

TIME SERIES

GOAL

the goal of time series data is to analyze trends in values that occurred at certain time intervals. This usually means that you want to predict future values or find the period of the data you have.

DEFINITIONS

Stationary: the time series is invariant to shifts in time, it is oscillating around some constant mean

Auto correlation: a measure of how related a time series is to itself at different lags

White noise: a time series that is random w/
constant mean & variance, it has no autocorrelation

Lag: how many samples of a difference from the one you're looking at. If measurements are taken daily a lag of 2 would be the measurement from 2 days ago.

Even: taken at evenly spaced times.

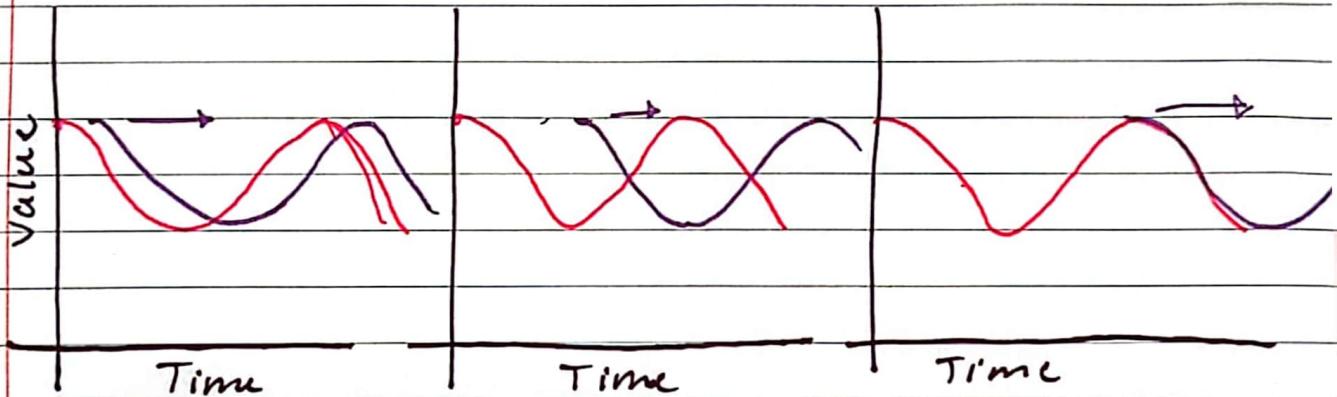
Uneven: taken at varying time intervals → many of the methods do not work for uneven data.

Methods

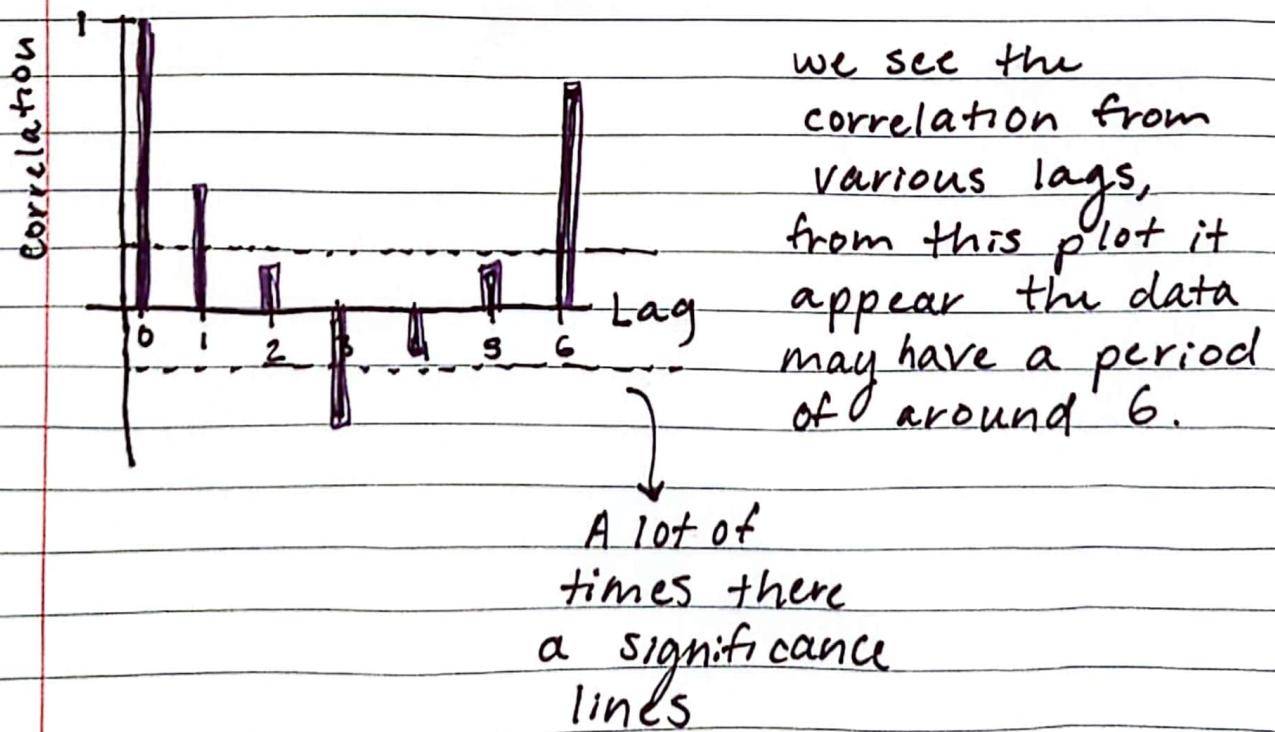
AutoCorrelation Function (ACF)

