathcal{L}(\theta) \\ &= \arg \max_{\theta} \ell(\theta).\end{aligned}$$

**Question 5.1.** Why are  $\arg \max_{\theta} \mathcal{L}(\theta)$  and  $\arg \max_{\theta} \ell(\theta)$  the same?

since  $\log$  is an increasing function  
and the density is non-negative.

- We can write  $\hat{\theta}_{MLE}$  to denote the MLE if we are considering various alternative estimation methods. However, in this course, we will most often be using maximum likelihood estimation so we let  $\hat{\theta}$  correspond to this approach.

## Standard errors for the MLE

- As statisticians, it would be irresponsible to present an estimate without a measure of uncertainty!
- Usually, this means obtaining a confidence interval, or an approximate confidence interval.
- It is good to say **approximate** when you present something that is not exactly a confidence interval with the claimed coverage. For example, remind yourself of the definition of a 95% confidence interval.
- Saying “approximate” reminds you that there is some checking that could be done to assess how accurate the approximation is in your particular situation.



# Three ways to quantify statistical uncertainty in an MLE

- ① Fisher information. This is computationally quick, but works well only when  $\hat{\theta}(Y_{1:N})$  is well approximated by a normal distribution.
- ② Profile likelihood estimation. This is a bit more computational effort, but generally is preferable to the Fisher information.
- ③ A simulation study, also known as a bootstrap.

# Standard errors via the observed Fisher information

- We suppose that  $\theta \in \mathbb{R}^D$  and so we can write  $\theta = \theta_{1:D}$ .
- The **Hessian matrix** of a function is the matrix of its second partial derivatives. We write the Hessian matrix of the log likelihood function as  $\nabla^2 \ell(\theta)$ , a  $D \times D$  matrix whose  $(i, j)$  element is

$$[\nabla^2 \ell(\theta)]_{ij} = \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ell(\theta).$$

- The observed Fisher information is

$$\hat{I} = -\nabla^2 \ell(\hat{\theta}).$$

- A standard asymptotic approximation to the distribution of the MLE for large  $N$  is

$$\hat{\theta}(Y_{1:N}) \approx N \left[ \theta, \hat{I}^{-1} \right],$$

where  $\theta$  is the true parameter value. This asserts that the MLE is asymptotically unbiased, with variance asymptotically attaining the Cramer-Rao lower bound.

"is approximately distributed as"

- Since the MLE attains the Cramer-Rao lower bound, under regularity conditions, we it is **asymptotically efficient**.
- We can interpret  $\approx$  in the above normal approximation to mean “one could write a limit statement formally justifying this approximation in a suitable limit.” Almost equivalently,  $\approx$  can mean “this approximation is useful in the finite sample situation at hand.”
- A corresponding approximate 95% confidence interval for  $\theta_d$  is  $\hat{\theta}_d \pm 1.96([\hat{I}^{-1}]_{dd})^{1/2}$ . The R function `arima` computes standard errors for the MLE of an ARMA model in this way.
- We usually only have one time series, with some fixed  $N$ , and so we cannot in practice take  $N \rightarrow \infty$ . When our time series model is non-stationary it may not even be clear what it would mean to take  $N \rightarrow \infty$ . These asymptotic results should be viewed as nice mathematical reasons to consider computing an MLE, but not a substitute for checking how the MLE behaves for our model and data.

# Confidence intervals via the profile likelihood

- We consider the problem of obtaining a confidence interval for  $\theta_d$ , the  $d$ th component of  $\theta_{1:D}$ .
- The **profile log likelihood function** of  $\theta_d$  is defined to be

$$\underline{\ell_d^{\text{profile}}(\theta_d) = \max_{\phi \in \mathbb{R}^D: \phi_d = \theta_d} \ell(\phi).}$$

In general, the profile likelihood of one parameter is constructed by maximizing the likelihood function over all other parameters.

- Check that  $\max_{\theta_d} \ell_d^{\text{profile}}(\theta_d) = \max_{\theta_{1:D}} \ell(\theta_{1:D})$ . Maximizing the profile likelihood  $\ell_d^{\text{profile}}(\theta_d)$  gives the MLE,  $\hat{\theta}_d$ .

