**What are the predictors? What are the outcomes?**

In general, the variables in a regression analysis are referred to as "dependent variables" and "independent variables", because a regression quantifies the the way that the dependent variables depend on the independent variables.

Start by assuming that the independent variables are the predictors; and that the dependent variables are the outcomes. Now, check if that makes sense: is it reasonable to say that the predictors will predict the outcome?

Not all regressions can be cast as a predictive model, which is why the standard terms for the variables are not "predictors" and "outcomes", but the more general "independent" and "dependent".

If the regression at hand is an explanatory one, and can't reasonably be viewed as a predictive model, it likely can't be incorporated into a management decision process. (Though, of course it can still be useful in other ways.)

**What is the total variation of outcome included in the regression data?**

Assuming that the regression being considered can be viewed as a predictive model, the next question is how large of a change in outcome can it predict?

**How much of the variation in the outcome is explained by the predictors?**

**How large was the sample size? What cut-off criteria did the researchers apply to size of effect?**

Typically, the only criteria that management researchers have for including a correlation in their discussion of the results of an analysis is statistical significance.

However, for practical application the effect size—how much of an impact a given correlation has on some outcome—is at least as important as statistical significance for determining what results are worth discussing. If not more important.

Whether or not a given correlation achieves statistical significance at some given confidence level is dependent on two things, and only these two things: the magnitude of the correlation, and the size of the sample. The way those two are related is that: while a large magnitude correlation will be statistically significant even with a small sample size, a small magnitude correlation will be statistically significant only with a large sample size.

Meanwhile the magnitude of a correlation is directly tied to the effect size, that is, to the impact of the outcome. Correlations with large magnitudes will have more of an impact on outcomes than correlations with a small magnitude.

For any given regression model, you can estimate what magnitude of correlation that would correspond to what amount of impact on the outcome. So, if you can say how much of a change in the outcome is worth discussing, you can also determine the minimum magnitude of correlation that is worth discussing.

For small sample size, whether or not a correlation is statistically significant is driven by the magnitude of the correlation. So, setting a criteria in addition to statistical significance for whether to include predictors in a discussion is not necessary. For large sample sizes, however, statistical significance is achieved by ever smaller magnitudes of correlation. In which case, it would be very helpful for researchers to distinguish between the correlations that are merely statistically significant, and those correlations that also have a non-negligible impact on the outcome.

As a reader, consider the sample size, and the smallest correlation that is being not just reported, but discussed in the paper. What kind of impact does that smallest, discussed correlation have on the outcome? Is it one that you think is worth discussing?

**Is the model simple enough to be useful? How many variables are there?**

Complex regression models, ones with more than, say, five variables, are probably not appropriate for practical application for a number of reasons.

Complex models are more difficult to understand, and therefore more likely to misunderstood. Complex models are more difficult to calculate, and therefore errors are more likely to pass unnoticed. Complex models are more likely to include Type I statistical errors

The last reason is the most important one, because it is why a simple model cannot be extracted from a more complex one in a reliably way.

You can estimate the number of Type I errors in an entire model based on the confidence level, a.k.a., the $\alpha$, of the analysis. If we say that the "error rate" of the analysis is the percentage excluded from our confidence, a.k.a., $1-\alpha$, then we can simply multiple the product of the number of predictors and the number of outcomes, by the that error rate. (Cohen 1990.)

For example, at a confidence level of 95% ($\alpha=0.95$), the error rate is 5% (1-0.95=0.05). If a regression has 14 predictors and 2 outcomes, we would expect at least one ($0.05 \cdot 14 \cdot 2 = 1.4$) erroneous correlation to be included. Indeed, at a 95% confidence level, any regression analysis with 10 or more independent/predictor variables is more likely than not to include one correlation that was erroneously identified as statistically significant.

In such a case, there is no way to know which correlation is the one that was incorrectly included among the statistically significant results. So, if you were to try to simplify a model with 10 predictors by only using three of them, there's a 30% chance that you're using a false predictor.

**What is the statistical power of the regression analysis?**