used for evenly spaced data to see if there is some sort of period or time correlation

The correlation at different lags is calculated



The purple moves over the actual data at lags that are equal to the time spacing. This results in a graph that looks like:



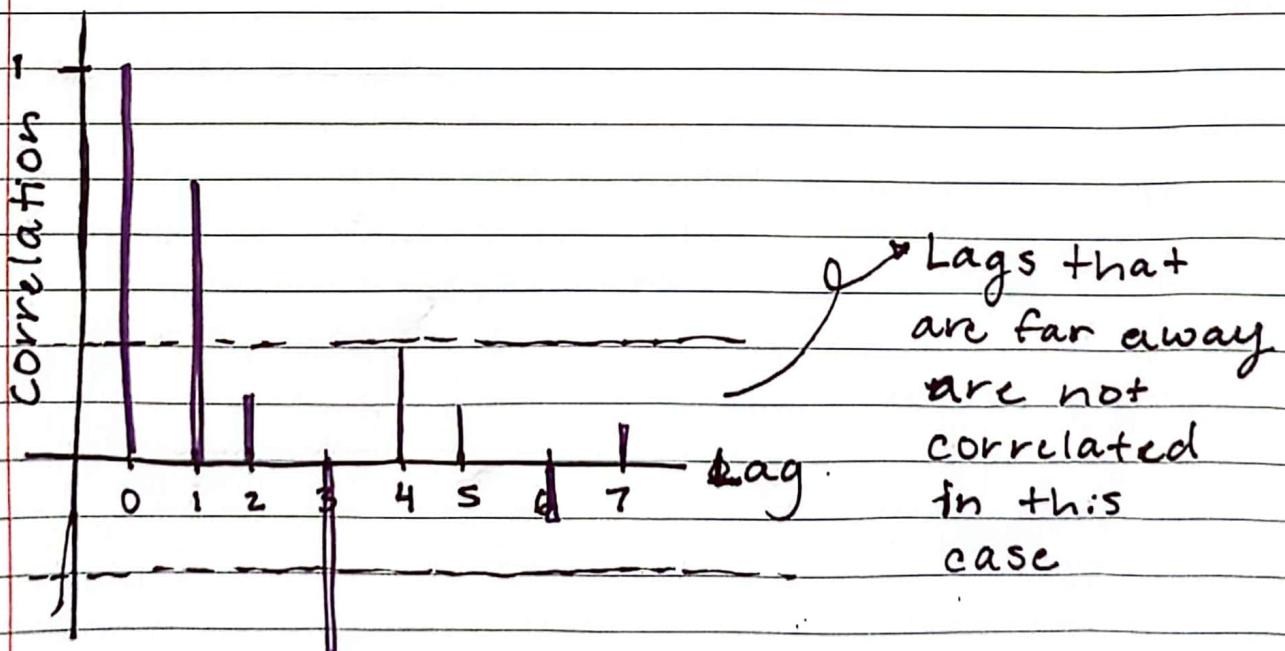
Partial Autocorrelation Function (PACF)

This is similar to ACF, but it gives the auto correlation at a value of k by removing the effects of shorter k s.

