

Central Limit Theorem for Linear Models

Fall 2014

The CLT you learn in Math Stat applies to sums or averages of independent R.V.s. We need a slightly different version to apply to coefficient estimates, which are much like means.

From Math Stat:

Lindeberg CLT (Arnold, 1981) says that if $Y_1, Y_2, \dots \sim (0, \sigma^2)$ with $\sigma^2 < \infty$, and c_1, c_2, \dots is a sequence of constants such that $\frac{\max c_j^2}{\sum c_j^2} \rightarrow 0$ as $n \rightarrow \infty$, then $\frac{\sum_{j=1}^n c_j Y_j}{\sqrt{\sum_{j=1}^n c_j^2}}$ converges in distribution to $N(0, \sigma^2)$. The condition keeps the c_j 's from increasing so fast that the last term dominates the sum.

Simple case:

Suppose we have a balanced ANOVA model (1, 2, or multiway) with n observations per cell, $Y_{ij} - \mu_i \sim \text{iid}(0, \sigma^2)$. Take the mean in group i to be $\frac{1}{n} \sum (Y_{ij} - \mu_i)$, so let $c_j = n^{-1}$. The condition is easily satisfied because $\max(c_j^2) = n^{-2}$ and $\sum c_j^2 = n \times n^{-2} = n^{-1}$ so the fraction is n^{-1} which goes to 0 nicely as $n \rightarrow \infty$. We conclude

$$\frac{n^{-1} \sum (Y_{ij} - \mu_i)}{n^{-1/2}} = n^{1/2}(\bar{Y}_i - \mu_i) \xrightarrow{\mathcal{D}} N(0, \sigma^2) \text{ or } n^{1/2}\bar{Y}_i \xrightarrow{\mathcal{D}} N(\mu_i, \sigma^2).$$

For any n , the cell means are independent of each other. Any linear combination of cell means will converge in distribution to a normal distribution.

Full linear model (Sen and Srivastava, 1997)

Change the condition for the CLT slightly, letting $a_i = c_i / \sqrt{\sum c_i^2}$ and require $\max |a_{n_i}| \rightarrow 0$ and $\sum_i a_{n_i}^2 \rightarrow 1$. This allows the use of a triangular array of constants instead of reusing the same values (Gnedenko and Kolmogorov, 1954). The conclusion is that $\sum a_{n_i} Y_i \xrightarrow{\mathcal{D}} N(0, \sigma^2)$ as $n \rightarrow \infty$.
Linear model:

Take the simple case: $\mathbf{y} \sim (\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$ with \mathbf{X} of full column rank, and consider the coefficient vector $(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$, or actually the distribution of

$$\begin{aligned} \sigma^{-2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top (\mathbf{X}^\top \mathbf{X})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) &= \sigma^{-2}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - \boldsymbol{\beta}]^\top (\mathbf{X}^\top \mathbf{X}) [(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - \boldsymbol{\beta}] \\ &= \sigma^{-2}[(\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta})]^\top \mathbf{X}^\top \mathbf{X} [(\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta})] \\ &= \sigma^{-2}(\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta})^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{X})(\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}) \\ &= \sigma^{-2}(\mathbf{y} - \mathbf{X} \boldsymbol{\beta})^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) \\ &= \sigma^{-2} \boldsymbol{\epsilon}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ &= \sigma^{-2} \boldsymbol{\epsilon}^\top \mathbf{H} \boldsymbol{\epsilon} = \sigma^{-2} \boldsymbol{\epsilon}^\top \mathbf{P}_{\text{ppo}} \boldsymbol{\epsilon} \end{aligned}$$

Use the full rank SVD or eigenvalue decomposition on \mathbf{H}_n to obtain $\mathbf{H} = \mathbf{L}_n^\top \mathbf{L}_n$ with $\mathbf{L}_n \mathbf{L}_n^\top = \mathbf{I}_r$. Take a vector \mathbf{b} s.t. $\mathbf{b}^\top \mathbf{b} = 1 = \sum b_i^2$ and let $a_n^{(i)} = \mathbf{b}^\top \mathbf{L}_n^{(i)}$ (a scalar), then look at $\mathbf{b}^\top \mathbf{L}_n \boldsymbol{\epsilon} = \sum_{i=1}^n a_n^{(i)} \epsilon_i$.

$$|a_n^{(i)}| = |\mathbf{b}^\top \mathbf{L}_n^{(i)}| \leq (\mathbf{b}^\top \mathbf{b})^{1/2} (\mathbf{L}_n^{(i)\top} \mathbf{L}_n^{(i)})^{1/2} = 1 \times h_{ii}$$