- An approximate 95% confidence interval for  $\theta_d$  is given by

*set of profile points within 1.92 log units of the maximized likelihood.*

$$\underline{\{\theta_d : \ell(\hat{\theta}) - \ell_d^{\text{profile}}(\theta_d) < 1.92\}.}$$

- This is known as a profile likelihood confidence interval.

## Where does the 1.92 cutoff come from

- The cutoff 1.92 is derived using **Wilks's theorem**, which we will discuss in more detail when we develop likelihood ratio tests.
- Note that  $1.92 = \frac{1.96^2}{2}$ .
- The asymptotic justification of Wilks's theorem is the same limit that justifies the Fisher information standard errors.
- Profile likelihood confidence intervals tend to work better than Fisher information confidence intervals when the log likelihood function is not close to quadratic near its maximum. This is more common when  $N$  is not large.

## A Simulation study, also called bootstrap

- If done carefully and well, this can be the best approach.
- A confidence interval is a claim about reproducibility. You claim, so far as your model is correct, that on 95% of realizations from the model, a 95% confidence interval you have constructed will cover the true value of the parameter.
- A simulation study can check this claim directly.
- The simulation study takes time to develop and debug, time to explain, and time for the reader to understand and check what you have done. We usually carry out simulation studies to check our main conclusions only.

# Bootstrap methods for constructing standard errors and confidence intervals

- Suppose we want to know the statistical behavior of the estimator  $\hat{\theta}(y_{1:N})$  for models in a neighborhood of the MLE.
- In particular, let's consider the problem of estimating uncertainty about  $\theta_1$ , the first component of the vector  $\theta$ .
- We use simulation to assess the behavior of the maximum likelihood estimator,  $\hat{\theta}_1(y_{1:N})$ , and possibly the coverage of an associated confidence interval estimator,  $[\hat{\theta}_{1,lo}(y_{1:N}), \hat{\theta}_{1,hi}(y_{1:N})]$ .
- The confidence interval estimator could be constructed using either the Fisher information method or the profile likelihood approach.

- We can design a simulation study to address the following goals:
  - (A) Evaluate the coverage of a proposed confidence interval estimator,  $[\hat{\theta}_{1,lo}, \hat{\theta}_{1,hi}]$ ,
  - (B) Construct a standard error for  $\hat{\theta}_1$ ,
  - (C) Construct a confidence interval for  $\theta_1$  with exact local coverage.



## A simulation study

1. Generate  $J$  independent Monte Carlo simulations,

$$Y_{1:N}^{[j]} \sim f_{Y_{1:N}}(y_{1:N}; \hat{\theta}) \text{ for } j \in 1 : J.$$

2. For each simulation, evaluate the maximum likelihood estimator,

$$\hat{\theta}^{[j]} = \hat{\theta}(Y_{1:N}^{[j]}) \text{ for } j \in 1 : J,$$

and, if desired, the confidence interval estimator,

$$[\hat{\theta}_{1,lo}^{[j]}, \hat{\theta}_{1,hi}^{[j]}] = [\hat{\theta}_{1,lo}(Y_{1:N}^{[j]}), \hat{\theta}_{1,hi}(Y_{1:N}^{[j]})].$$

3. For large  $J$ , the coverage of the proposed confidence interval is well approximated, for models in a neighborhood of  $\hat{\theta}$ , by the proportion of the intervals  $[\hat{\theta}_{1,lo}^{[j]}, \hat{\theta}_{1,hi}^{[j]}]$  that include  $\hat{\theta}_1$ .
4. The sample standard deviation of  $\{\hat{\theta}_1^{[j]}, j \in 1 : J\}$  is a natural standard error to associate with  $\hat{\theta}_1$ .

# Likelihood ratio tests for nested hypotheses

- The whole parameter space on which the model is defined is  $\Theta \subset \mathbb{R}^D$ .
- Suppose we have two **nested** hypotheses

$$\begin{aligned}H^{(0)} &: \theta \in \Theta^{(0)}, \\H^{(1)} &: \theta \in \Theta^{(1)},\end{aligned}$$

defined via two nested parameter subspaces,  $\Theta^{(0)} \subset \Theta^{(1)}$ , with respective dimensions  $D^{(0)} < D^{(1)} \leq D$ .

