

Midterm 1, STATS 531/631 W26

In class on 2/16

Name:

UMID:

This document produces different random tests each time the source code generating it is run. The actual midterm will be a realization generated by this random process, or something similar.

This version lists all the questions currently in the test generator. The actual test will have one question sampled from each of the 7 question categories.

Instructions. The test is closed book, and you are not allowed access to any notes. Any electronic devices in your possession must be turned off and remain in a bag on the floor.

For each question, circle one letter answer and provide some supporting reasoning.

Q1. Stationarity and unit roots.

Q1-01.

Suppose that a dataset $y_{1:N}^*$ is well described by the statistical model

$$Y_n = a + bn + \epsilon_n,$$

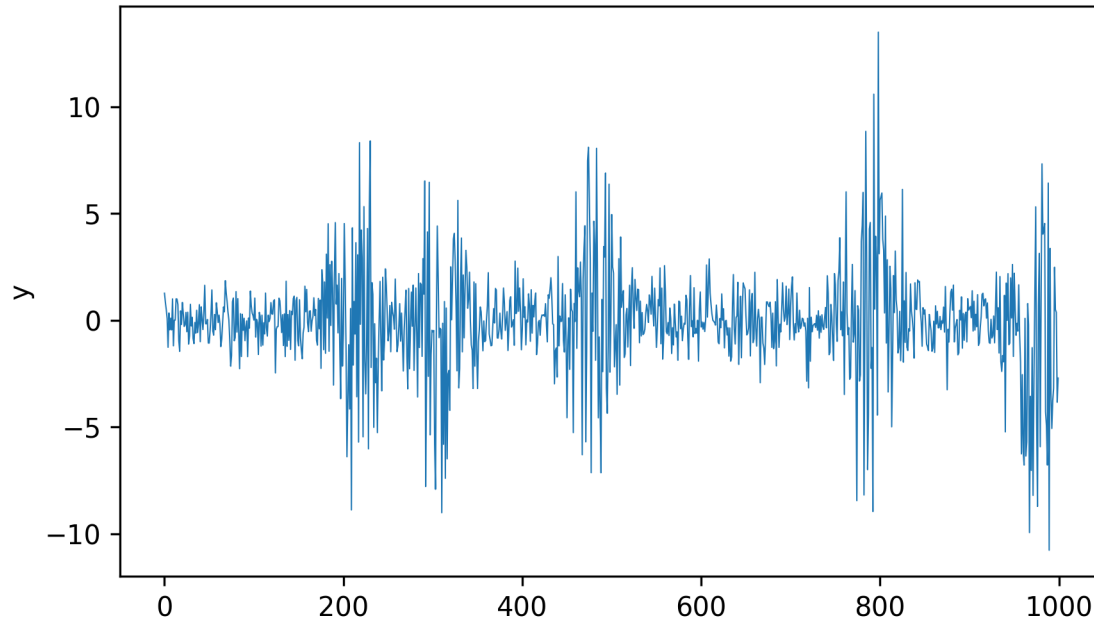
where ϵ_n is a Gaussian ARMA process and $b \neq 0$. Which of the following is the best approach to time series modeling of $y_{1:N}^*$?

- A. The data are best modeled as non-stationary, so we should take differences. The differenced data are well described by a stationary ARMA model.
- B. The data are best modeled as non-stationary, and we should use a trend plus ARMA noise model.
- C. The data are best modeled as non-stationary. It does not matter if we difference or model as trend plus ARMA noise since these are both linear time series models which become equivalent when we estimate their parameters from the data.
- D. We should be cautious about doing any of A, B or C because the data may have nonstationary sample variance in which case it may require a transformation before it is appropriate to fit any ARMA model.

Solution. B.

It does matter whether we take differences. For example, the differenced model is non-causal (has an MA root on the unit circle) so cannot be fitted by usual ARMA methods. D is not relevant since we are told that the data are well described by a model that rules out this possibility.

Q1-02.



Consider the time series plotted above. Which of the below is the most accurate statement about stationarity?

- A. The plot shows that the data are clearly non-stationary. We could make a formal hypothesis test to confirm that, but it would not be insightful. To describe the data using a statistical model, we will need to develop a model with non-constant variance.
- B. The sample variance is evidently different in different time intervals. However, we should not conclude that the underlying data generating mechanism is non-stationary before making a formal statistical test of equality of variances between the time regions that have lower sample variance and the regions that have higher sample variance. Visual impressions without a formal hypothesis test can be deceptive.
- C. A model with randomly changing variance looks appropriate for these data. Since the variance for such a model is time-varying, the model must be non-stationary.
- D. A model with randomly changing variance looks appropriate for these data. Despite the variance for such a model being time-varying, the model is stationary.
- E. The sample variance is evidently different in different time intervals. An appropriate next step to investigate stationarity would be to plot the sample autocorrelation function for different intervals to see if the dependence between time points is also time-varying.

Solution. D.

This is a subtle question, so let's discuss each option. The plotted time series is a realization of a stationary model:

```

import numpy as np
from scipy.stats import norm

N = 1000
sd1 = np.ones(N)
events = (np.random.uniform(size=N) < 5/N).astype(float)
sigma = 20
amplitude = 10
filter_length = int(5 * sigma)
filter_seq = np.linspace(-2.5*sigma, 2.5*sigma, filter_length)
filter_weights = norm.pdf(filter_seq, scale=sigma) * sigma * amplitude
sd2 = sd1 + np.convolve(events, filter_weights, mode='same')
Y = np.random.normal(0, sd2)

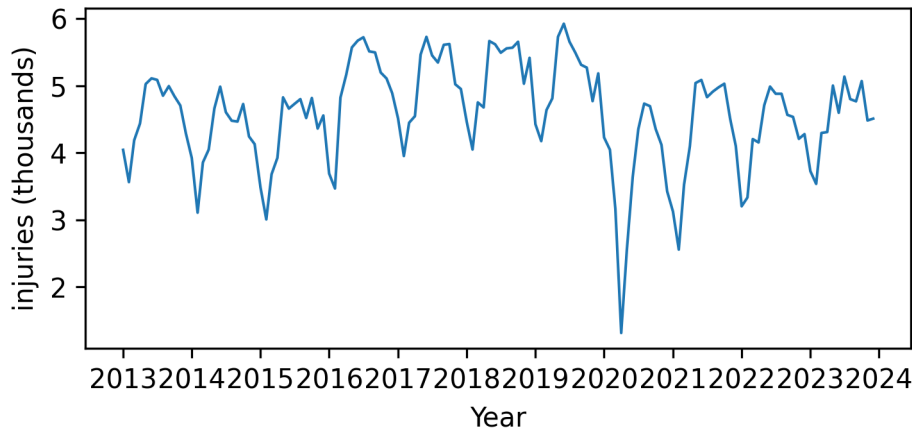
```

Hopefully, this suggests that it should not be clearly non-stationarity, ruling out A.

C and D contain a value judgement, “looks appropriate” which is hard to quantify but is (in this case) correct! “Randomly changing variance” is an informal description of a model with stochastic conditional variance. The sample variance estimates the variance conditional on the realization of the conditional variance. The actual variance is an expectation over possible values of the conditional variance. So, between C and D, only D can be correct.