This can help to better identify the underlying structure \rightarrow the direct relationship between lags.

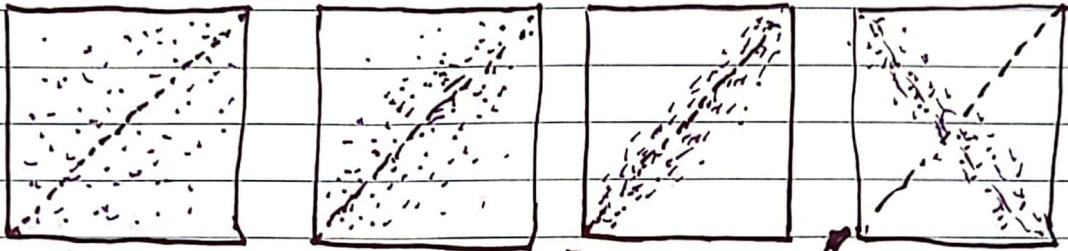
This is done by fitting a regressive model that accounts for each lag.

$$PACF(2) = \frac{ACF(2) - ACF(1)^2}{1 - ACF(1)^2}$$



Lag K Scatter Plot

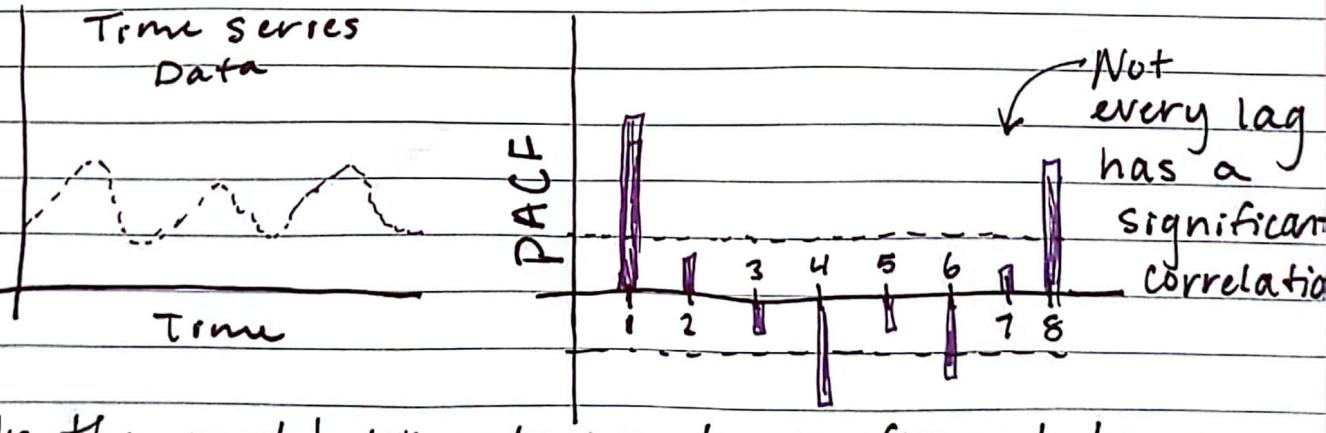
these plots plot the values of different lags
→ x_t vs x_{t+k} . This helps to visually determine patterns.



No correlation some + correlation strongly + correlation strong correlation

Autoregressive Model (AR)

Tries to predict future values based on past values (still for evenly spaced)



In the model we only want significant lags so it will look something like:

$$m_t = \beta_0 + \beta_1 m_{t-1} + \beta_2 m_{t-2} + \beta_6 m_{t-6} + \beta_8 m_{t-8}$$

Moving Average Model (MA)

Instead of basing your guess on the actual past values, MA bases it on the error. There's some value you're trying to predict, but at each time you're probably off by a little bit.

$$X_i = \varepsilon_i + \beta_1 \varepsilon_{i-1} + \dots + \beta_q \varepsilon_{i-q}$$

These are
the errors at
different lags

So if you're really far off a one lag, you'll want to change your answer more the next time.