- We consider the log likelihood maximized over each of the hypotheses,

$$\begin{aligned}\ell^{(0)} &= \sup_{\theta \in \Theta^{(0)}} \ell(\theta), \\ \ell^{(1)} &= \sup_{\theta \in \Theta^{(1)}} \ell(\theta).\end{aligned}$$

- A useful approximation asserts that, under the hypothesis  $H^{\langle 0 \rangle}$ ,

$$\ell^{\langle 1 \rangle} - \ell^{\langle 0 \rangle} \approx (1/2) \chi_{D^{\langle 1 \rangle} - D^{\langle 0 \rangle}}^2,$$

where  $\chi_d^2$  is a chi-squared random variable on  $d$  degrees of freedom and  $\approx$  means "is approximately distributed as."

- We will call this the **Wilks approximation**.
- The Wilks approximation can be used to construct a hypothesis test of the null hypothesis  $H^{\langle 0 \rangle}$  against the alternative  $H^{\langle 1 \rangle}$ .
- This is called a **likelihood ratio test** since a difference of log likelihoods corresponds to a ratio of likelihoods.
- When the data are iid,  $N \rightarrow \infty$ , and the hypotheses satisfy suitable regularity conditions, this approximation can be derived mathematically and is known as **Wilks's theorem**.
- The chi-squared approximation to the likelihood ratio statistic may be useful, and can be assessed empirically by a simulation study, even in situations that do not formally satisfy any known theorem.

# Using a likelihood ratio test to construct profile likelihood confidence intervals

- Recall the duality between hypothesis tests and confidence intervals:

The estimated parameter  $\theta^*$  does not lead us to reject a null hypothesis of  $\theta = \theta^{(0)}$  at the 5% level



$\theta^{(0)}$  is in a 95% confidence interval for  $\theta$ .

- We can check what the 95% cutoff is for a chi-squared distribution with one degree of freedom,

```
qchisq(0.95,df=1)
```

```
[1] 3.841459
```

- We can now see how the Wilks approximation suggests a confidence interval constructed from parameter values having a profile likelihood within 1.92 log units of the maximum.

# Akaike's information criterion (AIC)

- Likelihood ratio tests provide an approach to model selection for nested hypotheses, but how about when models are not nested?
- A more general approach is to compare likelihoods of different models by penalizing the likelihood of each model by a measure of its complexity.
- Akaike's information criterion **AIC** is given by

$$AIC = -2 \times \ell(\theta^*) + 2D$$

“Minus twice the maximized log likelihood plus twice the number of parameters.”

- We are invited to select the model with the lowest AIC score.
- AIC was derived as an approach to minimizing prediction error. Increasing the number of parameters leads to additional **overfitting** which can decrease predictive skill of the fitted model.

## A caution for using AIC

- Viewed as a hypothesis test, AIC may have weak statistical properties.
- It is a mistake to interpret AIC by making a claim that the favored model has been shown to provides a superior explanation of the data.
- However, viewed as a way to select a model with reasonable predictive skill from a range of possibilities, it is often useful.

## Comparing AIC with likelihood ratio tests

**Question 5.2.** Suppose we are in a situation in which we wish to choose between two nested hypotheses, with dimensions  $D^{(0)} < D^{(1)}$ . Suppose the Wilks approximation is valid. Consider the strategy of selecting the model with the lowest AIC value, and view this model selection approach as a formal statistical test.

(A) Find an expression for the size of this AIC test (i.e, the probability of rejecting the null hypothesis,  $H^{(0)}$ , when this null hypothesis is true).

(B) Evaluate this expression for  $D^{(1)} - D^{(0)} = 1$ .

# Likelihood-based inference for ARMA models in R

- The Great Lakes are an important resource for leisure, agriculture and industry in this region.
- A past concern has been whether human activities such as water diversion or channel dredging might be leading to a decline in lake levels.
- A current concern has been high levels leading to coastal erosion.
- Are lake levels affected by climate change?
- The physical mechanisms are not always obvious: for example, evaporation tends to be highest when the weather is cold but the lake is not ice-covered.
- We look at monthly time series data on the level of Lake Huron, which is essentially the same as Lake Michigan.