B and E acknowledge the variation in sample variance but do not provide useful ways to assess whether this variable sample variance comes about via a stationary stochastic conditional variance model or via a non-stationary model. In particular, if you follow the advice in B you would conclude that an appropriate model should be non-stationary, which would be incorrect in this case. In financial applications, it is common to fit stationary models to time series of financial returns that often resemble this model. For this particular case, E would not show significant autocorrelation in any time interval, but the same reasoning applies.

Q1-03.



Above are monthly injuries from motor vehicle collisions in New York City. An augmented Dickey-Fuller test, `adfuller(injuries)`, gives a p-value of 0.014528. Which is the best way to proceed:

A: The time plot indicates a non-constant mean function describing a major dip due to the COVID-19 pandemic and an increasing trend at other times. The ADF test does not support or refute that model.

B: The ADF test suggests the series is stationary, supporting a decision to fit a SARMA model.

C: The ADF test suggests the series is non-stationary; it should be differenced before fitting a SARMA.

D: The ADF test indicates that the series is non-stationary, supporting the use of a non-constant mean function to describe a major dip due to the COVID-19 pandemic and an increasing trend at other times.

Solution. A.

The ADF test has a null hypothesis of a unit root linear model and an alternative of a stationary linear model, so neither of these describes a nonlinear trend. Here, the role of the COVID-19 pandemic is large enough that it does not make much sense to build a model that omits it, or describes it as a large perturbation of a stationary process. The conventional interpretation of the ADF test is option B (we reject the unit root hypothesis, and so we are invited to fit a stationary model). Here, a nonlinear trend (which is neither a unit root nor a stationary model) makes sense.

Q2. Calculations for ARMA models

Q2-01.

Let $Y_n = \phi Y_{n-1} + \epsilon_n$ for $n = 1, 2, \dots$ with $\epsilon_n \sim \text{iid}N[0, \sigma^2]$ and $Y_0 = 0$. The covariance of Y_n with Y_{n+k} for $k \geq 0$ is

- A. $\sigma^2 \phi^k / (1 - \phi^2)$
- B. $\sigma^2 \phi^{2k} / (1 - \phi^2)$
- C. $\sigma^2 \phi^k / (1 - \phi)$
- D. $\sigma^2 \phi^{2k} / (1 - \phi)$
- E. None of the above.

Solution. E.

This model is not started in its stationary distribution, leading to a covariance that is not shift invariant. The exact calculation is not needed, but it is as follows.

$$\begin{aligned}\text{Cov}(Y_n, Y_k) &= \text{Cov}\left(\sum_{i=1}^n \phi^{n-i} \epsilon_i, \sum_{j=1}^{n+k} \phi^{n+k-j} \epsilon_j\right) \\ &= \sum_{i=1}^n \sum_{j=1}^{n+k} \phi^{2n+k-(i+j)} \text{Cov}(\epsilon_i, \epsilon_j) \\ &= \sigma^2 \phi^k \sum_{i=1}^n \phi^{2(n-i)} \\ &= \sigma^2 \phi^k \sum_{i=0}^{n-1} \phi^{2i} \\ &= \sigma^2 \phi^k (1 - \phi^{2n}) / (1 - \phi^2)\end{aligned}$$

If you didn't see this, you may feel you were tricked. However, it is a common mistake in practical data analysis to pay insufficient attention to initial conditions, so it is worth bringing this to your attention.

Q2-02.

Let Y_n be an ARMA model solving the difference equation

$$Y_n = (1/4)Y_{n-2} + \epsilon_n + (1/2)\epsilon_{n-1}.$$

This is equivalent to which of the following:

- A. $Y_n = (1/2)Y_{n-1} + \epsilon_n$
- B. $Y_n = -(1/2)Y_{n-1} + \epsilon_n$
- C. $Y_n = (1/2)Y_{n-2} - (1/16)Y_{n-4} + \epsilon_n + \epsilon_{n-1} + (1/4)\epsilon_{n-2}$
- D. $Y_n = -(1/2)Y_{n-2} - (1/16)Y_{n-4} + \epsilon_n + \epsilon_{n-1} + (1/4)\epsilon_{n-2}$
- E. None of the above

Solution. A.

Writing the model in terms of the lag operator, L , we get

$$(1 - (1/2)L)(1 + (1/2)L)Y_n = (1 + (1/2)L)\epsilon_n.$$

Canceling out a factor of $(1 + (1/2)L)$, we obtain

$$(1 - (1/2)L)Y_n = \epsilon_n.$$

Q2-03.

Is it possible for an $AR(2)$ model to have a finite moving average representation, so that it is equivalent to some $MA(q)$ model for $q < \infty$?

A. No. Any moving average representation of any $AR(2)$ model is $MA(\infty)$

B. Yes. Although it is not true for any $AR(2)$ process, it is possible to find particular choices of the autoregressive coefficients, p_1 and p_2 , that lead to a finite $MA(q)$ representation.

C. It is not possible for any real-valued p_1 and p_2 , but it is possible if you permit p_1 and p_2 to be complex-valued.

Solution. A.

For any $AR(p)$ model with $p \geq 1$, an $MA(q)$ representation always has $q = \infty$. One way to see this is to argue by contradiction, by supposing there is a value of $q < \infty$. Then, setting $\phi(x)$ and $\psi(x)$ as the AR and MA polynomials, we can write

$$\frac{1}{\phi(B)} = \psi(B)$$

Thus, $\phi(x)\psi(x) = 1$. But $\phi(x)\psi(x)$ has a nonzero x^{p+q} coefficient so it cannot equal 1. This argument applies for either real-valued or complex-valued coefficients.

Q2-04.

Different criteria for selecting a time series model include (i) Akaike's Information Criterion; (ii) leave-one-out cross-validation; (iii) out-of-sample k -step-ahead prediction error; (iv) holding out the most recent 20% of the data for testing, while fitting to the first 80%. Suppose our goal is to make predictions for a collection of forecasting windows, and that the model fits the data fairly well. Which of the following are correct:

A. AIC is based on one-step prediction, so for longer-term forecasts it is less reliable than fitting by k -step prediction error for $k > 1$.

B. Leave-one-out cross-validation is not designed for dependent data.

C. Ideally, we should use a different model for each time in the forecasting window, fitting using k -step prediction error when forecasting k steps ahead.

D. The model that best predicts the most recent data when fitted to earlier data (i.e., (iv)) is more reliable for subsequent forecasting than methods which evaluate based on the whole time series history (i, ii, iii).

E. More than 1 of (A, B, C, D)

Note: You can suppose that the time series model is ARMA, but that is unimportant. All we need is that the model has a likelihood that can be computed, and a conditional expectation that can be evaluated to give a prediction rule.

Solution. B.

Leave-one-out cross-validation considers the situation where we predict time t given all the data both before and after t . That is a different problem from the forecasting problem of predicting at t given previous data. Without time dependence, the problems are equivalent, but time series data generally have time dependence.

D errs because the held-out test set is used to select the model but not to estimate its parameters. Losing access to the most recent data for model fitting may lead to inaccuracy. Also, using only 20% of the data to select the model may be problematic unless there is an abundance of data; a poor model could be good by chance on a small time interval.

A and C suffer from a similar error. If the model is not too grossly violated, the theoretical large-sample optimality properties of maximum likelihood for parameter estimation and AIC for model selection become relevant. Although the likelihood has a one-step factorization, it concerns the full joint distribution and so there is nothing to gain (and something to be lost) by using other criteria such as k -step prediction.

The question asserts that the model class fits reasonably well. If the model is a poor fit, time may be better spend looking for a better model rather than using a different criterion for an inferior model.

