

Introduction and Review

February 7, 2020

The material in this set of notes was initially based on Chapter 2 of Robert Shumway and David Stouffer's Time Series Analysis and Its Applications: With R Examples, but has since grown to include additional material. Notes on the parametric bootstrap are based on Efron and Tibshirani's 1993 textbook An Introduction to the Bootstrap. Chapter 7 of Simon Wood's Core Statistics is also a helpful and relevant reference.

Notation

- $\mathbb{E}[x]$ refers to the expectation of x ;
- $\mathbb{V}[x]$ refers to the variance of x ;
- \mathbb{Z} refers to the set of all positive and negative integers, $\mathbb{Z} = \{0, \pm 1, \pm 2, \dots\}$;
- Bolded lowercase letters denote column vectors;
- Bolded uppercase letters denote matrices;
- $x \sim \mathcal{N}(\mu, \sigma^2)$ indicates that x is normally distributed with mean μ and variance σ^2 ;
 - Another way of writing this is $x = \sigma v + \mu$, where $v \sim \mathcal{N}(0, 1)$. I will sometimes use this notation to describe the distribution of x .

- $x \sim \mathcal{T}_k$ indicates that x is central t -distributed with k degrees of freedom;
- $x \sim \chi_k^2$ indicates that x is chi-square distributed with k degrees of freedom;
- $x \sim \mathcal{F}_{k,j}$ indicates that x is central F -distributed with k and j degrees of freedom.
- $\|\mathbf{x}\|_2^2 = \sum_{i=1}^n x_i^2$;
- $x \approx y$ denotes that x is approximately equal to y .
- $\mathbb{1}_{\{x=1\}}$ and $(x=1)$ are two different ways that I will denote an indicator function, which is equal to 1 if the statement contained in the brackets $\{\cdot\}$ or parentheses (\cdot) is true and equal to 0 otherwise.

Basic Idea!

Most (univariate) time series analysis problems boil down to observing an $n \times 1$ vector $\mathbf{y} = (y_1, \dots, y_n) = \boldsymbol{\mu} + \boldsymbol{\epsilon}$, where $\boldsymbol{\mu}$ is a fixed but unknown mean and $\boldsymbol{\epsilon}$ are mean zero random errors, and:

- **Estimating $\boldsymbol{\mu}$;**
- **Predicting** future values y_{n+1}, \dots, y_{n+k} .

Time series analysis problems differ from classical statistical problems because elements of \mathbf{y} are ordered in time. Several examples of time series data and problems are given in Chapter 1 of S&S, pages 4-11.

Because elements of \mathbf{y} are ordered in time, consecutive elements of \mathbf{y} may be correlated and classical statistical methods may not work well. This is easiest to see via example. Suppose we assume $\mu_{x,t} = \mu$ for all $t = 1, \dots, n$, and we are interested in estimating μ . Ignoring the time series aspect of \mathbf{y} , we assume ϵ_j are independent and identically distributed with known variance σ^2 . The classical approach would be to compute a point estimate of μ , $\hat{\mu} = \sum_{t=1}^n y_t / n$ and corresponding standard error, σ^2 / n . Is this accurate?

The classical approach gives a **incorrect** standard error if elements of \mathbf{y} are correlated. What would be the correct standard error?

$$\begin{aligned}\mathbb{E}[(\hat{\mu} - \mu)^2] &= \mathbb{E}\left[\left(\sum_{t=1}^n y_t - n\mu\right)^2 / n^2\right] \\ &= \sigma^2/n + \sum_{t=1}^n \sum_{t'=1, t' \neq t}^n \mathbb{E}[(y_t - \mu)(y_{t'} - \mu)] / n^2.\end{aligned}$$

The correct standard error depends on covariances of elements of \mathbf{y} ,

$$\mathbb{E}[(y_t - \mu)(y_{t'} - \mu)] = \mathbb{E}[\epsilon_t \epsilon_{t'}],$$

which may be nonzero if elements of \mathbf{y} are ordered in time!