## Reading in the data

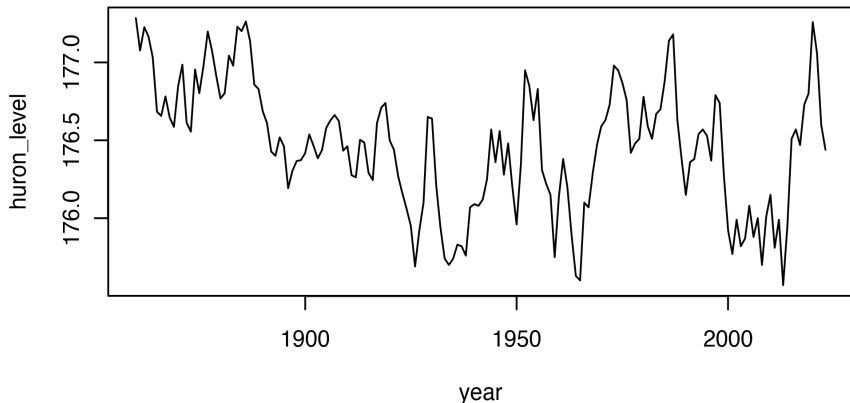
The file `huron_level.csv` gives monthly water level, in meters, for Lakes Michigan and Huron from 1860 to 2023.

```
dat <- read.table(file="huron_level.csv", sep=",", header=TRUE)
head(dat[,1:7], 2)
```

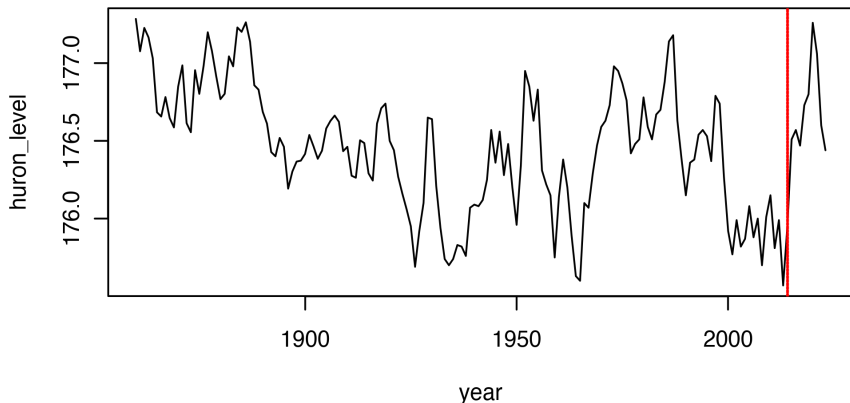
Year	Jan	Feb	Mar	Apr	May	Jun
1860	177.285	177.339	177.349	177.388	177.425	177.461
1861	177.077	177.105	177.224	177.254	177.382	177.431

For now, we avoid monthly seasonal variation by considering an annual series of January depths. We will investigate seasonal variation later in the course, but sometimes it is best avoided.

```
huron_level <- dat$Jan  
year <- dat$Year  
plot(huron_level~year,type="l")
```



- Until the recent surge in water level, there was concern about a long-run decline in lake level due to dredging or water diversion or climate change.
- We put ourselves back in 2014 and temporarily ignore subsequent data



# Fitting an ARMA model

- Later, we will consider hypotheses of trend. For now, let's start by fitting a stationary ARMA( $p, q$ ) model under the null hypothesis that there is no trend. This hypothesis, which asserts that nothing has substantially changed in this system over the last 160 years, is not entirely unreasonable from looking at the data.
- We seek to fit a stationary Gaussian ARMA( $p, q$ ) model with parameter vector  $\theta = (\phi_{1:p}, \psi_{1:q}, \mu, \sigma^2)$  given by

$$\phi(B)(Y_n - \mu) = \psi(B)\epsilon_n,$$

where

$$\begin{aligned}\mu &= \mathbb{E}[Y_n] \\ \phi(x) &= 1 - \phi_1 x - \cdots - \phi_p x^p, \\ \psi(x) &= 1 + \psi_1 x + \cdots + \psi_q x^q, \\ \epsilon_n &\sim \text{iid } N[0, \sigma^2].\end{aligned}$$