Q3. Likelihood-based inference for ARMA models

Q3-01.

The following table of AIC values results from fitting ARMA(p,q) models to a time series $y_{1:415}$ where y_n is the time, in milliseconds, between the n th and $(n + 1)$ th firing event for a monkey neuron. The experimental details are irrelevant here. You are asked to check how many adjacent pairs of AIC values in this table are inconsistent, such that they could mathematically arise only from a numerical error? Adjacent pairs of models are those directly above or below or left or right of each other in the table.

	MA0	MA1	MA2	MA3
AR0	3966.0	3961.5	3962.7	3964.7
AR1	3961.1	3962.6	3964.6	3966.6
AR2	3962.7	3960.5	3959.8	3961.7
AR3	3964.6	3965.5	3962.6	3968.4

A: 0, so the table is mathematically plausible.

B: 1

C: 2

D: 3

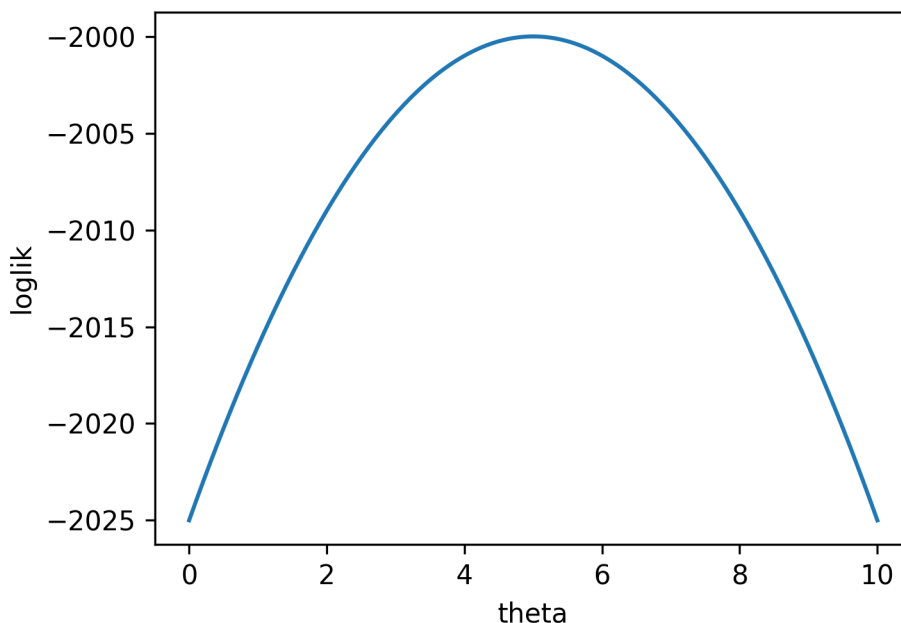
E: 4 or more

Solution. E.

Adding one parameter in a nested model cannot decrease the maximized log-likelihood, so it can increase the AIC by at most 2 units. Adjacent pairs $\{(p, q), (p', q')\}$ inconsistent with this are $\{(3, 2), (3, 3)\}$, $\{(2, 3), (3, 3)\}$, $\{(2, 1), (3, 1)\}$, $\{(2, 2), (3, 2)\}$.

Q3-02.

The Python function `statsmodels.tsa.arima.model.ARIMA.fit` and the R function `arima()` provide standard errors calculated by observed Fisher information. This question tests your understanding of what that means. Suppose a parametric model has a single parameter, θ , and the log-likelihood function when fitting this model to dataset is as follows:



What is the observed Fisher information (I_{obs}) for θ ?

Hint 1. The observed Fisher information is accumulated over the whole dataset, not calculated per observation, so we don't have to know the number of observations, N .

Hint 2. Observations in time series models are usually not independent, so the log-likelihood is not the sum of the log-likelihoods for each observation. Its calculation will involve consideration of the dependence, and usually the job of calculating the log-likelihood is left to a computer.

Hint 3. The usual variance estimate for the maximum likelihood estimate, $\hat{\theta}$, is $\text{Var}(\hat{\theta}) \approx 1/I_{obs}$.

A: $I_{obs} = 2$

B: $I_{obs} = 1$

C: $I_{obs} = 1/2$

D: $I_{obs} = 1/4$

E: None of the above

Solution. A.

The log-likelihood here is a quadratic function. We can see by inspection that this quadratic is given by

$$\ell(\theta) = -2000 - (\theta - 5)^2.$$

The observed Fisher information is the negative of the second derivative of the log-likelihood at the MLE, so $I_{obs} = 2$. Thus, the standard error is $1/\sqrt{2} = 0.707$

Q3-03.

```

SARIMAX Results
=====
Dep. Variable:          y      No. Observations:      165
Model:                ARIMA(2, 0, 1)  Log Likelihood      21.419
Date:                Thu, 12 Feb 2026  AIC              -32.838
Time:                17:16:02  BIC              -17.308
Sample:              0      HQIC              -26.534
                        - 165
Covariance Type:      opg
=====

```

```

SARIMAX Results
=====
Dep. Variable:          y      No. Observations:      165
Model:                ARIMA(2, 0, 2)  Log Likelihood      22.626
Date:                Thu, 12 Feb 2026  AIC              -33.251
Time:                17:16:02  BIC              -14.616
Sample:              0      HQIC              -25.686
                        - 165
Covariance Type:      opg
=====

```

The Python output above uses `ARIMA` from `statsmodels` to fit `ARMA(2,1)` and `ARMA(2,2)` models to the January level (in meters above sea level) of Lake Huron from 1860 to 2024. Residual diagnostics (not shown) show no major violation of model assumptions. We aim to choose one of these as a null hypothesis of no trend for later comparison with models including a trend.

Which is the best conclusion from the available evidence:

A: The `ARMA(2,2)` model has a lower AIC so it should be preferred.

B: We cannot reject the null hypothesis of `ARMA(2,1)` since the `ARMA(2,2)` model has a likelihood less than 1.92 log units higher than `ARMA(2,1)`. Since there is not sufficient evidence to the contrary, it is better to select the simpler `ARMA(2,1)` model.

C: Since the comparison of AIC values and the likelihood ratio test come to different conclusions in this case, it is more-or-less equally reasonable to use either model.

D: When the results are borderline, numerical errors in the `stats::arima` optimization may become relevant. We should check using optimization searches from multiple starting points in parameter space, for example, using `arima2::arima`.

Solution. D.

All the answers are fairly reasonable here! Perhaps the most unreasonable thing would be to be sure there's only one reasonable answer.

However, a more careful optimization using multiple starting values (e.g., the R package `arma2::arima`) shows us that ARMA(2,1) actually has a higher AIC than ARMA(2,2) so all lines of evidence suggest ARMA(2,1) is a better choice. The differences are small, so the choice is unlikely to be highly consequential.

```
#| label: q3_03_arma2
#| echo: false
dat <- read.table(file="data/huron_level.csv",sep="," ,header=TRUE)
huron_level <- dat$Jan
set.seed(28)
arma2.2.1 <- arma2::arima(huron_level,order=c(2,0,1),
  max_iters=200,max_repeats=20)
arma2.2.2 <- arma2::arima(huron_level,order=c(2,0,2),
  max_iters=200,max_repeats=20)
arma2.2.1
arma2.2.2
```

Sometimes the multiple starts used by `arma2::arima` make a difference, sometimes the results from `stats::arima` are unchanged. In this case, it happens to make a difference.

