Time Series Analysis Basics of Time Series Analysis

Nicoleta Serban, Ph.D.

Professor

Stewart School of Industrial and Systems Engineering

Basic Statistical Modeling Concepts



About This Lesson





Multiple Linear Regression: Model

```
Data: \{(x_{1,1}, ..., x_{1,p}), y_1\}, ..., \{(x_{n,1}, ..., x_{n,p}), y_n\}

Model: Y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \cdots + \beta_p x_{i,p} + \varepsilon_i, i = 1, ..., n
```

Assumptions:

- Linearity/Mean Zero Assumption: $E(\varepsilon_i) = 0$
- Constant Variance Assumption: $Var(\varepsilon_i) = \sigma^2$
- Independence Assumption: $\{\varepsilon_1,...,\varepsilon_n\}$ are independent random variables
- ε_i ~ Normally distributed for confidence/prediction intervals, hypothesis testing



Multiple Linear Regression: Model

Data:
$$\{(x_{1,1},...,x_{1,p}), y_1\},...,\{(x_{n,1},...,x_{n,p}), y_n\}$$

Model: $Y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \cdots + \beta_p x_{i,p} + \varepsilon_i, i = 1,...,n$

Assumptions:

- Linearity/Mean Zero Assumption: $E(\varepsilon_i) = 0$
- Constant Variance Assumption: $Var(\varepsilon_i) = \sigma^2$
- Independence Assumption: $\{\varepsilon_1,...,\varepsilon_n\}$ are independent random variables
- ε_i ~ Normally distributed for confidence/prediction intervals, hypothesis testing

The model parameters are: β_0 , β_1 , ..., β_p , σ^2

- Unknown regardless how much data are observed
- Estimated given the model assumptions
- · Estimated based on data



Parameter Estimation $(\beta_0, \beta_1, ..., \beta_p), \sigma^2$

To estimate $(\beta_0, \beta_1, ..., \beta_p)$, find values $(\hat{\beta}_0, \hat{\beta}_1, ..., \hat{\beta}_p)$ that minimize squared error:

$$\sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{n} (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i,1} + \dots + \hat{\beta}_p x_{i,p}))^2 = (y - X \hat{\beta})^T (y - X \hat{\beta})$$

By linear algebra (Orthogonal Decomposition Theorem) or differentiation:

$$X^{\mathrm{T}}(y - \widehat{y}) = X^{\mathrm{T}}(y - X\widehat{\beta}) = 0$$

So

$$X^{\mathrm{T}}X\widehat{\boldsymbol{\beta}} = X^{\mathrm{T}}y$$

If X^TX is invertible,

$$\widehat{\boldsymbol{\beta}} = (\boldsymbol{X}^{\mathrm{T}}\boldsymbol{X})^{-1}\boldsymbol{X}^{\mathrm{T}}\boldsymbol{y}$$



Parameter Estimation (cont'd)

The fitted values are $\hat{y} = X\hat{\beta}$, and $\hat{\beta} = (X^TX)^{-1}X^Ty$, so

$$\widehat{y} = X\widehat{\beta} = X(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y = \mathrm{H}y$$

where $\mathbf{H} \equiv \mathbf{X}(\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathrm{T}}$ is called the **hat matrix** because multiplying \mathbf{y} by \mathbf{H} gives $\hat{\mathbf{y}}$.

The residuals are:

$$\hat{\varepsilon} = y - \hat{y} = y - X\hat{\beta} = y - Hy = (I - H)y$$

To estimate
$$\sigma^2$$
,
$$\widehat{\sigma}^2 = \widehat{\boldsymbol{\varepsilon}}^T \widehat{\boldsymbol{\varepsilon}} / (n - p - 1)$$



Model Interpretation

The Least Squares estimated coefficients have specific interpretations:

- $\hat{\beta}_0$ The estimated expected value of the response variable when all predicting variables equal zero.
- The estimated expected change in the value of the response variable associated with one unit of change in the value of the i^{th} predicting variable (i.e., associated with a one-unit change in x_i , where i is any of 1, ..., p), holding all other predictors in the model fixed (i.e., holding fixed x_i for j = 1, ..., p where $j \neq i$).



Sampling Distribution

$$E(\hat{\beta}) = \beta$$

$$V(\hat{\beta}) = \sigma^{2} (X^{T}X)^{-1} = \Sigma$$

Furthermore, $\hat{\beta}$ is a linear combination of $\{y_1,...,y_n\}$. If we assume that $\varepsilon_i \sim \text{Normal}$ $(0, \sigma^2)$, then $\hat{\beta}$ is also distributed as $\hat{\beta} \sim \text{Normal}$ $(0, \sigma^2)$.

$$\widehat{\beta} \sim N(\beta, \Sigma)$$



Sampling Distribution (cont'd)

$$\widehat{\beta} \sim N(\beta, \Sigma)$$

 σ^2 is unknown! Replace σ^2 with $\hat{\sigma}^2 = MSE$

$$\widehat{\sigma}^2 = \frac{\sum \widehat{\sigma}_i^2}{n - p - 1} \sim \chi_{n - p - 1}^2$$

$$(\text{chi-squared distribution with } n - p - 1 \text{ degrees of freedom})$$

$$(t - distribution with n - p - 1 \text{ degrees of freedom})$$



Confidence Interval Estimation

We can derive confidence intervals for β_i using this t sampling distribution:

$$\hat{\beta}_j \pm (t_{\alpha/2, n-p-1})(SE(\hat{\beta}_j))$$

Is β_i statistically significant?

Check whether zero is in the confidence interval

Why is this a *t*-interval?



