

Multivariate Regression

April 21, 2020

Multivariate Regression Review (S&S 5.7)

Many methods for multivariate time series analysis build on **multivariate linear regression**, also known as **general linear regression** (not to be confused with generalized linear regression!). When we perform multivariate linear regression, we jointly model r $n \times 1$ response vectors $\mathbf{y}_1, \dots, \mathbf{y}_r$ arranged as an $n \times r$ matrix $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_r]$ as a linear function of the same $n \times 1$ covariate vectors $\mathbf{x}_1, \dots, \mathbf{x}_q$ arranged as an $n \times q$ matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_q]$. We want to find an $q \times r$ matrix of regression coefficients \mathbf{B} such that $\mathbf{Y} \approx \mathbf{X}\mathbf{B}$ by solving:

$$\min_{\mathbf{B}} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_2^2, \quad (1)$$

where $\|\mathbf{Y}\|_2^2 = \sum_{i=1}^n \sum_{j=1}^r y_{ij}^2$ gives the sum of squared elements of the matrix \mathbf{Y} .

We still refer to the quantity $\|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_2^2$ 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 also still minimize (1) by differentiating; the minimizing value $\hat{\mathbf{B}}$ will satisfy:

$$\mathbf{X}'\mathbf{X}\hat{\mathbf{B}} - \mathbf{X}'\mathbf{Y} = \mathbf{0} \implies \mathbf{X}'\mathbf{X}\hat{\mathbf{B}} = \mathbf{X}'\mathbf{Y}.$$

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

$$\hat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}. \quad (2)$$

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

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{W}. \quad (3)$$

If we assume:

- $\mathbb{E}[\mathbf{W}] = \mathbf{0}$, then $\hat{\mathbf{B}}$ is **unbiased**, i.e. $\mathbb{E}[\hat{\mathbf{B}}] = \mathbf{B}$.
- $\mathbf{w}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \Sigma_w)$, where \mathbf{w}_i are columns of the remainder \mathbf{W} , then:
 - (\star) $\hat{\mathbf{B}}$ is the maximum likelihood estimator of \mathbf{B} ;
 - (\ast) Elements of $\hat{\mathbf{B}}$ are normally distributed, with $\mathbb{V}[\hat{\mathbf{b}}_i] = \sigma_{ii}(\mathbf{X}'\mathbf{X})^{-1}$ and $\text{Cov}[\hat{\mathbf{b}}_i, \hat{\mathbf{b}}_j] = \sigma_{ij}(\mathbf{X}'\mathbf{X})^{-1}$ where $\hat{\mathbf{b}}_i$ be the i -th column of $\hat{\mathbf{B}}$;
 - (\dagger) The residuals $\mathbf{R} = \mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}$ are normally distributed, with $\mathbb{E}[\mathbf{R}] = \mathbf{0}$, $\mathbb{V}[\mathbf{r}_i] = \sigma_{ii}(\mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')$, and $\text{Cov}[\mathbf{r}_i, \mathbf{r}_j] = \sigma_{ij}(\mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')$ where \mathbf{r}_i be the i -th column of \mathbf{R} ;
 - (\circ) $\hat{\mathbf{B}}$ and $\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}$ are independent.

We're not going to derive (\star) this time around. Standard practice for constructing standard errors and confidence intervals is to use (\ast), plugging in an unbiased estimator of the variance-covariance matrix Σ_w :

$$\mathbf{S}_w = \frac{\mathbf{R}'\mathbf{R}}{n - q}. \quad (4)$$

Note that this is *not* the maximum likelihood estimate of Σ_w - the maximum likelihood estimator $\hat{\sigma}_w^2 = \mathbf{R}'\mathbf{R}/n$ is biased.

It follows from (\ast), (\dagger), and (\circ) that

$$t_{n-q} = \frac{\hat{b}_{ij} - b_{ij}}{\sqrt{S_{w,jj}} \sqrt{(\mathbf{X}'\mathbf{X})_{ii}^{-1}}} \sim \mathcal{T}_{n-q}. \quad (5)$$

This gives us a way of testing the null hypothesis that b_{ij} is exactly equal to a specific value because it tells us the approximate distribution of \hat{b}_{ij} for specific values of b_{ij} . We call such tests **t-tests**.

F-tests are a bit trickier to derive for multivariate linear models, so we'll just talk about performing model selection (choosing the covariates or columns of \mathbf{X} to use) using AIC, AICc and SIC. Letting \mathbf{X}_k refer to a matrix containing k covariates and \mathbf{B}_k and $\hat{\mathbf{B}}_k$ the corresponding regression coefficients and their linear regression estimates, several popular methods for performing model selection are:

(★) Compute **Akaike's Information Criterion (AIC)**

$$AIC = \ln \left(\left| \frac{(\mathbf{Y} - \mathbf{X}_k \hat{\mathbf{B}}_k)' (\mathbf{Y} - \mathbf{X}_k \hat{\mathbf{B}}_k)}{n} \right| \right) + \frac{2}{n} \left(rk + \frac{r(r+1)}{2} \right) \quad (6)$$

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

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

$$AICc = \ln \left(\left| \frac{(\mathbf{Y} - \mathbf{X}_k \hat{\mathbf{B}}_k)' (\mathbf{Y} - \mathbf{X}_k \hat{\mathbf{B}}_k)}{n} \right| \right) + \frac{r(n+q)}{n-r-q-1} \quad (7)$$

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

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

$$SIC = \ln \left(\left| \frac{(\mathbf{Y} - \mathbf{X}_k \hat{\mathbf{B}}_k)' (\mathbf{Y} - \mathbf{X}_k \hat{\mathbf{B}}_k)}{n} \right| \right) + (kr + r(r+1)/2) \left(\frac{\log(n)}{n} \right) \quad (8)$$

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

Recall that 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. Because

including one additional covariate (column of \mathbf{X}) yields r additional regression coefficients when we are performing multivariate linear regression, we may tend to prefer SIC/BIC.