Regression Review (S&S 2.1-2.2)

Many methods for time series analysis build on **linear regression**. We perform linear regression when we are interested in expressing an $n \times 1$ response vector \mathbf{y} as a linear function of r $n \times 1$ covariate vectors $\mathbf{x}_1, \dots, \mathbf{x}_r$, i.e. we want to find regression coefficients β_1, \dots, β_r such that $\mathbf{y} \approx \beta_1 \mathbf{x}_1 + \dots + \beta_r \mathbf{x}_r$. If \mathbf{y} is a time series, then covariates might include:

- Indicators for distinct time periods different elements of \mathbf{y} belong to;
- A vector \mathbf{t} , where t_i is the time y_i was observed or the order of y_i in the sequence;
- Nonlinear functions of elements of \mathbf{t} , e.g. $z_{ij} = \sin(t)$ for some $j \in \{1, \dots, p\}$;
- Lagged values of \mathbf{y} ;
- Lagged values of a different but related time series.

We will very rarely be able to describe \mathbf{y} as an exactly linear function of $\mathbf{x}_1, \dots, \mathbf{x}_r$. Instead, we try to find the “best” way of writing \mathbf{y} as a nearly linear function of $\mathbf{x}_1, \dots, \mathbf{x}_r$

by computing the regression coefficients β that solve:

$$\min_{\beta} \|\mathbf{y} - \beta_1 \mathbf{x}_1 - \cdots - \beta_r \mathbf{x}_r\|_2^2. \quad (1)$$

This is easier to express concisely in matrix form. Letting $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_r]$ be the $n \times r$ matrix of regression coefficients, β equivalently solves:

$$\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2. \quad (2)$$

We refer to the quantity $\|\mathbf{y} - \mathbf{X}\beta\|_2^2$ as the **residual sum of squares (RSS)**, as it measures how much of the variability of \mathbf{y} remains after subtracting off a linear function of the covariates. We can minimize (2) by differentiating; the minimizing value $\hat{\beta}$ will satisfy:

$$\mathbf{X}'\mathbf{X}\hat{\beta} - \mathbf{X}'\mathbf{y} = \mathbf{0} \implies \mathbf{X}'\mathbf{X}\hat{\beta} = \mathbf{X}'\mathbf{y}.$$

If the matrix \mathbf{X} is full rank with rank r , then the minimizing value is

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}. \quad (3)$$

If we want to say more about $\hat{\beta}$, we need to make some more assumptions. First, note that we can always decompose the observed response \mathbf{y} into a linear part $\mathbf{X}\beta$ and a remainder ϵ :

$$\mathbf{y} = \mathbf{X}\beta + \epsilon. \quad (4)$$

If we assume:

- $\mathbb{E}[\epsilon] = \mathbf{0}$, then $\hat{\beta}$ is **unbiased**, i.e. $\mathbb{E}[\hat{\beta}] = \beta$.
- $\epsilon_j \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$, then:

(\star) $\hat{\beta}$ is the maximum likelihood estimator of β ;

(\ast) $\hat{\beta} \sim \text{normal}(\beta, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1})$;

(\dagger) $\mathbf{y} - \mathbf{X}\hat{\beta} \sim \text{normal}(\mathbf{0}, \sigma^2 (\mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'))$;

(\circ) $\hat{\beta}$ and $\mathbf{y} - \mathbf{X}\hat{\beta}$ are independent.

Time series methods rely extensively on likelihood based inference, so we pause to derive (★). If $\epsilon_j \stackrel{i.i.d.}{\sim} \text{normal}(0, \sigma^2)$ then rearranging (4) gives $\mathbf{y}_i - \mathbf{x}_i' \boldsymbol{\beta} \stackrel{i.i.d.}{\sim} \text{normal}(0, \sigma^2)$, where \mathbf{x}_i is the i -th row of \mathbf{X} . This yields the likelihood:

$$l(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \right\}.$$

Finding the value of $\boldsymbol{\beta}$ that maximizes the likelihood is equivalent to finding the value of $\boldsymbol{\beta}$ that minimizes the negative log likelihood, which corresponds to a constant plus the residual sum of squares (2). As a side note - we can actually eliminate σ^2 from the log-likelihood by *profiling* it out, thus making the connection between maximizing the log likelihood and computing the residual sum of squares even clearer. The negative log-likelihood is given by:

$$\frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2.$$

We want to minimize this with respect to $\boldsymbol{\beta}$ and σ^2 , and this can be achieved by *first* minimizing over σ^2 for fixed $\boldsymbol{\beta}$, and plugging in the minimizing value of σ^2 to get a minimization problem that depends on $\boldsymbol{\beta}$ alone. This takes advantage of the observation that:

$$\min_{\sigma^2, \boldsymbol{\beta}} \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 = \max_{\boldsymbol{\beta}} \left(\max_{\sigma^2} \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \right).$$