Q3-04.

Suppose model M_0 is nested within a larger model M_1 which has one additional parameter. Suppose that the AIC for M_1 is 0.5 units lower than the AIC for M_0 . Which of the following is a correct expression for the p-value of a likelihood ratio test for M_1 against the null hypothesis M_0 , supposing that a Wilks approximation is accurate? Here, χ_1^2 is a chi-square random variable on 1 degree of freedom.

A: $P(\chi_1^2 > 0.5)$

B: $P(\chi_1^2 > 1)$

C: $P(\chi_1^2 > 1.5)$

D: $P(\chi_1^2 > 2)$

E: $P(\chi_1^2 > 2.5)$

F: $P(\chi_1^2 > 3)$

Solution. E.

The AIC for each model $k \in \{0, 1\}$ is $AIC_k = -2\ell_k + 2D_k$, where ℓ_k is the log-likelihood for

M_k and D_k is the number of parameters. Thus, $AIC_0 - AIC_1 = 2(\ell_1 - \ell_0) - 2 = 0.5$, and so $2(\ell_1 - \ell_0) = 2.5$. Under M_0 , according to Wilks' approximation,

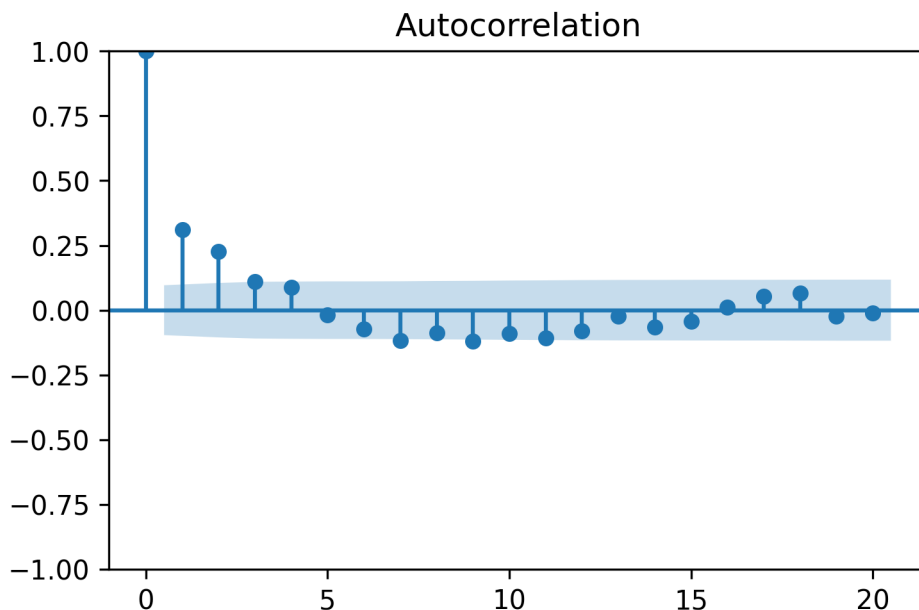
$$2(\ell_1 - \ell_0) \sim \chi_1^2.$$

Therefore, the p-value is $P(\chi_1^2 > 2.5) = 0.11$.

Q4. Interpreting diagnostics

Q4-01.

We consider data $y_{1:415}$ where y_n is the time, in milliseconds, between the n th and $(n+1)$ th firing event for a monkey neuron. Let $z_n = \log(y_n)$, with log being the natural logarithm. The sample autocorrelation function of $z_{1:415}$ is shown below.



We are interested about whether it is appropriate to model the time series as a stationary causal ARMA process. Which of the following is the best interpretation of the evidence from these plots:

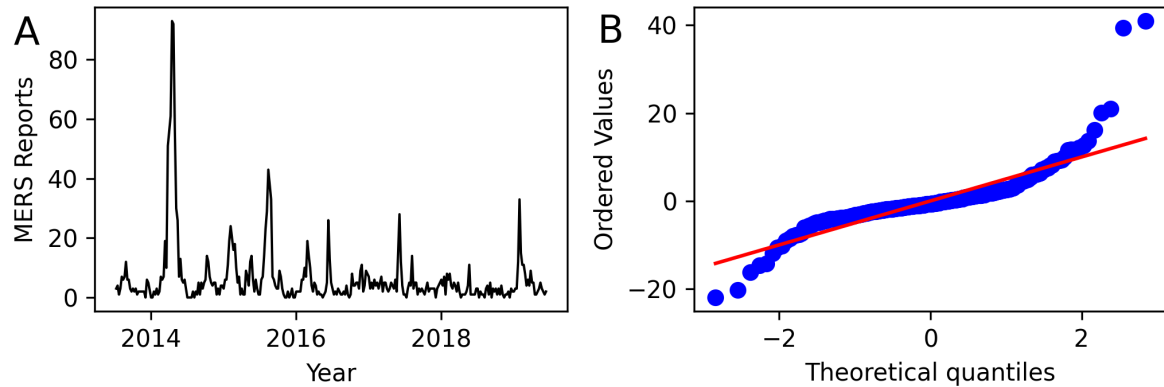
- A. There is clear evidence of a violation of stationarity. We should consider fitting a time series model, such as ARMA, and see if the residuals become stationary.
- B. This plot suggests there would be no benefit from detrending or differencing the time series before fitting a stationary ARMA model. It does not rule out a sample covariance that varies with time, which is incompatible with ARMA.
- C. This plot is enough evidence to demonstrate that a stationary model is reasonable. We should proceed to check for normality, and if the data are also not far from normally distributed then it is reasonable to fit an ARMA model by Gaussian maximum likelihood.

Solution. B.

The usual interpretation of the sample ACF assumes that variance and covariance depend only on lag (i.e., the data are well modeled by a shift-invariant autocovariance) but the plot does not check this.

There is clear evidence that the process is not well modeled by white noise. There is also clear evidence against a trend in the mean, which would show up as a slowly decaying sample ACF.

Q4-02.



(A) Weekly cases of Middle East Respiratory Syndrome (MERS) in Saudi Arabia. (B) a normal quantile plot of the residuals from fitting an ARMA(2,2) model to these data using `arma()`. What is the best interpretation of (B)?

A: We should consider fitting a long-tailed error distribution, such as the t distribution.

B: The model is missing seasonality, which could be critical in this situation.

C: For using ARMA methods, these data should be log-transformed to make a linear Gaussian approximation more appropriate.

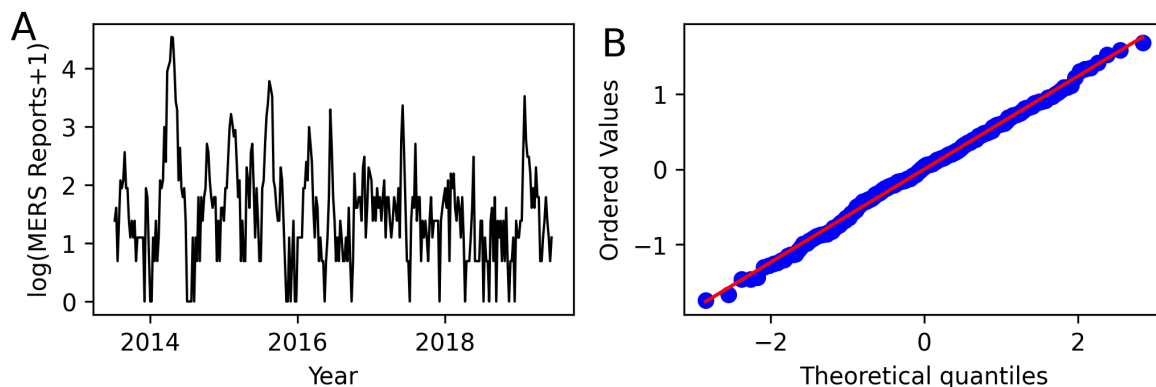
D: The normal quantile plot shows a long-tailed distribution, but this is not a major problem. We have over 300 data points, so the central limit theorem should hold for parameter estimates.

