

Linear Prediction

Introduction

STAT 464 / 864 | Fall 2024

Discrete Time Series Analysis

Skyepaphora Griffith, Queen's University

Nerd time!

[don't write this down]



Hey Skye, when would we actually want to model an ARMA?

ARMA(p,q) are just really good at modelling (3.1 of Brocky & Davey):

- 🕒 ARMA ACVFs can effectively approximate a huge class of ACVFs
- 🕒 For any positive integer K , there's an ARMA process $\{X_t\}$ such that

$$\gamma_X(h) = \gamma(h) \text{ for } h = 0, 1, \dots, K.$$

We can leverage cool properties

Given AR(p) and MA(q) parameters, say we build an ARMA(p,q) model.

Describes **(unique) stationary** $\{X_t\}$ if and only if
[no roots on unit circle \rightarrow can model as stationary]

$$\phi(z) = 1 - \phi_1 z - \dots - \phi_p z^p \neq 0 \quad \text{for all } |z| = 1.$$

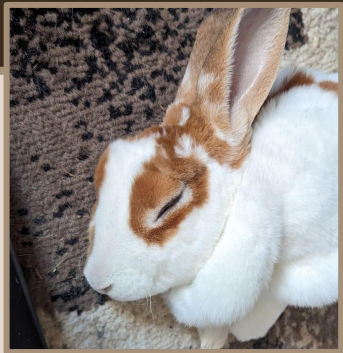
ARMA(p,q) process $\{X_t\}$ is **causal** if and only if
[no roots in unit disk \rightarrow can model using only past values]

$$\phi(z) = 1 - \phi_1 z - \dots - \phi_p z^p \neq 0 \quad \text{for all } |z| \leq 1.$$

- 🕒 Later, we'll see similar tools regarding something called **invertibility**
- 🕒 Also we'll see data-driven examples in workshop eventually

Nerd time!

[don't write this down]



Also uhh what about Yule-Walker?



Not time for that yet (that's chapter 5) but nice pull



We're wrapping up a unit on linear filters → thesis of the lesson:

We can model an ARMA(1,1) process using a linear filter

$$\gamma_X(h) = \sigma^2 \sum_{j=0}^{\infty} a_j a_{j+h} = \begin{cases} \sigma^2 \left(1 + \frac{(\phi+\theta)^2}{1-\phi^2} \right) & h = 0 \\ \sigma^2 (\phi + \theta) \phi^{|h|-1} + \frac{\sigma^2 (\phi+\theta)^2 \phi^{|h|}}{1-\phi^2} & |h| > 0 \end{cases}$$

⌚ Filter method's *derivation* is messy

⌚ Easier to compute than Y.W. for large lags

⌚ In practice, without infinitely many past values, Filter method is an *approximation*
(The past is far behind us, but not that far)

Further reading: Great discussion in the intro to section 5.1 (Brocky & Davey)

Problem Setup | Review from Friday

Let Y be an RV with finite variance. We want to:

- 🕒 Predict Y based on some sequence of RVs, $\mathbf{W} = \{W_N, \dots, W_1\}$
- 🕒 Find a function $g(W_1, \dots, W_N)$ that gives a “good” prediction of Y

We'll measure that “goodness” using **Mean Squared Error (MSE)**

$$MSE \stackrel{\text{def}}{=} E \left[\left(Y - \underbrace{g(W_1, \dots, W_N)}_{\text{prediction}} \right)^2 \right]$$

Linearity of the Predictor | Review from Friday

Optimal Prediction:

Conditional expectation $E[Y | W_1, \dots, W_N]$

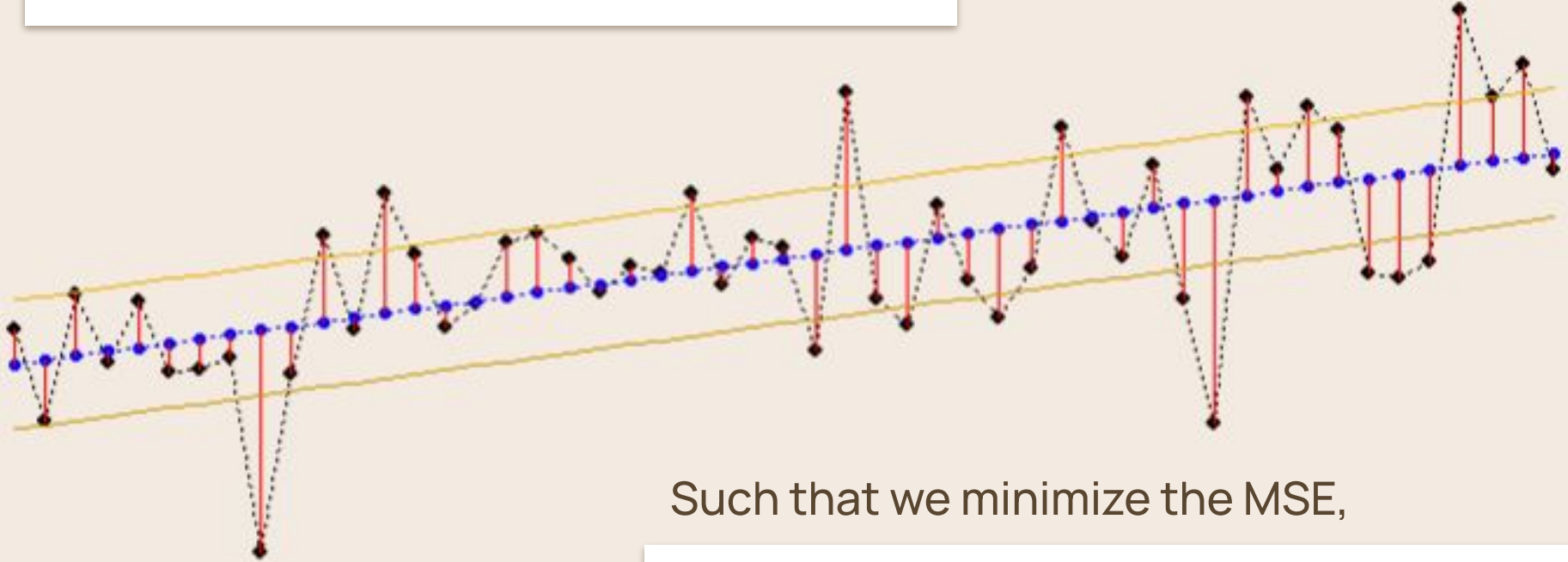
Not computable without joint distribution of Y and W

In time series, we typically only specify up to 2nd order properties (Var/Cov)

If we specify the prediction to be a linear function g of W ,
we can compute optimal prediction + MSE using only 2nd order properties

We wish to find the best linear combination

Remember that
“residual sum of squares”
idea, from linear regression?



Such that we minimize the MSE,

Notation

Note that we reverse the time indices in \mathbf{W} , by convention



In the spirit of filtering



convenient linear algebra, down the road

Minimizing the MSE:

$$E \left[\left(Y - \left(a_0 + \underbrace{a_1 W_N + \cdots + a_N W_1}_{\mathbf{a}^T \mathbf{W}} \right) \right)^2 \right]$$

Assume we can differentiate MSE with respect to each a_i by taking the derivative inside the expectation

Set equal to 0 for all i ,
solve for
coefficients $\{a_i\}_{i=0}^N$

$\mathbf{i} = 0$ Divide both sides by -2 $\rightarrow E[Y - (\mathbf{a}_0 + \mathbf{a}^T \mathbf{W})] = 0 \rightarrow$ Use linearity of E !

Minimizing the MSE ($i > 0$) $0 = E \left[-2(Y - (a_0 + a^T W)) W_{N+1-i} \right]$

Minimizing the MSE

$$0 = \mathbb{E}[(Y - \mu_Y)W_{N+1-i}] - a^T \mathbb{E}[(W - \mu_W)W_{N+1-i}]$$

$$0 = \mathbb{E}[(Y - \mu_Y)(W_{N+1-i} - \mu_{N+1-i})] - a^T \mathbb{E}[(W - \mu_W)(W_{N+1-i} - \mu_{N+1-i})]$$

Why can we subtract μ_{N+1-i} ?

When you take the expectation through,

Then you're left with



This is , the i^{th} component of the vector γ

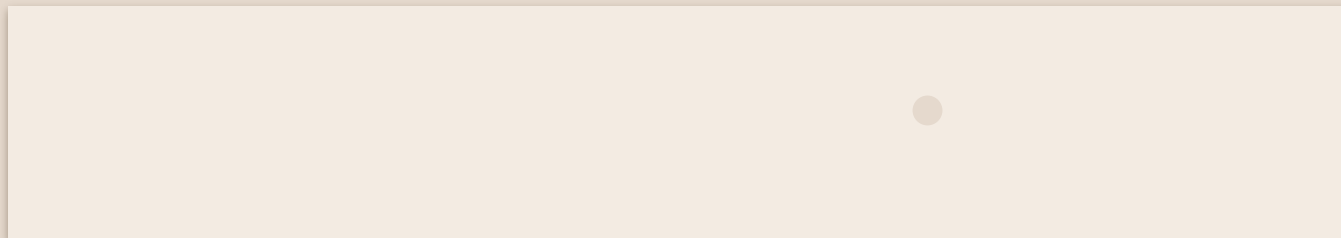


This is the vector of covariances between W_{N+1-i} and components of W
ie) the i^{th} row (or column) of