There's a nice closed form solution to the minimization with respect to σ^2 for fixed $\boldsymbol{\beta}$. Taking derivatives with respect to σ^2 and rearranging, we get:

$$\frac{n}{2\sigma^2} - \frac{1}{2\sigma^4} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 = 0 \implies \sigma^2 = \frac{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2}{n}.$$

Then plugging this expression for the minimizing value of σ^2 for fixed $\boldsymbol{\beta}$, we get a minimiza-

tion problem over just β :

$$\begin{aligned}
\min_{\sigma^2, \beta} \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 &= \min_{\beta} \left(\min_{\sigma^2} \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \right) \\
&= \min_{\beta} \left(\frac{n}{2} \right) \log \left(2\pi \left(\frac{\|\mathbf{y} - \mathbf{X}\beta\|_2^2}{n} \right) \right) + \frac{\|\mathbf{y} - \mathbf{X}\beta\|_2^2}{2 \left(\frac{\|\mathbf{y} - \mathbf{X}\beta\|_2^2}{n} \right)} \\
&= \min_{\beta} \left(\frac{n}{2} \right) \log \left(2\pi \left(\frac{\|\mathbf{y} - \mathbf{X}\beta\|_2^2}{n} \right) \right) + \frac{n}{2}.
\end{aligned}$$

Now only one term depends on the residual sum of squares $\|\mathbf{y} - \mathbf{X}\beta\|_2^2$, and since $\log(\cdot)$ is a strictly increasing function, this is equivalent to minimizing the residual sum of squares. This isn't very useful to us now, because the log-likelihood is already easy to maximize, and because the maximum likelihood estimate of β does not depend on σ^2 . However we will find the idea of profiling very useful when we start considering maximum likelihood estimation for time series models, in which elements of the error vector ϵ are correlated and we have more parameters to optimize over.

Returning to the claims we can make about the least squares estimate $\hat{\beta}$ when we assume that $\epsilon_j \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$, $(*)$, (\dagger) , and (\circ) are very useful. They allow us not only to compute standard errors and confidence intervals for $\hat{\beta}$ but also to test the null hypothesis that β_i is exactly equal to a specific value or that a subset of $p - p_1$ elements $\beta_{-1} = (\beta_{t_1}, \dots, \beta_{t_{p-p_1}})$ are jointly exactly equal to $\mathbf{0}$.

Standard practice for constructing standard errors and confidence intervals is to use $(*)$, plugging in an unbiased estimator of the variance:

$$s^2 = \frac{\left\| \mathbf{y} - \mathbf{X}\hat{\beta} \right\|_2^2}{n - p}. \quad (5)$$

Note that this is *not* the maximum likelihood estimate of σ^2 - the maximum likelihood estimator $\hat{\sigma}^2 = \left\| \mathbf{y} - \mathbf{X}\hat{\beta} \right\|_2^2 / n$ is biased.

It follows from $(*)$, (\dagger) , and (\circ) that

$$\frac{\hat{\beta}_i - \beta_i}{s\sqrt{(\mathbf{X}'\mathbf{X})_{ii}^{-1}}} \sim \mathcal{T}_{n-p}. \quad (6)$$

This gives us a way of testing the null hypothesis that β_i is exactly equal to a specific value because it tells us the approximate distribution of $\hat{\beta}_i$ for specific values of β_i . We call such tests **t-tests** or **z-tests**.