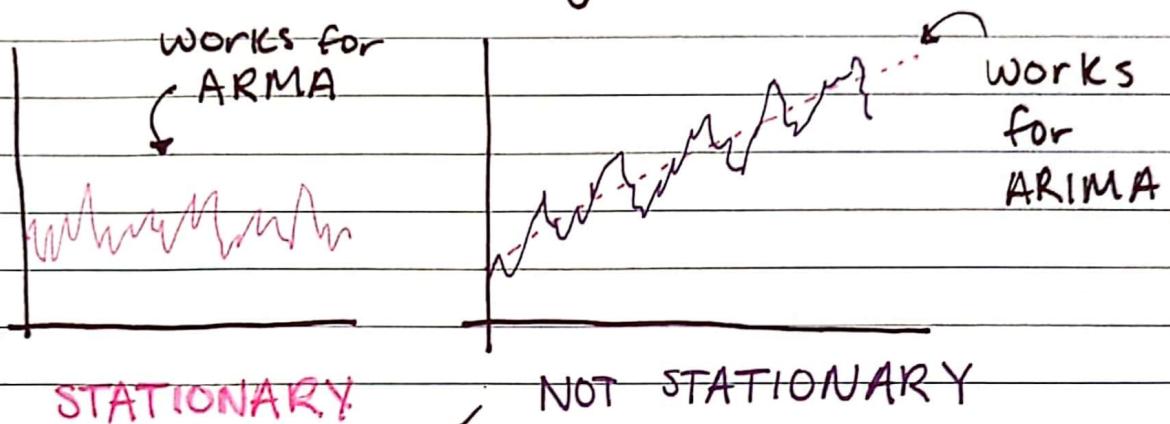
ARMA (Autoregressive - Moving Average)

This model combines AR & MA into one. It uses both the past values & the errors to make future predictions.

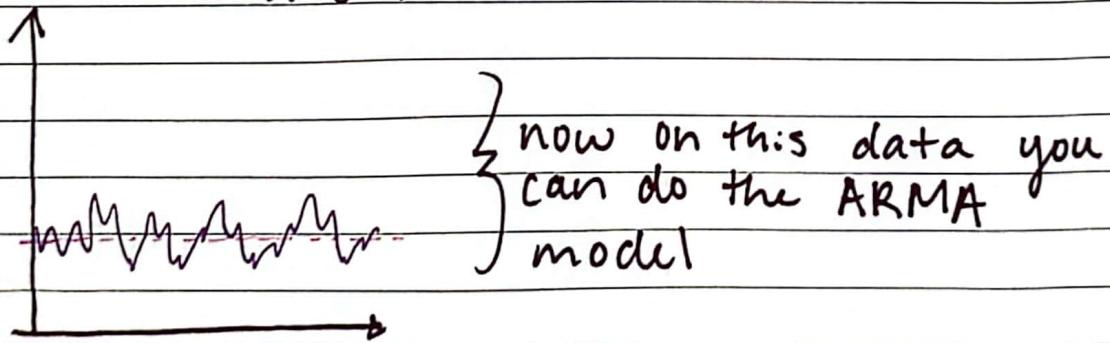
$$X_i = \underbrace{\alpha_1 X_{i-1} + \dots + \alpha_p X_{i-p}}_{\text{Autoregressive}} + \varepsilon + \underbrace{\beta_1 \varepsilon_{i-1} + \dots + \beta_q \varepsilon_{i-q}}_{\text{Moving Average}}$$

ARIMA Model (Autoregressive integrated moving average)

Still used for predicting the future, with evenly spaced data, but this time the data does not have to be stationary.



This data is transformed by take $x_{t+1} - x_t$ to get the "average" change which is linear



Three parameters

p: autoregressive Order

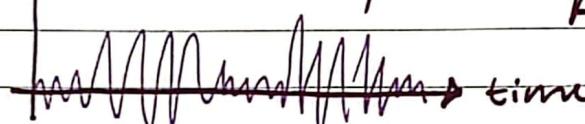
d: differencing to account for the change in the mean

q: order for MA model

ARCH Model (autoregressive conditional heteroscedastic)

This model is used to better understand errors in the other AR models

Error Model Errors

 time

if your other AR model
has patterns in the errors
ARCH may be useful

For ARCH(1) :

$$\text{Var}(\varepsilon_t) = \sigma_t^2 = \alpha_0 + \alpha_1 \sigma_{t-1}^2$$

so your variance
relates to the
variance
of previous
day

$$\varepsilon_t = w_t \sqrt{\alpha_0 + \alpha_1 \varepsilon_{t-1}^2}$$

So, essentially this just models your error.