## Choosing $p$ and $q$

- We need to decide where to start in terms of values of  $p$  and  $q$ .
- We tabulate AIC values for a range of different choices of  $p$  and  $q$ .

```
aic_table <- function(data,P,Q){
  table <- matrix(NA,(P+1),(Q+1))
  for(p in 0:P) {
    for(q in 0:Q) {
      table[p+1,q+1] <- arima(data,order=c(p,0,q))$aic
    }
  }
  dimnames(table) <- list(paste("AR",0:P, sep=""),
    paste("MA",0:Q,sep=""))
  table
}
huron_aic_table <- aic_table(huron_level,4,5)
require(knitr)
kable(huron_aic_table,digits=2)
```

	MA0	MA1	MA2	MA3	MA4	MA5
AR0	166.78	46.98	7.71	-13.70	-17.62	-24.89
AR1	-37.25	-36.62	-34.74	-33.13	-33.14	-31.18
AR2	-36.52	-37.41	-35.89	-33.89	-33.24	-31.91
AR3	-34.79	-34.43	-32.44	-31.89	-32.05	-32.14
AR4	-33.19	-33.91	-33.48	-33.54	-30.15	-29.52

**Question 5.3.** What do we learn by interpreting the results in the above table of AIC values?

**Question 5.4.** In what ways might we have to be careful not to over-interpret the results of this table?

- Let's fit the ARMA(2,1) model recommended by consideration of AIC.

```
huron_arma21 <- arima(huron_level, order=c(2,0,1))
huron_arma21
```

Call:

```
arima(x = huron_level, order = c(2, 0, 1))
```

Coefficients:

	ar1	ar2	ma1	intercept
	-0.0561	0.7935	1.0000	176.4591
s.e.	0.0521	0.0525	0.0257	0.1209

sigma^2 estimated as 0.04217: log likelihood = 23.71, aic = -37.41

- We can examine the roots of the AR polynomial,

```
AR_roots <- polyroot(c(1,-coef(huron_arma21)[c("ar1","ar2")]))
AR_roots

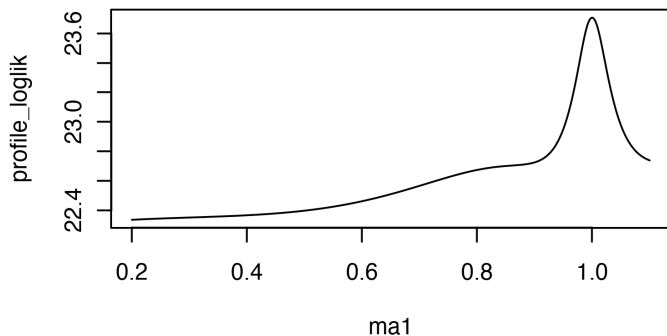
[1] 1.158532-0i -1.087774+0i
```

- The roots are just outside the unit circle, suggesting we have a stationary causal fitted ARMA.
- However, the MA root is  $-1$ , showing that the fitted model is at the threshold of non-invertibility.
- Do we have a non-invertibility problem? We investigate this using profile and bootstrap methods. The claimed standard error on the MA1 coefficient, from the Fisher information approach used by `arma`, is small.



- First, we can see if the approximate confidence interval constructed using profile likelihood is in agreement with the approximate confidence interval constructed using the observed Fisher information.
- To do this, we need to maximize the ARMA likelihood while fixing the MA1 coefficient at a range of values. This is done using `arma` in the code below.
- Note that the `fixed` argument expects a vector of length  $p + q + 1$  corresponding to a concatenated vector  $(\phi_{1:p}, \psi_{1:q}, \mu)$ . Somehow, the Gaussian white noise variance,  $\sigma^2$ , is not included in this representation. Parameters with NA entries in `fixed` are estimated.

```
K <- 500
ma1 <- seq(from=0.2,to=1.1,length=K)
profile_loglik <- rep(NA,K)
for(k in 1:K){
  profile_loglik[k] <- logLik(arma(huron_level,order=c(2,0,1),
    fixed=c(NA,NA,ma1[k],NA)))
}
plot(profile_loglik~ma1,ty="l")
```