Similarly, letting $\mathbf{X}_1 = [\mathbf{x}_{t_1}, \dots, \mathbf{x}_{t_{p_1}}]$ be the design matrix containing the p_1 columns corresponding to elements of $\boldsymbol{\beta}_1$ and letting $\hat{\boldsymbol{\beta}}_1$ be the linear regression estimate of $\boldsymbol{\beta}_1$ from regressing \mathbf{y} on just the p_1 columns of \mathbf{X} contained in \mathbf{X}_1 , it follows from $(*)$, (\dagger) , and (\circ) that

$$F_{p-p_1, n-p} = \left(\frac{\|\mathbf{y} - \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1\|_2^2 - \|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}\|_2^2}{\|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}\|_2^2} \right) \left(\frac{n-p}{p-p_1} \right) \sim \mathcal{F}_{p-p_1, n-p}. \quad (7)$$

This gives us a way of testing the null hypothesis that the $p - p_1$ elements of $\boldsymbol{\beta}_{-1}$ are jointly exactly equal to $\mathbf{0}$ by giving us an approximate distribution of $F_{p-p_1, n-p}$ under the null. We call such tests **F-tests**.

Parametric Bootstrap for Standard Errors, Confidence Intervals, and Testing

The derivation of standard errors, confidence intervals, and testing given above are all based on the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \text{normal}(\mathbf{0}, \sigma^2 \mathbf{I}_n). \quad (8)$$

In all cases, we are computing quantities that tell us about the expected variability of a function of the data \mathbf{y} and \mathbf{X} across repeated realizations of data \mathbf{y} , assuming a model of the form (8) holds. In general, we'll denote the quantity we are interested in understanding the expected variability of as $c = g(\mathbf{y}, \mathbf{X})$.

For example, when we are interested in computing the standard error of $\hat{\beta}_j$, we would be interested in understanding the variability of $c = g(\mathbf{y}, \mathbf{X}) = \hat{\beta}_j$. As another example,

if we were interested in testing the null hypothesis that $\beta_j = 0$, we would be interested in understanding the variability of $c = g(\mathbf{y}, \mathbf{X}) = \frac{\hat{\beta}_j}{s\sqrt{(\mathbf{X}'\mathbf{X})_{jj}^{-1}}}$.

A natural approach is to simulate synthetic realizations of the observed data from the approximate model,

$$\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \text{normal}(\mathbf{0}, s^2\mathbf{I}_n), \quad (9)$$

plugging in our estimates $\hat{\boldsymbol{\beta}}$ and s^2 for the unknown parameters $\boldsymbol{\beta}$ and σ^2 . This approach is commonly referred to as the **Parametric Bootstrap**, and it is especially useful when we are interested in understanding the expected variability of functions of the unknown parameters for which the distribution under (8) is either only known approximately, e.g. for large n relative to p , or is not known at all. Intuitively, the parametric bootstrap will work well as long as $\hat{\boldsymbol{\beta}}$ and s^2 are “good enough” estimates of $\boldsymbol{\beta}$ and σ^2 and as long as the distribution of the quantity we are interested c in does not vary “too much” as a function of $\boldsymbol{\beta}$ and σ^2 .

In practice, the parametric bootstrap proceeds as follows:

- Estimate $\boldsymbol{\beta}$ and σ^2 from the model given by (8), and compute $c = g(\mathbf{y}, \mathbf{X})$.
- Set a desired number of bootstrap samples, n_{boot} .
- For k in $1, \dots, n_{boot}$:
 - Draw $\mathbf{y}^{(k)} \sim \text{normal}(\mathbf{X}\hat{\boldsymbol{\beta}}, s^2\mathbf{I}_n)$;
 - Compute $c^{(k)} = g(\mathbf{y}^{(k)}, \mathbf{X})$.

This procedure produces n_{boot} simulated values of the quantity of interest $c^{(1)}, \dots, c^{(n_{boot})}$ from the approximate model given by (9), which can be used for many purposes!

- An estimate of the standard error of c can be obtained by computing the sample standard deviation of the simulated values $c^{(1)}, \dots, c^{(n_{boot})}$.
- An approximate level- α test of the null hypothesis that c is equal to zero under the

model (9) can be obtained by comparing c to the $\alpha/2$ and $1 - \alpha/2$ quantiles of the simulated values $c^{(1)}, \dots, c^{(n_{boot})}$.

- This is not the only or best way to construct a $1 - \alpha$ interval for a quantity c based on bootstrap samples $c^{(1)}, \dots, c^{(n_{boot})}$, but it is very intuitive so we will use it here. You can find more information about alternative methods in Bootstrap Methods and their Applications (1997) by Davison and Hinkley, or other more recent textbooks or papers on the construction of bootstrap-based confidence intervals.

