## Data Cleaning

Data Cleaning is the process of turning raw data into a clean and analyzable data set. "Garbage in, garbage out." Make sure garbage doesn't get put in.

### Errors vs. Artifacts
1. **Errors:** information that is lost during acquisition and can never be recovered e.g. power outage, crashed servers
2. **Artifacts:** systematic problems that arise from the data cleaning process. these problems can be corrected but we must first discover them

### Data Compatibility
Data compatibility problems arise when merging datasets. Make sure you are comparing "apples to apples" and not "apples to oranges". Main types of conversions/unifications:
- **units** (metric vs. imperial)
- **numbers** (decimals vs. integers),
- **names** (John Smith vs. Smith, John),
- **time/dates** (UNIX vs. UTC vs. GMT),
- **currency** (currency type, inflation-adjusted, dividends)

### Data Imputation
Process of dealing with missing values. The proper methods depend on the type of data we are working with. General methods include:
- Drop all records containing missing data
- Heuristic-Based: make a reasonable guess based on knowledge of the underlying domain
- Mean Value: fill in missing data with the mean
- Random Value
- Nearest Neighbor: fill in missing data using similar data points
- Interpolation: use a method like linear regression to predict the value of the missing data

### Outlier Detection
Outliers can interfere with analysis and often arise from mistakes during data collection. It makes sense to run a "sanity check".

### Miscellaneous
Lowercasing, removing non-alphanumeric, repairing, unidecode, removing unknown characters

*Note:* When cleaning data, always maintain both the raw data and the cleaned version(s). The raw data should be kept intact and preserved for future use. Any type of data cleaning/analysis should be done on a copy of the raw data.

## Feature Engineering

Feature engineering is the process of using domain knowledge to create features or input variables that help machine learning algorithms perform better. Done correctly, it can help increase the predictive power of your models. Feature engineering is more of an art than science. FE is one of the most important steps in creating a good model. As Andrew Ng puts it:

*"Coming up with features is difficult, time-consuming, requires expert knowledge. 'Applied machine learning' is basically feature engineering."*

### Continuous Data
**Raw Measures**: data that hasn't been transformed yet
**Rounding**: sometimes precision is noise; round to nearest integer, decimal etc..
**Scaling**: log, z-score, minmax scale
**Imputation**: fill in missing values using mean, median, model output, etc..
**Binning**: transforming numeric features into categorical ones (or binned) e.g. values between 1-10 belong to A, between 10-20 belong to B, etc.
**Interactions**: interactions between features: e.g. subtraction, addition, multiplication, statistical test
**Statistical**: log/power transform (helps turn skewed distributions more normal), Box-Cox
**Row Statistics**: number of NaN's, 0's, negative values, max, min, etc
**Dimensionality Reduction**: using PCA, clustering, factor analysis etc

### Discrete Data
**Encoding**: since some ML algorithms cannot work on categorical data, we need to turn categorical data into numerical data or vectors
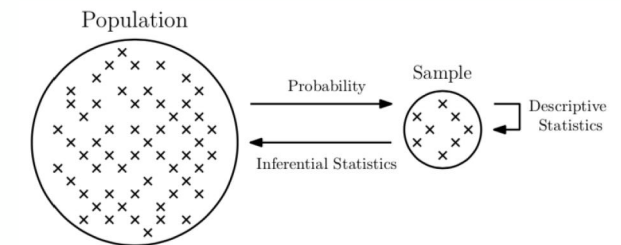**Ordinal Values**: convert each distinct feature into a random number (e.g. [r,g,b] becomes [1,2,3])
**One-Hot Encoding**: each of the m features becomes a vector of length m with containing only one 1 (e.g. [r, g, b] becomes [[1,0,0],[0,1,0],[0,0,1]])
**Feature Hashing Scheme:** turns arbitrary features into indices in a vector or matrix
**Embeddings**: if using words, convert words to vectors (word embeddings)

## Statistical Analysis

Process of statistical reasoning: there is an underlying population of possible things we can potentially observe and only a small subset of them are actually sampled (ideally at random). Probability theory describes what properties our sample should have given the properties of the population, but ***statistical inference*** allows us to deduce what the full population is like after analyzing the sample.



### Sampling From Distributions
**Inverse Transform Sampling** Sampling points from a given probability distribution is sometimes necessary to run simulations or whether your data fits a particular distribution. The general technique is called *inverse transform sampling* or Smirnov transform. First draw a random number $p$ between [0,1]. Compute value x such that the CDF equals $p$: $F_X(x) = p$. Use x as the value to be the random value drawn from the distribution described by $F_X(x)$.

**Monte Carlo Sampling** In higher dimensions, correctly sampling from a given distribution becomes more tricky. Generally want to use Monte Carlo methods, which typically follow these rules: define a domain of possible inputs, generate random inputs from a probability distribution over the domain, perform a deterministic calculation, and analyze the results.