E: The normal quantile plot shows long tails, but with the right tail noticeably longer than the left tail. We should consider an asymmetric error distribution.

F: We should not interpret (B) before testing for stationarity. First make an ADF test and, if the null hypothesis is not rejected, recalculate (B) when fitting to the differenced data.

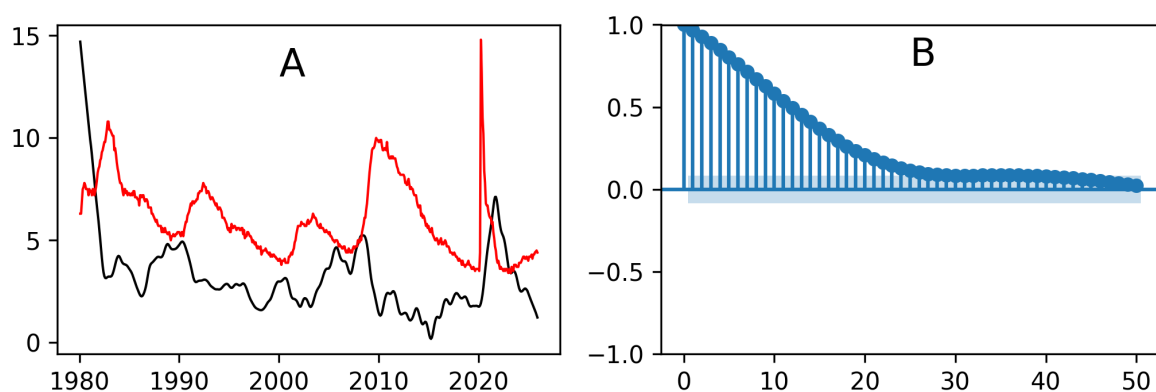
Solution. C.

Here's what happens when we take a $\log(x+1)$ transform, fitting ARMA(2,2) and checking the residuals as before.



It is often a good idea to log-transform non-negative quantities, and failure to do this can show up as long tailed residuals. Fitting long-tailed ARMA models is possible, but non-standard and not necessary here. There is seasonality, but an ARMA(2,2) model can already explain some periodicity so including a seasonal term in the model is not critical. There may be some non-stationarity here, but nothing that resembles the null hypothesis of the Augmented Dickey-Fuller test, so that is not relevant here.

Q4-03.



(A) Inflation (black) and unemployment (red) for the USA, 1980-2024. (B) Sample autocorrelation function of the residuals from a least square regression ($y=\text{inflation}$, $x=\text{unemployment}$), with estimated coefficients below. Which is the best interpretation of these graphs and fitted model?

	coef	std err	t	P> t	[0.025	0.975]
const	2.8109	0.293	9.577	0.000	2.234	3.387
x1	0.0724	0.047	1.550	0.122	-0.019	0.164

A: 0.072 is a reasonable estimate for the additional unemployment caused by one percentage point of additional inflation. We should not trust the uncertainty estimate, since our model does not allow for autocorrelation of the residuals.

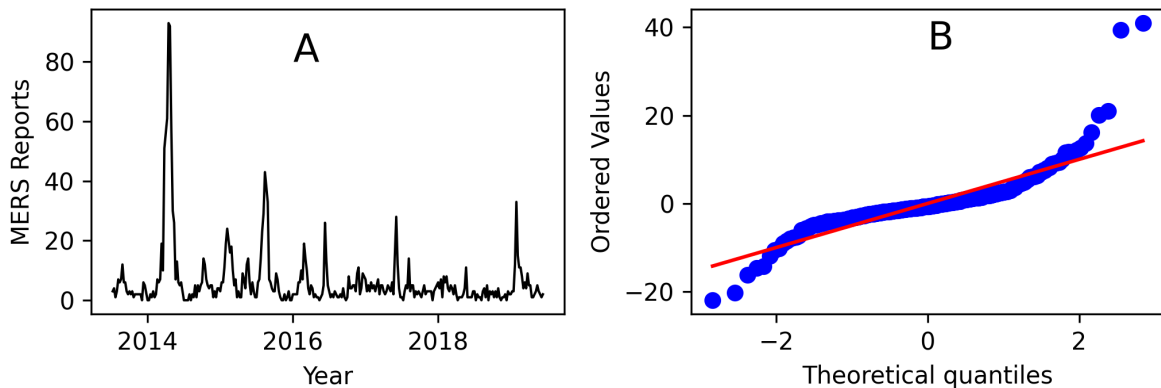
B: 0.072 is a reasonable estimate for the association between inflation and unemployment. We should not assume there is a causal relationship. We should not trust the uncertainty estimate, since our model does not allow for autocorrelation of the residuals.

C: 0.072 is a reasonable estimate for the association between inflation and unemployment. We should make an additional assumption that there are no confounding variables, and then we can interpret this association to be causal. We should not trust the uncertainty estimate, since our model does not allow for autocorrelation of the residuals.

Solution. B.

Association is not causation. Inflation and unemployment are two aspects of a complex system. They may have some direct effect on each other, but they are also both driven by other aspects of the economy such as interest rates, consumer spending, and international trade. Mathematically, we could wish away these confounding variables (as proposed in answer C) but for practical purposes it does not make sense to do so.

Q4-04.



(A) Weekly cases of Middle East Respiratory Syndrome (MERS) in Saudi Arabia. (B) a normal quantile residual plot for ARMA(2,2). We can formally test for non-normality of these residuals by a Shapiro-Wilk test ($p\text{-value}=4.8 \times 10^{-21}$). What best describes the value added by presenting the Shapiro-Wilk test here?

A. We should always be alert for the danger of seeing patterns in noise. The Shapiro-Wilk test is useful to confirm our assessment that the normal quantile plot shows long tails.

B. Presenting the Shapiro-Wilk test here is not very insightful here, since the long tails are obvious from the normal quantile plot. However, adding this test demonstrates technical competence so it is better to include it than to omit it.

C. The long tails are established from the normal quantile plot. We could consider a log

transform, or a long-tailed model, or a bootstrap simulation study to investigate whether the conclusions are sensitive to non-normality. Adding a fairly uninformative test instead of investigating the consequences of the error distribution could be a distraction from good data analysis.

D. The Shapiro-Wilk test is useful, but has the problem that it only tells us about lack of normality, not whether the non-normality is due to skew or kurtosis. We should supplement with a Jarque-Bera test to assess those.

Solution. C.

The key step after plotting B is to notice that the residuals have a long tail, especially to the right. This may remind us to try a log transform, which is successful here.

It is rather uninteresting to test a null hypothesis that is no longer plausible after plotting the data. If the normal quantile plot shows little deviation from normal, then it may be interesting to test whether that is enough to reject a Gaussian null, though in that case we probably do not have to modify the model (recall that statistical significance and practical significance can differ).