Model Selection Tests are very useful for **model selection**, i.e. for choosing the covariates to include in our model. Model selection is especially relevant in linear regression methods for time series analysis, e.g. we may need to decide which lagged values of \mathbf{y} to include as covariates. Letting \mathbf{X}_k refer to a matrix containing k covariates and $\boldsymbol{\beta}_k$ and $\hat{\boldsymbol{\beta}}_k$ the corresponding regression coefficients and their linear regression estimates, several popular methods for performing model selection when performing linear regression are:

- (*) Perform an t - or z -test comparing *nested* models with k and $k' = k + 1$ covariates.
- (*) Perform an F -test comparing *nested* models with k and k' covariates.
- (*) Compute **Akaike's Information Criterion (AIC)**

$$AIC = \ln \left(\frac{\left\| \mathbf{y} - \mathbf{X}_k \hat{\boldsymbol{\beta}}_k \right\|_2^2}{n} \right) + \frac{n + 2k}{n} \quad (10)$$

for models with k and k' covariates, and choose the model with the lower AIC value.

- (*) Compute **AIC, Bias Corrected (AICc)**

$$AICc = \ln \left(\frac{\left\| \mathbf{y} - \mathbf{X}_k \hat{\boldsymbol{\beta}}_k \right\|_2^2}{n} \right) + \frac{n + k}{n - k - 2} \quad (11)$$

for models with k and k' covariates, and choose the model with the lower $AICc$ value.

(★) Compute **Schwarz's/Bayesian Information Criterion (SIC/BIC)**

$$SIC = \ln \left(\frac{\left\| \mathbf{y} - \mathbf{X}_k \hat{\boldsymbol{\beta}}_k \right\|_2^2}{n} \right) + \frac{k \log(n)}{n} \quad (12)$$

for models with k and k' covariates, and choose the model with the lower SIC value.

(★) Choose any measure of model performance, denoted as $f(\mathbf{y}, \mathbf{X}, \boldsymbol{\beta})$, define several equally sized, possibly overlapping subsets of the data (training datasets), for each subset fit the models and evaluate the performance of each model on the remaining data (test data), and choose the model that performs best on average across all of the training datasets. This is commonly referred to as **cross-validation**. There are many different ways of performing cross-validation:

- Leave- k -out cross-validation refers to a type of cross-validation that constructs l training data sets by choosing k elements of the full data, and constructing the training data to be the full data with the k chosen elements removed, and the test data to be the chosen k elements. In this case, the average measure of performance is an average across l training data sets.

- * When $k = 1$, we will often set $l = n$ because there are n ways to remove a single observation from the full data. When $k > 1$, we will often set $\binom{l}{n,k}$ and choose the l subsets of k elements each at random, because it can become computationally prohibitive to average across all $\binom{n}{k}$ possible ways to remove k observations from the full data.

- k -fold cross-validation refers to a type of cross-validation that forms k training datasets by dividing the data into k subsets of size n/k , and then defining each test dataset to be one of the subsets of n/k observations in the full data, and each training dataset as the remaining observations. In this case, the average measure

of performance is an average across k training data sets.

- * This corresponds to leave- $\binom{n}{n/k}$ -out cross-validation with $l = k$ total subsets/training datasets.

- A common measure of performance is squared error loss, $f(\mathbf{y}, \mathbf{X}, \boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$.

Note that the methods marked with (*) require that the two models be *nested*, i.e. the columns in \mathbf{X}_k must be a subset of the columns in $\mathbf{X}_{k'}$ or vice versa. The procedures denoted with (★) are not. Whether AIC, AICc, or BIC is most appropriate for a given problem is problem-specific; AICc can perform better than AIC when n is relatively small, and SIC/BIC can perform better than AIC when the number of covariates k is relatively large. Cross-validation also does not require that the two models be nested, and can be performed in many different ways depending on the choice of k and the choice of a measure of model performance. A very popular choice is 10-fold cross-validation using squared error loss to measure performance on the test data, $f(\mathbf{y}, \mathbf{X}, \boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$.