Minimizing the MSE

$$0 = \text{E}[(Y - \mu_Y)W_{N+1-i}] - a^T \text{E}[(W - \mu_W)W_{N+1-i}]$$

$$0 = \text{E}[(Y - \mu_Y)(W_{N+1-i} - \mu_{N+1-i})] - a^T \text{E}[(W - \mu_W)(W_{N+1-i} - \mu_{N+1-i})]$$



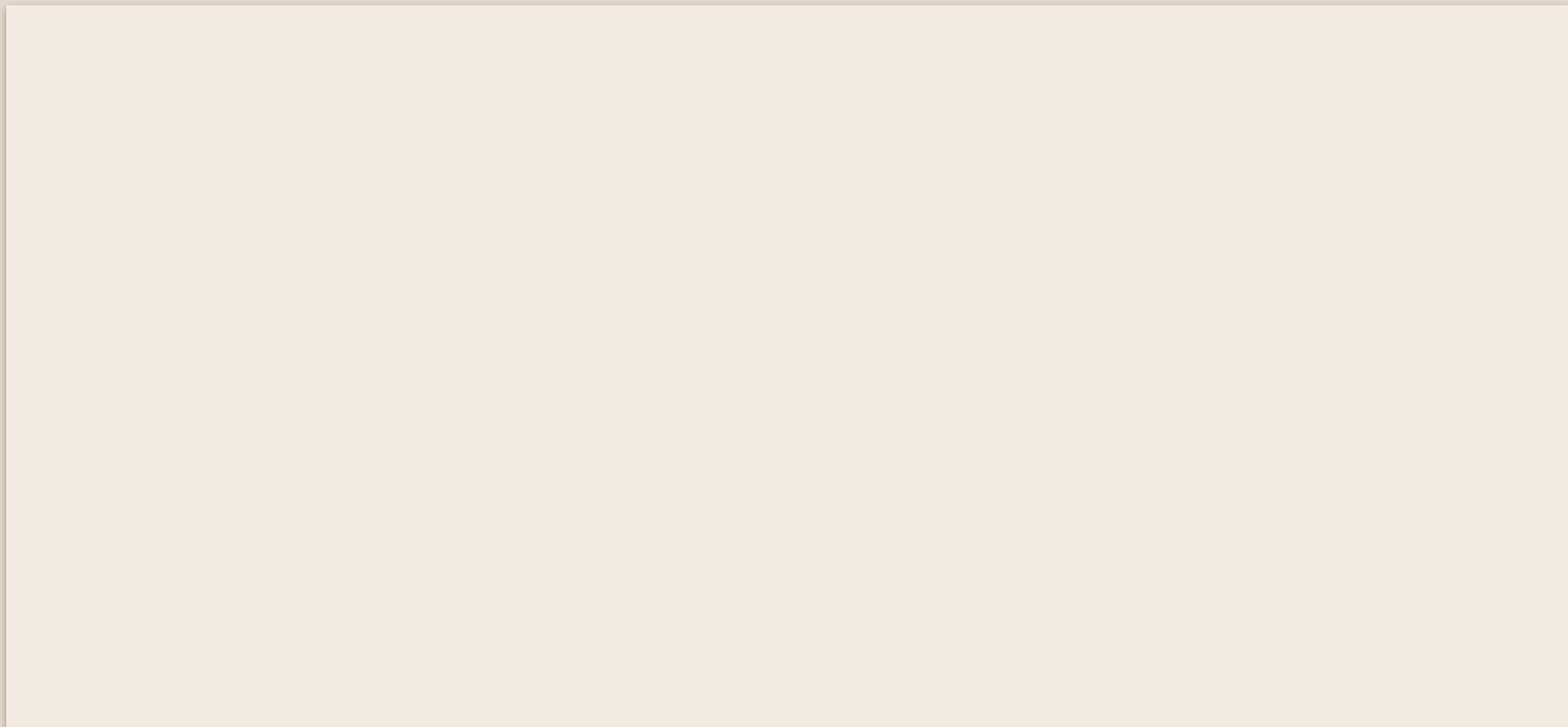
This is $\text{Cov}(Y, W_{N+1-i})$, the i^{th} component of the vector γ



This is the vector of covariances between W_{N+1-i} and components of W
ie) the i^{th} row (or column) of Γ

Minimum MSE

[warning: skye straight up didn't proof-read the LaTeX]



Conclusion

The best linear predictor of Y , based on W , that minimizes MSE is:

$$\hat{Y} = \mu_Y + \gamma(W - \mu_W)$$

Where the vector a satisfies

$$a = \frac{\gamma}{\Gamma}$$

Remark:

This solution is theoretical. Assumes we know μ_Y, μ_W, γ , and Γ
In practice, we must estimate these quantities.