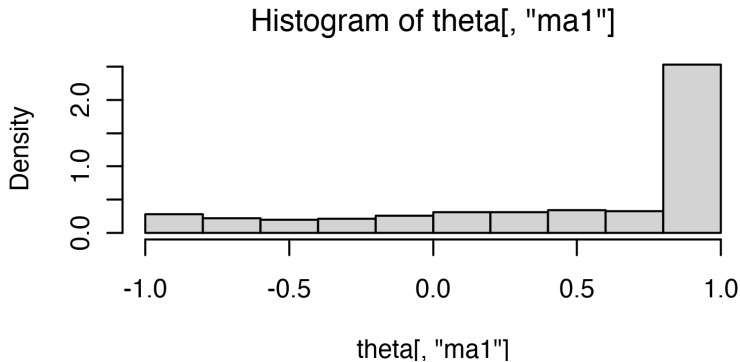
**Question 5.5.** Interpret the profile likelihood plot for  $\psi_1$ .

**Question 5.6.** What do you conclude about the Fisher information confidence interval proposed by arima?

**Question 5.7.** In what situations is the Fisher information confidence interval reliable?

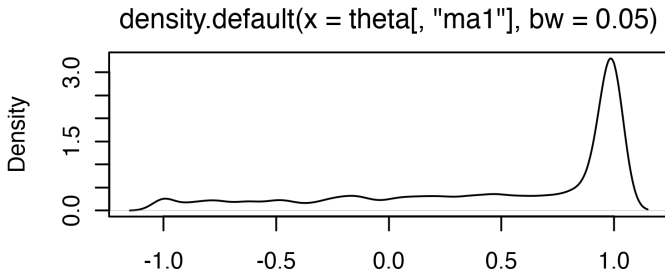
**Question 5.8.** Is this profile likelihood plot, and its statistical interpretation, reliable? How could you support your opinion on this?

```
set.seed(578922)
J <- 1000
params <- coef(huron_arma21)
ar <- params[grepl("^ar",names(params))]
ma <- params[grepl("^ma",names(params))]
intercept <- params["intercept"]
sigma <- sqrt(huron_arma21$sigma2)
theta <- matrix(NA,nrow=J,ncol=length(params),
  dimnames=list(NULL,names(params)))
for(j in 1:J){
  try({
    Y_j <- arima.sim(
      list(ar=ar,ma=ma),
      n=length(huron_level),
      sd=sigma
    )+intercept
    theta[j,] <- coef(arima(Y_j,order=c(2,0,1)))
  })
}
theta <- na.omit(theta)
hist(theta[, "ma1"],freq=FALSE)
```



- This seems consistent with the profile likelihood plot.
- A density plot shows this similarity even more clearly.

```
plot(density(theta[, "ma1"], bw=0.05))
```



N = 1000 Bandwidth = 0.05

- Here, we look at the raw plot for instructional purposes. For a report, one should improve the default axis labels and title.
- Note that arima transforms the model to invertibility. Thus, the estimated value of  $\theta_1$  can only fall in the interval  $[-1, 1]$ .

```
range(theta[, "ma1"])
```

```
[1] -1 1
```

- A minor technical issue: estimated densities outside  $[-1, 1]$  are artifacts of the density estimation procedure.

**Question 5.9.** How would you refine this density estimation procedure to respect the range of the parameter estimation procedure?

- We do a simulation study for which we fit ARMA(2,1) when the true model is AR(1).

## Using multiple cores for simulation studies

- When doing simulation studies, **multicore computing** is helpful. All modern computers have multiple cores.
- A basic approach to multicore statistical computing is to tell R you want it to look for available processors, using the `doParallel` package.
- We can use `foreach` in the `doParallel` package to carry out a parallel for loop where jobs are sent to different processors.

```
library(doParallel)  
registerDoParallel()
```



```
J <- 1000
huron_ar1 <- arima(huron_level, order=c(1,0,0))
params <- coef(huron_ar1)
ar <- params[grepl("^ar", names(params))]
intercept <- params["intercept"]
sigma <- sqrt(huron_ar1$sigma2)
t1 <- system.time(
  huron_sim <- foreach(j=1:J) %dopar% {
    Y_j <- arima.sim(list(ar=ar), n=length(huron_level),
      sd=sigma)+intercept
    try(coef(arima(Y_j, order=c(2,0,1))))
  }
)
```

- Some of these `arma` calls did not successfully produce parameter estimates. The `try` function lets the simulation proceed despite these errors. Let's see how many of them fail:

