

# Kickstarter Project: Data Cleaning Steps

## 1.1 General Data Cleaning Steps:

- Convert launched\_date and deadline\_date into month and year
- Convert state (whether project failed or succeeded) into a numeric variable with failed as zero and successful as 1
- Calculate the difference between start launch\_date and deadline\_date
- Create new variable goal\_per\_day or divide the usd\_goal by number of days between the deadline\_date and launch\_date

## 1.2 Imputation using KNN

- For the countries variable, around 3K observations are missing
  - Replace 'N,0' variable with 'nan'
  - Convert all the countries into label encoding (replace category name with number; this is the only way KNN will process the data)
  - Apply KNN algorithm

## 1.3 Part of Speech Tagging

- In order to extract information about part of each of each word, we would need to use nltk's part of speech tagger (pos\_tag)
- Extract the part of speech tag from the output of pos\_tag
- Then , apply set to get only unique post\_tags and join them
- Remove any extraneous symbols (:, #)
- Create dummies for each unique post\_tag