If it is not obvious to you that the deviation shown in the normal quantile plot is entirely incompatible with a Gaussian model, train your intuition by plotting some simulated normal quantile plots for data generated using `rnorm`.

If making useless tests is harmless, it might not matter whether we make a formal test for normality here. However, in practice, useless analysis distract from useful analysis.

The normal quantile plot has diagnostic value beyond simply rejecting the null of normality. For example, we can see whether there are outliers, and we can compare the left and right tails. Tests such as Shapiro-Wilk and Jarque Bera may be useful in some situations. However, they often add little to a normal quantile plot.

Q4-05.

	AIC MA0	AIC MA1	AIC MA2	AIC MA3	AIC MA4	LBT MA0	LBT MA1	LBT MA2	LBT MA3	LBT MA4
AR0	174.82	48.81	9.61	-14.62	-17.98	7.6e- 63	2.84e- 26	2.28e- 11	0.002	0.0409
AR1	-33.05	-34.44	-32.68	-30.69	-30.31	0.567	0.95	0.925	0.918	0.999
AR2	-34.07	-32.84	-33.25	-29.27	-28.51	0.951	0.926	0.98	0.954	0.996
AR3	-32.67	-33.2	-29.21	-29.91	-27.38	0.922	0.985	0.967	0.981	0.95
AR4	-30.77	-31.36	-30.11	-28.94	-24.92	0.92	0.967	0.989	0.978	0.994

The [Ljung-Box test \(LBT\)](#) provides an alternative approach to comparison of AIC values for selecting ARMA models. Whereas the standard sample autocorrelation function (ACF) residual plot tests each ACF component $\hat{\rho}_k$ under a null hypothesis of white noise, LBT tests

$\sum_{k=1}^h \hat{\rho}_k^2$. Here, we present an AIC table and an LBT table (for $h = 5$). This course have favored AIC, with visual inspection of ACF and checking whether residual patterns appear in the frequency domain. There may be reasons to prefer LBT. Which of the following are good reasons to use LBT?

- (i). LBT provides a p-value which is more formal than the comparison of AIC values.
- (ii). Numerical issues involved in fitting an ARMA model may cause problems for comparing AIC values.
- (iii). The LBT gives insights into what model to investigate next if the null hypothesis is rejected.
- (iv). The LBT is useful in conjunction with AIC and ACF, since it provides an alternative perspective.

- A. (i) only
- B. (i, ii, iv)
- C. (i,iii, iv)
- D. (ii, iii, iv)
- E. None of the above

Solution. E.

“None of the above” is correct in two senses, since none of the reasons proposed are strong. Perhaps LBT had more value in an era of less computational power, before maximum likelihood estimation became routine for ARMA models. LBT usually adds little or nothing to the methods covered in class. We consider each in turn.

AIC values are a formal quantitative measure, just like p-values. They measure different things, as explained in the notes. Using p-values for model diagnostics is actually informal, since a formal p-value should correspond to a hypothesis specified before examining the data. So, (i) is not a good answer.

It is correct that numerical optimization can be problematic for fitting ARMA models. However, both AIC and LBT assess the same fitted model and so they share any consequences of imperfect optimization. So, (ii) is not a good answer.

Null hypothesis tests have the feature that they are relatively weak at telling you what to do if the null hypothesis is rejected, especially when they test against a very general alternative as for LBT. AIC provides a ranking of the models under investigation (though there may be reasons not to proceed with the top ranked model according to AIC). By contrast, it is unclear from the LBT table which model to choose. We can see that AR(0) models are inappropriate, but that is also clear from AIC. So, (iii) is not a good answer.

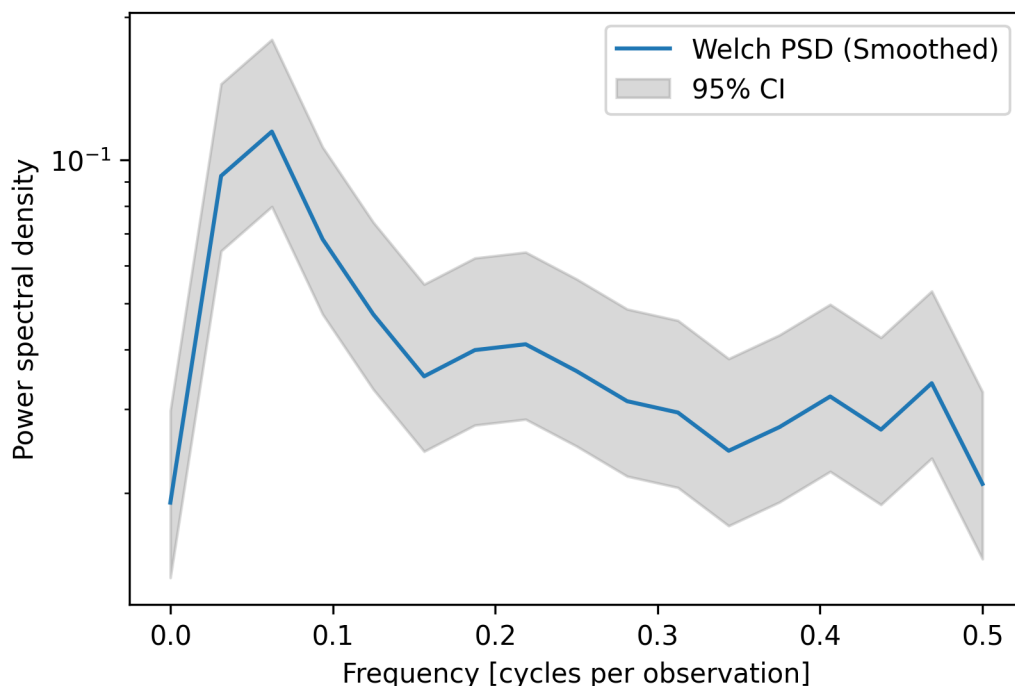
If LBT added anything substantial beyond AIC, it could be a useful extra component to the analysis. However, in this example, we see that every model with reasonable AIC values does not reject the LBT null. This is commonly the case; if there is substantial autocorrelation in the residuals then a larger ARMA model will have better AIC. It can be useful to supplement AIC with a likelihood ratio test, since this provides a complementary perspective: AIC asks

which model has better estimated predictive skill, whereas the null hypothesis test asks whether the simpler model is statistically plausible against the alternative of the more complex model. LBT does not address that. So, (iv) is not a good answer.

Q5. The frequency domain

Q5-01.

We consider data $y_{1:415}$ where y_n is the time interval, in milliseconds, between the n th and $(n + 1)$ th firing event for a monkey neuron. Let $z_n = \log(y_n)$, with \log being the natural logarithm. A smoothed periodogram of $z_{1:415}$ is shown below. Units of frequency are the default value in R, i.e., cycles per unit observation. We see a peak at a frequency of approximately 0.07.



Which if the following is the best inference from this figure

- A. Transitions between rapid neuron firing (short intervals between firing) and slow neuron firing (long intervals between firing) occur every $1/0.07 \approx 14$ firing events.
- B. The neuron has a characteristic duration between firing events of $1/0.07 \approx 14$ milliseconds.
- C. The neuron has a characteristic duration between firing events of $1/\exp(0.07) \approx 0.9$ milliseconds.