Ljung-Box Statistic

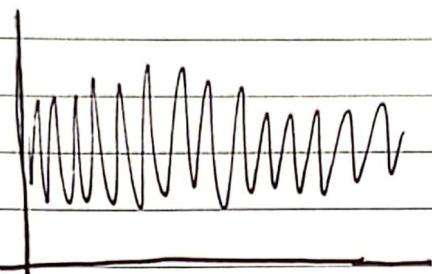
H₀: the time series is random or white noise

H₁: there is serial correlation

$$Q(m) = n(n+2) \sum_{j=1}^m \frac{r_j^2}{n-j}$$

→ m's are the different lags

Power Spectral Density



→ for a stationary process

mean squared error power

$$S_x(f) = \lim_{T \rightarrow \infty} \frac{E[|F_{x_T}(f)|^2]}{2T}$$

$$\bar{x}^2 = \int_{-\infty}^{\infty} S_x(f) df$$

$$X_T(t) = \begin{cases} X(t) & |t| \leq T \\ 0 & |t| > T \end{cases}$$

$$F_{x_T}(f) = \int_{-\infty}^{\infty} X_T(t) e^{-j2\pi ft} dt \quad T < \infty$$

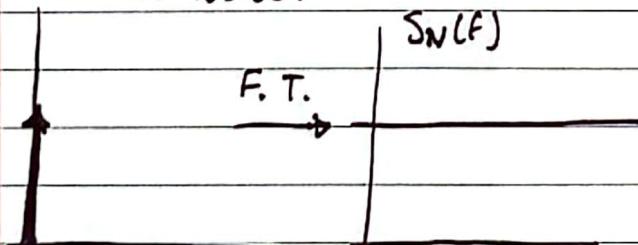
→ take the fourier transform of

$$\int_{-\infty}^{\infty} x^2(t) dt = \int_{-\infty}^{\infty} |F_x(f)|^2 df$$

$$S_x(f) = \text{F.T.} [R_x(\tau)]$$

Fourier transform auto correlation function

ACF of Noise:



ACF

$S_x(f)$



Periodogram

Used to identify dominant period of a time series, which have cyclical behavior.
This falls under spectral analysis. Even data.

Consider data w/ n points, represent the time series w/ sine & cosine

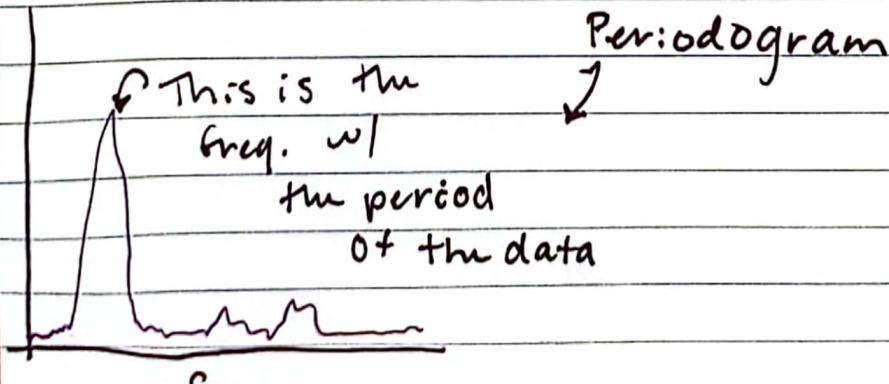
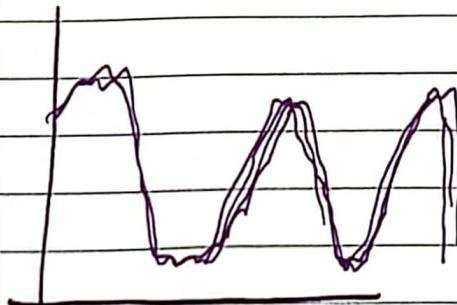
$$x_t = \sum_{j=1}^{n/2} [\beta_1 \left(\frac{j}{n}\right) \cos(2\pi w_j t) + \beta_2 \left(\frac{j}{n}\right) \sin(2\pi w_j t)]$$

estimate β_1 & β_2

$$P\left(\frac{j}{n}\right) = \hat{\beta}_1^2 \left(\frac{j}{n}\right) + \hat{\beta}_2^2 \left(\frac{j}{n}\right)$$

