

1. Techniques:
  - a. Plotting: Scatter plot (2D, 3D, Pair plot, PDF, CDF, box-plots, violin plots, Contour plots)  
Refer: <https://seaborn.pydata.org/generated/seaborn.violinplot.html>
  - b. Probability and Stats:
    - i. Counts, mean, std-dev, median, percentiles, IQR
    - ii. Distributions & skewness
    - iii. Correlations
    - iv. Hypothesis testing
  - c. Visualizing High Dimensional data: PCA, t-SNE
  - d. Model based: Rule based, Linear and logistic regression, Feature importance, Feature collinearity.
  - e. Cluster analysis.
2. Questions [Sherlock holmes]
  - a. Varies widely based on the dataset and problem being solved.
  - b. ART vs Science, {Practice, Practice, Practice}
  - c. Always look at the raw data and not just the aggregate numbers.
  - d. High level stats:
    - i. Number of data points & features.
    - ii. Imbalance in labels for classification.
    - iii. Distribution (skewness) of  $y_i$  for regression.
  - e. Feature wise analysis:
    - i. Categorical feature: distribution of categories
    - ii. Real-valued: distribution of the feature.
    - iii. Missing values & imputation.
    - iv. Outliers in feature values.
  - f. Feature vs output
    - i. Correlation
    - ii.  $P(y_i=1|f_j=k)$
    - iii. Build a model with just one feature.
  - g. What other features might work well?
    - i. Group-by + count on raw data.
    - ii. Binning features using Decision Trees
    - iii. Interaction variables using Decision Trees.
    - iv. Mathematical transforms: log, exp, sqrt, ^2, box-cox
    - v. Normalization, Standardization, one-hot encoding
    - vi. Matrix factorization based features.
    - vii. Autoencoder based features.
    - viii. Clustering of features
    - ix. Different encodings: Word2Vec, TF-IDF etc
  - h. High dimensional viz:
    - i. Are there clusters of points in a region where we are performing badly.
    - ii. Why I am observing more error while classifying class\_i and class\_j?

- iii. What types of examples are messing up my model?

Refer:

[https://www.google.com/search?q=t+SNE+for+model+diagnosis&safe=active&source=Inms&tbm=isch&sa=X&ved=0ahUKEwjDopr0mIjhAhUBPY8KHZIIIBdgQ\\_AUIDygC&biw=1309&bih=725#imgsrc=WMgzF6Rp-381yM:](https://www.google.com/search?q=t+SNE+for+model+diagnosis&safe=active&source=Inms&tbm=isch&sa=X&ved=0ahUKEwjDopr0mIjhAhUBPY8KHZIIIBdgQ_AUIDygC&biw=1309&bih=725#imgsrc=WMgzF6Rp-381yM:)

[https://cs.stanford.edu/people/karpathy/cnnembed/cnn\\_embed\\_1k.jpg](https://cs.stanford.edu/people/karpathy/cnnembed/cnn_embed_1k.jpg)

<https://lvdmaaten.github.io/tsne/>

### 3. EDA on

- a. Text data: word counts,  $P(y_i|w_j)$  like in Naive Bayes, linear models and feature importance.
- b. Time-series data: repetitions, fourier transforms, moving averages.

Refer:

<https://image.slidesharecdn.com/timeseriesrtalk-160329213057/95/time-series-analysis-fpp-package-2-638.jpg?cb=1459287213>

- c. Image data: CNN featurizations + tSNE+ look at the raw images.