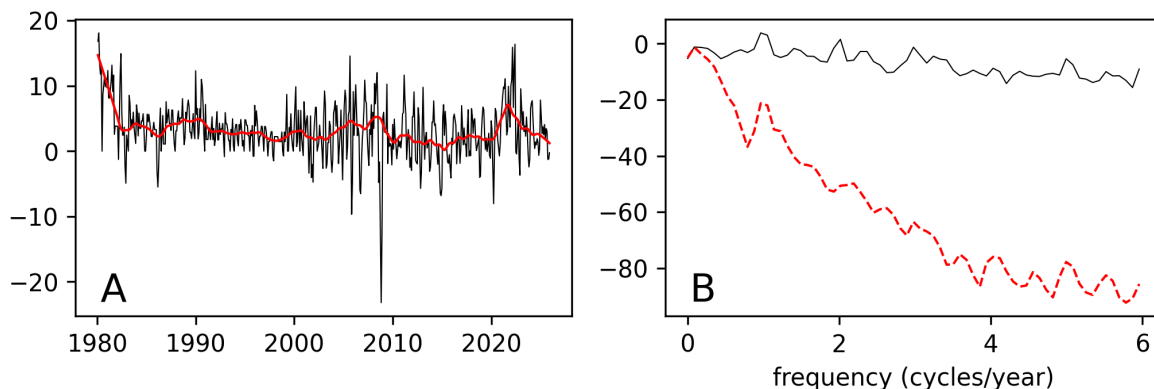
Solution. A.

In this example, the units of the “time” variable n are a dimensionless count of firing events. So, the units of frequency are cycles per firing event. The units of period are firing events per cycle. The peak in the smoothed periodogram corresponds to a cycle length of about

$1/0.07 \approx 14$ observations. These cycles correspond to oscillations between low firing rate (i.e., long time intervals between firing) and high firing rate (short time intervals between firing).

The units of y_n are time, in milliseconds. It is unusual, and therefore somewhat counter-intuitive, to have a time series with units of time for the measured variable and a dimensionless “time” variable.

Q5-02.



The monthly US consumer price index (CPI) combines the price of a basket of products, such as eggs and bread and gasoline. (A) Annualized monthly percent inflation, i.e., the difference of log-CPI multiplied by 12×100 (black line); a smooth estimate via local linear regression (red line). (B) The periodogram of inflation and its smooth estimate. Which best characterizes the behavior of the smoother?

- A: Cycles longer than 2 months are removed
- B: Cycles shorter than 2 months are removed
- C: Cycles longer than 2 year are removed
- D: Cycles shorter than 2 year are removed
- E: Cycles longer than $(1/2)$ year are removed
- F: Cycles shorter than $(1/2)$ year are removed

Solution. D.

The units are in cycles per year; if we doubt this, we can tell because the largest peak at 1 corresponds to annual seasonality. At about $(1/2)$ cycle per year, the power in the smooth fit (the red dashed line) rapidly falls several log units below the power in the raw monthly inflation data. $(1/2)$ cycle per year corresponds to cycles of 2 year.

Q6. Scholarship for time series projects

Q6-01.

This question on citing references applies to any statistics report, but it is particularly relevant here since we are learning proper use of sources in order to write open-access midterm and final projects.

Suppose that the midterm project P1 cites a past project, P2, in the reference list. P1 references P2 at one point, mentioning that the projects have similarities. When you look at the source code and the writing, you find various points where P1 and P2 are almost identical, though at other points the projects are entirely different. What do you infer?

A: The authors of P1 have done enough to honestly disclose the relationship with P2. After all, there is sufficient information provided for any reader to track down the exact relationship.

B: The authors of P1 have misrepresented the relationship with P2 by appearing to take credit for some original work which was in fact heavily dependent on a source. This is a serious offence which should be reported to Rackham and/or the Associate Chair for Graduate Programs in Statistics as a violation of academic integrity.

C: There is not enough information to tell the actual story for certain. The authors of P1 may or may not have done something wrong, depending on information that is not available to us, but they did cite P2 so they should be given the benefit of the doubt and should not lose any scholarship points.

D: The authors of P1 have misrepresented the relationship with P2 by appearing to take credit for some original work which was in fact heavily dependent on a source. This is a moderately severe offence, partly offset by including P2 in the reference list. A substantial number of scholarship points should be subtracted.

E: P1 evidently has not shown perfect scholarship, but this is a small issue that could easily be an honest mistake given that the authors were not trying to hide the fact that they had studied P2. It is appropriate to subtract, say, 1 point for scholarship for this mistake.

Solution. D.

If occurrences such as this were reported to the authorities, this would fill up too much time for the deans and chairs. This is a fairly serious issue, and a responsible student should not usually do this by mistake. It wastes everyone's time if proper credit is not clearly assigned to sources and if the grader has to track down the contribution of the authors.

Q6-02.

Four people in a team collaborate on a project. After the project is submitted, a reader identifies that part of the project is adapted from an unreferenced source, i.e., it has been plagiarized. The team worked using git and cooperates on tracking down the issue, and the

commit history clearly reveals who wrote the problematic part of the project. What is the most appropriate course of action:

- A. The guilty coauthor should be penalized heavily for poor scholarship, and the other coauthors should have a minor penalty for failing to check their colleague's work.
- B. All coauthors should share the same penalty, since this is a team project and all coauthors share equal responsibility for the submitted report.
- C. The guilty coauthor should be penalized heavily for poor scholarship. The other coauthors have demonstrated strong scholarship by following good transparent working practices that enabled this issue to get quickly resolved, so they should not receive any penalty.
- D. It is necessary to collect more information before coming to a decision. For example, the team may argue that the source is well known to all readers so did not have to be cited.

Solution. A.

Everyone should take responsibility for checking work submitted under their name. However, when it is possible to isolate the misconduct to one individual, and the rest of the team has demonstrated strong scholarship, they have some protection from the heavy scholarship consequences for the serious error.

One can always propose collecting new information, as in D, but in practice you often have to act on the information at hand. Parties can request a regrade if they think the decision is too far from justice.

Q6-03.

You discover that your team-mate is using Google Translate to carry out their share of the writing. The translation looks poorly done, similar in quality to ChatGPT, and does not use technical time series terminology correctly. What is the best course of action among the options below

- A. Alert the instructor that you have a team mate adopting questionable scholarship strategies, in order to make sure you are not personally held responsible.
- B. Ask ChatGPT to rewrite this problematic section to improve its quality
- C. Help your team mate to rewrite the section in their own voice (shared with your voice).

Solution. C.

Different team mates bring different skills to the project, and perhaps you are more fluent at writing in English than one of your team mates. Major team problems can be discussed with the instructor, but this issue is best caught and corrected early and solved within the team.

Q6-04.

Why is it helpful for a course such as DATASCI/STATS 531, that permits the use of internet resources including GenAI and past solutions, to require students to say explicitly say when they do not use sources?

