

How do we select variables?

I can build a model

1. $\text{salary} = \beta_0 + \beta_1 \cdot \text{SEC} + \epsilon$
2. $\text{salary} = \beta_0 + \beta_1 \cdot \text{Career} + \epsilon$
3. $\text{salary} = \beta_0 + \beta_1 \cdot \text{Loss} + \epsilon$
4. $\text{salary} = \beta_0 + \beta_1 \cdot \text{SEC} + \beta_2 \cdot \text{Career} + \epsilon$
5. $\text{salary} = \beta_0 + \beta_1 \cdot \text{Loss} + \beta_2 \cdot \text{Career} + \beta_3 \cdot \text{Loss} \cdot \text{Career} + \epsilon$

In fact, I have $2^4 = 16$ possible models. Which one should I tell R to estimate?

What is the trade-off?

Adding more predictors to your model

1. Increase the fit of the model and reduce the total error, even if that variable is irrelevant.
2. Reduce the information because you need to use your data to estimate an extra parameter. Your parameters will be less precisely estimated, adding extra uncertainty to your forecasts!

Total Error + Penalty for complicated models

Criteria

$$\text{AIC} = N \log\left(\frac{\text{Sum of Squared Residuals}}{N}\right) + (k + 2) \cdot 2$$

$$\text{BIC} = N \log\left(\frac{\text{Sum of Squared Residuals}}{N}\right) + (k + 2) \cdot \log(N)$$

Rule: Find a model with the smallest value!