# Exercises:

1. For each part, indicate whether we would generally expect the performance of a flexible model (*i.e.* one with more degrees of freedom) to be better or worse than an inflexible method. Justify your answer.

    (a) The sample size (number of observations) $n$ is extremely large, and the number of predictors $p$ is small.

    (b) The number of predictors $p$ is extremely large, and the number of observations $n$ is small.

    (c) The relationship between the predictors and responses is highly non-linear.

    (d) The variance of the error terms $\epsilon$, is extremely high.

2. Explain whether each scenario is a classification or regression problem, and indicate whether we are most interested in inference or prediction. Finally, provide $n$ and $p$.

    (a) We collect a set of data on the top 500 firms in the US. For each firm we record profit, number of employees, industry and the CEO salary. We are interested in understanding which factors affect CEO salary.

    (b) We are considering launching a new product and wish to know whether it will be a *success* or a *failure*. We collect data on 20 similar products that were previously launched. For each product we have recorded whether it was a success or failure, price charged for the product, marketing budget, competition price, and ten other variables.

    (c) We are interest in predicting the % change in teh USD/Euro exchange rate in relation to the weekly changes in the world stock markets. Hence we collect weekly data for all of 2012. For each week we record the % change in the USD/Euro, the % change in the US market, the % change in the British market, and the % change in the German market.