By Hölder's inequality, where h_{ii} is the i th diagonal of the hat matrix. We need: $\max h_{ii} = \max \mathbf{L}_n^{(i)\top} \mathbf{L}_n^{(i)} \rightarrow 0$, so the condition needed is that diagonals of the hat matrix, the leverages, must go to zero as $n \rightarrow \infty$. Then $\sum_i (a_n^{(i)})^2 = \mathbf{b}^\top \mathbf{L}_n \mathbf{L}_n^\top \mathbf{b} = \mathbf{b}^\top \mathbf{b} = 1$ so we conclude that

$\mathbf{b}^\top \mathbf{L}_n \boldsymbol{\epsilon} \rightarrow N(0, \sigma^2)$ for every norm one vector \mathbf{b} . By properties of MVN, $\mathbf{L}_n \boldsymbol{\epsilon}$ has an asymptotic $N(\mathbf{0}, \sigma^2 \mathbf{I})$ distribution. Hence

$$\sigma^{-2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top (\mathbf{X}^\top \mathbf{X})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{\mathcal{D}} \chi_r^2.$$

In practice, we use t and F tests, so we need to ensure that these distributions also hold. We know that s^2 converges to σ^2 in probability so

$$s^{-2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top (\mathbf{X}^\top \mathbf{X})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{\mathcal{D}} \chi_r^2.$$

In general if \mathbf{C} is a p by q matrix of rank q and $\mathbf{C}^\top \boldsymbol{\beta}$ is estimable, then

$$s^{-2}(\mathbf{C}^\top \hat{\boldsymbol{\beta}} - \mathbf{C}^\top \boldsymbol{\beta})^\top [\mathbf{C}^\top (\mathbf{X}^\top \mathbf{X})^g \mathbf{C}]^{-1} (\mathbf{C}^\top \hat{\boldsymbol{\beta}} - \mathbf{C}^\top \boldsymbol{\beta}) \xrightarrow{\mathcal{D}} \chi_q^2.$$

However, the above assumes that s^2 is a “perfect” estimate of the variance. We generally prefer to not make that assumption and use the F distribution instead.

$$\frac{1}{qs^2}(\mathbf{C}^\top \hat{\boldsymbol{\beta}} - \mathbf{C}^\top \boldsymbol{\beta})^\top [\mathbf{C}^\top (\mathbf{X}^\top \mathbf{X})^g \mathbf{C}]^{-1} (\mathbf{C}^\top \hat{\boldsymbol{\beta}} - \mathbf{C}^\top \boldsymbol{\beta}) \xrightarrow{\mathcal{D}} F_{q, n-r}.$$

The critical assumption is that no single point is entirely “influential”, but the leverage of each is going to zero. We shall see that $\sum h_{ii} = \text{trace}(\mathbf{H}) = \text{rank}(\mathbf{X})$ which does not increase as n does.

More generally, the above also works with $\text{Var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{V}$ for known \mathbf{V} , we just need to do a little transformation. Take \mathbf{L} to be the Cholesky decomposition of \mathbf{V}^{-1} . That means $\mathbf{V}^{-1} = \mathbf{L}\mathbf{L}^\top$, with $\text{rank}(\mathbf{V}^{-1}) = \text{rank}(\mathbf{V}) = \text{rank}(\mathbf{L})$, since we cannot increase rank through multiplication. \mathbf{L} and \mathbf{L}^\top are therefore non-singular, and $\mathbf{L}^{\top-1} \mathbf{L}^{-1} = (\mathbf{L}\mathbf{L}^\top)^{-1} = \mathbf{V}$

Premultiply our linear model by \mathbf{L}^\top to get:

$$\mathbf{L}^\top \mathbf{y} = \mathbf{L}^\top \mathbf{X} \boldsymbol{\beta} + \mathbf{L}^\top \boldsymbol{\epsilon} = \mathbf{X}^* \boldsymbol{\beta} + \boldsymbol{\epsilon}^*$$

The mean of $\boldsymbol{\epsilon}^*$ is still $\mathbf{0}$, and it has variance-covariance matrix $\sigma^2 \mathbf{L}^\top \mathbf{V} \mathbf{L} = \sigma^2 \mathbf{L}^\top \mathbf{L}^{\top-1} \mathbf{L}^{-1} \mathbf{L} = \sigma^2 \mathbf{I}$, so the second formulation of our model fits the conditions for the CLT and as n gets big, the sampling distribution of $\hat{\boldsymbol{\beta}}$ converges to a normal distribution.

In practice, we never know exactly how big n needs to be to make the sampling distribution close enough to normality so that our inferences based on t and F distributions are valid. Problems in convergence occur when we have skewed residuals and/or outliers in the data. If residuals have short tails, the CLT will work well even if sample sizes are only moderate ($n - r = 30?$), otherwise we might need several hundred points to apply the CLT.

References

- Arnold, S. *The Theory of Linear Models and Multivariate Analysis*. Wiley (1981).
- Sen, A. K. and Srivastava, M. S. *Regression Analysis: Theory, Methods, and Applications*. Springer-Verlag Inc (1997).