↳ periodogram value at a specific frequency
↳ large P's indicate significant periods

So the data may look like:

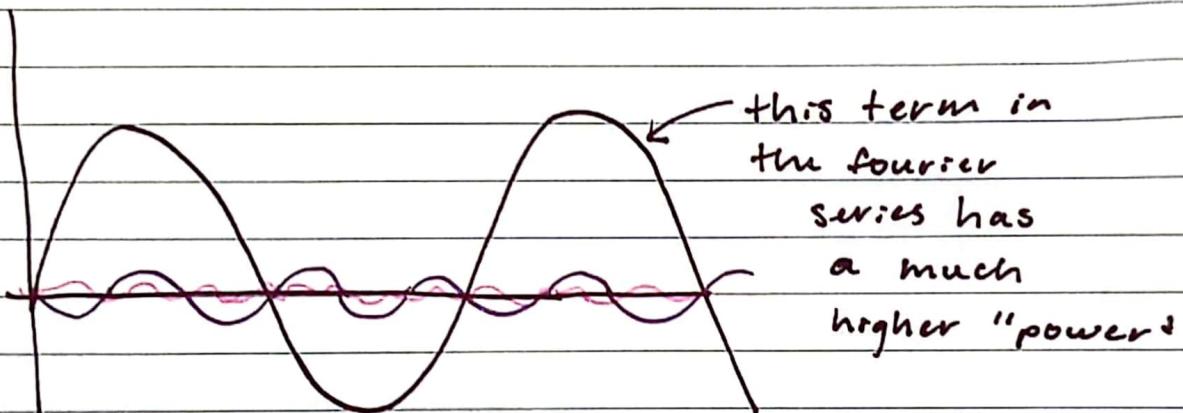


Lomb-Scargle Periodogram

Useful for finding the period of unevenly spaced time series data.

$$f(x) = \frac{1}{2} a_0 + \sum a_n \cos(nx) + \sum b_n \sin(nx)$$

the different "n's" have different powers



So you can keep track of the a_n & b_n for each freq. to see what freq. are "strong".

Now you want to take signals/data points at different times then recover the period

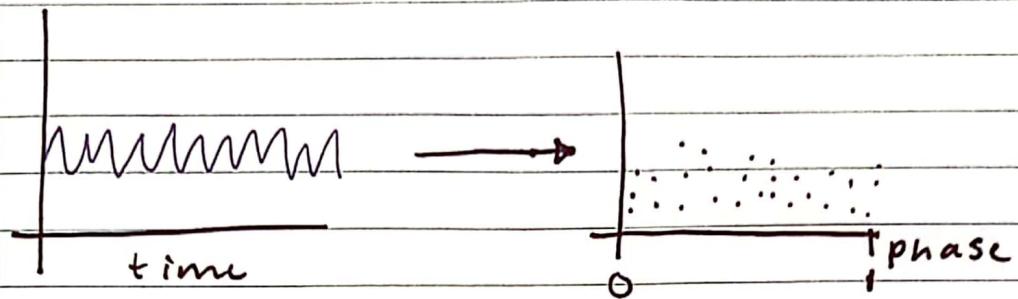
Could try fitting different curves to reconstruct signal. !

LSP is very similar to a normal periodogram but it is modified to deal w/ uneven data.

$$P_{LS}(v) = \frac{1}{2\sigma^2} \left[\frac{\left[\sum_{i=1}^n X_i \cos(2\pi v t_i) \right]^2}{\left[\sum \cos^2(2\pi v(t_i - \tau(v))) \right]} + \frac{\left[\sum X_i \sin(2\pi v t_i) \right]^2}{\left[\sum \sin^2 2\pi v(t_i - \tau(v)) \right]} \right]$$

Minimum String Length

Another way to find the period of the data, by using the folded modulo.



Each of the time values are converted into phase values. You have to cycle through different potential periods. So whatever period you're trying to find the phase.