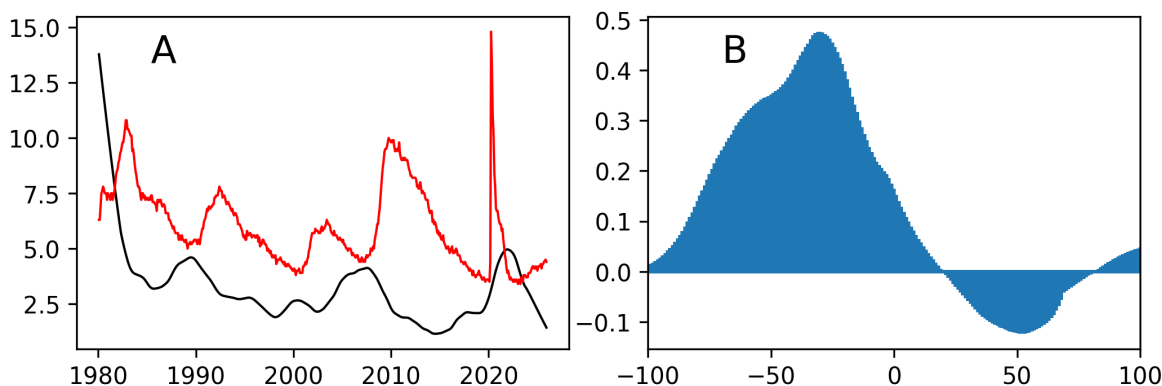
- A. Failure to give credit to sources is against the academic integrity rules of Rackham, the graduate school at University of Michigan.
- B. It helps the GSI to grade the homework when they know exactly what sources have been used and for what question.
- C. Students whose solution is more dependent on sources than they want to admit are reluctant to explicitly deny using sources.
- D. The GSI has the task of evaluating whether the student has demonstrated thought about the homework task beyond collecting material from sources into a solution. This is not an easy task even when the sources are clearly listed and referenced at the point (or points) where they are used.

Solution. C.

The points made in A, B, D are all correct but are not directly relevant to the question. In an ideal world, the absence of a list of sources would be logically equivalent to an explicit statement saying that no sources were used. However, in practice that is not the case. We need a system that is robust against the natural tendency to hide information that could lead us to get a lower grade (because the grader would be able to see that our own contribution was smaller than it might otherwise appear). Most of us realize that explicitly saying we did not consult a source, when in fact we did, is a lie and amounts to academic misconduct. Failure to mention the source sounds like a milder misdemeanor. Therefore, to help the grader distinguish these things, it is important to be explicit about the lack of sources when indeed you did not need to consult any. It is appropriate for the GSI to award points for turning in solutions that are well-written and easier for the grader to evaluate.

Q7. Data analysis

Q7-01.



(A) Inflation (black) and unemployment (red) for the USA, 1980-2024. (B) Cross-correlation function, `plt.xcorr(inflation,unemployment)`. What is the best interpretation of this plot?

A: High inflation generally led high unemployment, with a lag of about 4 yr.

B: High inflation generally followed high unemployment, with a lag of about 4 yr.

C: Association is not causation, so we should not interpret a cross-correlation plot in terms of lead and lag relationships.

Solution. A.

The cross-correlation, $\rho_{XY}(h)$, is the correlation between X_{n+h} and Y_n where $X_{1:N}$ and $Y_{1:N}$ are jointly stationary. The sample cross-correlation is the standard estimator of this. The large peak at $h = -48\text{month} = -4\text{yr}$ means that X_n can predict Y_{n+h} , so X_n leads Y_n .

For answer C, lead and lag relationships are generally used to refer to associations. The cross-correlation plot is a measure of lagged association, and we should not give it a causal interpretation without carefully explaining the justification.

Q7-02.

Which of the below best explains the role of models in time series analysis (and statistics more broadly)? Pick one, and explain your choice.

A. A statistical inference requires a statistical model. Investigating a hypothesis via a p-value requires a model to obtain the distribution of the p-value. Investigation of Bayesian posterior probabilities requires a model.

B. Parametric statistical tests require a model. However, where possible, we should use non-parametric methods (e.g., rank tests, or cross-validation) that do not require a model and so are more robust.

C. Modeling is of intrinsic scientific value, as a way of understanding the data-generating mechanism. We do not have to have specific statistical tests in mind when developing a model.

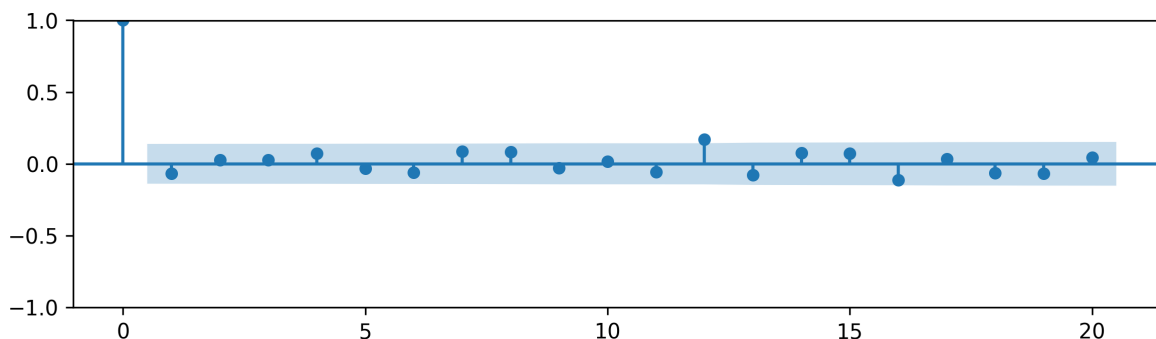
Solution. A.

So-called nonparametric tests always depend critically on some modeling assumption, often an assumption of independence. Nonparametric models are flexible classes of models that can be expressed via infinite-dimensional parameter spaces rather than the finite-dimensional parameter spaces of regular parametric models. Sometimes the model assumptions are not so evident, or not so well-known, for nonparametric models, but they must be there in order to obtain inferences such as p-values or posterior probabilities. This discredits answer B.

Models have many purposes, not all statistical. Simple conceptual models are widely used in science (e.g., ideal springs, biochemical interactions within a cell) but these are not usually designed to give quantitative statistical understanding of data. A statistical model is one that is developed for the purposes of statistical inferences, and the model is useful so far as those inferences are reliable and pertinent to the scientific study. This discredits answer C.

We are left with answer A. This is the traditional framework for the introductory theory and practice of statistics. It remains true in more complex time series situations, for which we need models describing dependence through time.

Q7-03.



The plot above is the sample autocorrelation function (ACF) for a time series $y_{1:N}$. What is the best conclusion to draw from this evidence?

A. The sample ACF is statistically consistent with an iid model. Therefore, standard statistical reasoning lets us conclude that the time series is iid.

B. The time series passed this test for being iid, but it might fail other tests. In practice, we can never make all possible tests so we cannot reliably conclude that the time series is iid.

C. It is meaningless to say that a time series is iid, since a time series is a sequence of numbers and iid is a property of a sequence of random variables.

D. The statistical evidence in the sample ACF is consistent with using an iid model to describe the data.

E. There is an oscillating pattern in the sample ACF. Even though no individual lag contradicts an iid model assumption at the 5% level, the overall pattern is clear evidence against iid.

Solution. D.

First, note that C is a correct statement, though it is not a conclusion drawn from the data so it cannot be the best answer to the question.

Both A and B suppose that it would make sense to say that the time series is iid. This is wrong, as pointed out by C.

Note that some time series authors (e.g., the textbook by Huang & Petukhina, 2022) attempt to avoid this difficulty by asserting that a time series is a collection of jointly distributed random variables, which in this course we call a time series model. However, Huang & Petukhina then refer to datasets as “time series” which leads to the same inconsistency by a different route.

D explains what can be concluded by a statistical test for an iid model, no more and no less.

For E, visual inspection might suggest a wave pattern in the residuals. Because the data here are simulated iid Normal, any pattern your eye might see is spurious.

License: This material is provided under a [Creative Commons license](#)