## Exploratory Data Analysis (EDA)- Data Exploration

| Emp_Code | Gender | Date | New_Day | New_Month | New_Year |
|----------|--------|------|---------|-----------|----------|
| A001 | Male | 21-Sep-11 | 21 | 9 | 2011 |
| A002 | Female | 27-Feb-13 | 27 | 2 | 2013 |
| A003 | Female | 14-Nov-12 | 14 | 11 | 2012 |
| A004 | Male | 07-Apr-13 | 7 | 4 | 2013 |
| A005 | Female | 21-Jan-11 | 21 | 1 | 2011 |
| A006 | Male | 26-Apr-13 | 26 | 4 | 2013 |
| A007 | Male | 15-Mar-12 | 15 | 3 | 2012 |

There are various techniques to create new features. Let's look at the some of the commonly used methods:

- **Creating derived variables:** This refers to creating new variables from existing variable(s) using set of functions or different methods. Let's look at it through "**Titanic - Kaggle competition**". In this data set, variable age has missing values. To predict missing values, we used the salutation (Master, Mr, Miss, Mrs) of name as a new variable. How do we decide which variable to create? Honestly, this depends on business understanding of the analyst, his curiosity and the set of hypothesis he might have about the problem. Methods such as taking log of variables, binning variables and other methods of variable transformation can also be used to create new variables.
- **Creating dummy variables:** One of the most common application of dummy variable is to convert categorical variable into numerical variables. Dummy variables are also called Indicator Variables. It is useful to take categorical variable as a predictor in statistical models.  Categorical variable can take values 0 and 1. Let's take a variable 'gender'. We can produce two variables, namely, "**Var_Male**" with values 1 (Male) and 0 (No male) and "**Var_Female**" with values 1 (Female) and 0 (No Female). We can also create dummy variables for more than two classes of a categorical variables with n or n-1 dummy variables.

| Emp_Code | Gender | Var_Male | Var_Female |
|----------|--------|----------|------------|
| A001 | Male | 1 | 0 |
| A002 | Female | 0 | 1 |
| A003 | Female | 0 | 1 |
| A004 | Male | 1 | 0 |
| A005 | Female | 0 | 1 |
| A006 | Male | 1 | 0 |
| A007 | Male | 1 | 0 |

COMPLETE & CONTINUE  →