Testing Statistical Significance

To test for statistical significance of β_j given all other predicting variables in the model, use a *t*-test for H_0 and H_a :

$$H_0$$
: $\beta_j = 0$ vs. H_a : $\beta_j \neq 0$

$$t-\text{value} = \frac{\hat{\beta}_j - 0}{\text{SE}(\hat{\beta}_j)}$$

- Reject H₀ if |t-value| gets too large
- Interpret rejecting the null hypothesis as β_i being statistically significant



General Hypothesis Testing

How will the procedure change if

$$H_0$$
: $\beta_j = b$ vs. H_a : $\beta_j \neq b$

for some known null value b?

we test
$$t$$
-value $=\frac{\hat{\beta}_j - b}{\text{SE}(\hat{\beta}_j)}$

- Reject H₀ if | t-value | is large
 - For significance level α , if $|t-value| > t_{\alpha/2, n-p-1} \longrightarrow \text{reject } H_0$
- Alternatively, compute a p-value based on the probability that the t distribution is greater than the *t*-value:

p-value =
$$2 \text{Prob}(T_{n-p-1} > |t\text{-value}|)$$

If p-value is small (e.g., < 0.01) \longrightarrow reject H₀



Testing Subsets of Coefficients

Partial F-test:

• Consider a full model with two sets of predictors, X_1 , ..., X_p (perhaps controlling factors) and $(Z_1, ..., Z_q)$ (perhaps additional explanatory factors):

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \alpha_1 Z_1 + \dots + \alpha_q Z_q + \varepsilon$$

Test whether any of the Z factors add explanatory power to the model:

$$\mathbf{H_0}$$
: $\alpha_1 = \alpha_2 = \cdots = \alpha_q = 0$ vs. $\mathbf{H_a}$: $\alpha_i \neq 0$ for at least one α_i , $i = 1, \dots, q$

F-statistic =
$$F_{partial} = \frac{SSReg(Z_1, ..., Z_q | X_1, ..., X_p)/q}{SSE(Z_1, ..., Z_q, X_1, ..., X_p)/(n-p-q-1)}$$

- Reject H_0 if F-statistic is large (F-statistic > $F_{\alpha, q, n-p-q-1}$)
 - At least one coefficient is different from zero at the α significance level



Normality Transformation

Problem: Constant variance or/and normality assumption

Solution: Transform the response variable from y to \mathring{y} via

$$\dot{y} = y^{\lambda}$$

where the value of λ depends on how Var(y) changes as x changes.

$$\sigma_{y}(x) \propto const \qquad \lambda = 1 \qquad \text{(don't transform)}$$

$$\sigma_{y}(x) \propto \sqrt{\mu_{x}} \qquad \lambda = 1/2 \qquad \mathring{\boldsymbol{y}} = \sqrt{\boldsymbol{y}}$$

$$\sigma_{y}(x) \propto \mu_{x} \qquad \lambda = 0 \qquad \mathring{\boldsymbol{y}} = \ln(\boldsymbol{y})$$

$$\sigma_{y}(x) \propto \mu_{x} \qquad \lambda = -1 \qquad \mathring{\boldsymbol{y}} = \frac{1}{y}$$



Outliers in Regression

A data point far from the majority of the data (in *y* and/or any *x*) may be called an *outlier*, especially if it does not follow the general trend of the rest of the data.

- Data points that are far from the means of the Xs or near the edge of the observation space are called leverage points.
- A data point that is far from the means of y and/or an x is called an influential point if it influences the fit of the regression.
- Excluding a leverage point may or may not the regression fit significantly, thus
 a leverage point may or may not be an influential point.

The upshot: Sometimes there are good reasons to exclude subsets of data (e.g., errors in data entry or experimental errors). Sometimes an outlier belongs in the data. Outliers should always be examined.



R²: Coefficient of Determination

A measure that efficiently summarizes how well the Xs can be used to predict Y is R² (called *R-squared* or the *coefficient* of *determination*):

$$R^2 = 1 - SSE/SST$$

where

SSE =
$$\sum_{i=1}^{n} \hat{\varepsilon}_{i}^{2} = \sum_{i=1}^{n} (y_{i} - \hat{y}_{i})^{2}$$

SST = $\sum_{i=1}^{n} (y_{i} - \bar{y})^{2}$

R² is interpreted as the proportion of total variability in Y that can be explained by the linear regression model.



Multicollinearity

Recall that finding the ordinary least squares estimator of $\widehat{m{\beta}}$

$$\widehat{\boldsymbol{\beta}} = \left(\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X} \right)^{-1} \boldsymbol{X}^{\mathrm{T}} \boldsymbol{Y}$$

depends on X^TX being invertible (nonsingular or nondegenerate). From linear algebra, a square matrix is invertible if and only if its columns are linearly independent (i.e., no column is a linear combination of the others).

If that doesn't hold, the ordinary least squares estimator of $\hat{\beta}$ doesn't exist. That's probably due to a specification error where one or more predictors should be eliminated as redundant (e.g., if years and number of rings were included in a model for trees).

Even if the columns of X^TX are linearly independent, some problems might arise if the value of one predictor can be closely estimated from the other predictors. We call this condition *multicollinearity* or *near collinearity*.

Multicollinearity (cont'd)

- Indications that near collinearity is present:
 - The estimated coefficients $\hat{\beta}$ are unstable: When the value of one predictor changes slightly, the fitted regression coefficients change dramatically
 - The standard error of $\hat{\beta}$ is artificially large
 - The overall F statistic is significant, but individual *t*-statistics are not
- Prediction may be affected
 - The relationship to the response may change widely
- Some computational algorithms are sensitive to multicollinearity
- But no inflation or deflation in R²



Summary