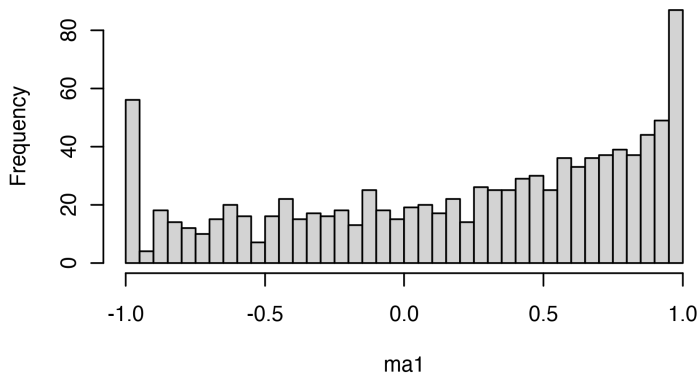
```
sum(sapply(huron_sim, function(x) inherits(x, "try-error")))
```

```
[1] 3
```

- Now, for the remaining ones, we can look at the resulting estimates of the MA1 component:

```
ma1 <- unlist(lapply(huron_sim,function(x)  
  if(!inherits(x,"try-error"))x["ma1"] else NULL ))  
hist(ma1,breaks=50)
```

Histogram of ma1



- When the true model is AR1 and we fit ARMA(2,1), it seems that we often obtain a model with estimated MA1 coefficient on the boundary of invertibility.
- Thus, we cannot reject an AR1 hypothesis for the Huron data, even though the Fisher information based analysis appears to give strong evidence that the data should be modeled with a nonzero MA1 coefficient.
- It may be sensible to avoid fitted models too close to the boundary of invertibility. This is a reason not to blindly accept whatever model AIC might suggest.

**Question 5.10.** What else could we look for to help diagnose, and understand, this kind of model fitting problem? Hint: pay some more attention to the roots of the fitted ARMA(2,1) model.

# Assessing the numerical correctness of evaluation and maximization of the likelihood function

- We can probably suppose that `arima()` has negligible numerical error in evaluating the likelihood.
- Likelihood evaluation is a linear algebra computation which should be numerically stable away from singularities.
- Possibly, numerical problems could arise for models very close to reducibility (canceling AR and MA roots).
- Numerical optimization is more problematic.
- `arima` calls the general purpose optimization routine `optim`.
- The likelihood surface can be multimodal and have nonlinear ridges, when AR and MA roots almost cancel.
- No optimization procedure is reliable for maximizing awkward, non-convex functions.
- Evidence for imperfect maximization (assuming negligible likelihood evaluation error) can be found in the AIC table, copied below.

	MA0	MA1	MA2	MA3	MA4	MA5
AR0	166.8	47.0	7.7	-13.7	-17.6	-24.9
AR1	-37.2	-36.6	-34.7	-33.1	-33.1	-31.2
AR2	-36.5	-37.4	-35.9	-33.9	-33.2	-31.9
AR3	-34.8	-34.4	-32.4	-31.9	-32.0	-32.1
AR4	-33.2	-33.9	-33.5	-33.5	-30.1	-29.5

**Question 5.11.** How is this table inconsistent with perfect maximization?

- Hint: recall that, for nested hypotheses  $H^{\langle 0 \rangle} \subset H^{\langle 1 \rangle}$ , the likelihood maximized over  $H^{\langle 1 \rangle}$  cannot be less than the likelihood maximized over  $H^{\langle 0 \rangle}$ .
- Recall also the definition of AIC,  
$$\text{AIC} = -2 \times \text{maximized log likelihood} + 2 \times \text{number of parameters}$$

## Further reading

- Section 3.5 of Shumway and Stoffer (2017) gives a complementary discussion of parameter estimation for ARMA models.
- Section 3.7 of Shumway and Stoffer (2017) takes a different perspective on selecting ARMA models, putting less emphasis on likelihood. Both perspectives can be valuable.



# References and Acknowledgements

Shumway RH, Stoffer DS (2017). *Time Series Analysis and its Applications: With R Examples*. 4th edition. Springer.

- Compiled on January 28, 2024 using R version 4.2.3.
- Licensed under the [Creative Commons Attribution-NonCommercial license](#). Please share and remix non-commercially, mentioning its origin.
- We acknowledge [previous versions of this course](#).

