**Prof. Pedram Jahangiry**
Utah State University: **Huntsman School of Business**

| Machine Learning Type | Model | Model Type and use case | Description | Pros | Cons | Hyperparameters |
|---|---|---|---|---|---|---|
| **Supervised** | **Linear regression** | Linear - Parametric<br><br>Used for regression only | Finds the "best fit" through all the data points. | - highly interpretable (giving significance results)<br>- very fast training because of closed form solution<br>- no hyperparameter tuning required | - validity of linear regression assumptions<br>- cannot capture complex relationships | none |
| | **Polynomial regression** | Linear - Parametric<br><br>Used for regression only | Extending linear regression model to capture non-linearities | - interpretable for low values of d (giving significance results)<br>- Can capture polynomial relationships | - need to choose the right polynomial degree<br>- notorious tail behavior (sensitive to outliers) | d: degree of polynomial |
| | **Penalized regression: Ridge, LASSO and Elastic Net** | Linear - Parametric<br><br>Used for regression only | Linear method that penalizes irrelevant features using regularization<br>L1 regularization: LASSO<br>L2 regularization: Ridge<br>combining L1 and L2: Elastic net | - Can be used for feature selection (reducing the dimension of the feature space)<br><br>- interpretable | - Requires feature scaling | penalty (how much to penalize the parameters)<br><br>L1 ratio: ration between L1 and L2 regularization |
| | **Logistic regression** | Linear - Parametric<br><br>Used for classification only | Basically the adaptation of linear regression to classification problems. | - probabilitsitc model (the outputs are probabilities)<br>- highly interpretable (giving significance results)<br>- easy to understand<br>- fast and efficient | - validity of linear regression assumptions<br>- sensitive to extreme values<br>- cannot capture complex relationships | - the same as penalized regression if regularization is used |
| | **KNN** | Non-linear - Non-parametric<br><br>Used for both regression and classification | Make prediction for a new observation by finding similarities ("nearness") between it and its k nearest neighbors in the existing dataset. | - Intuitive and simple<br>- Easy to implement for multi class problem<br>- Few parameters/hyper parameters<br>- No assumption (non parametric) | - Choice of K<br>- Slow (memory based approach)<br>- Curse of dimensionality<br>- Hard to interpret<br>- Requires feature scaling<br>- Not good with multiple categorical features | - K value<br>- distance metrics |
| | **SVM** | Kernel basis (non-linear)<br>Linear SVM is parametric<br>Kernel SVM is non-parametric<br><br>Used for both regression and classification | Uses a kernel to transform the feature space to linearly separable boundaries | - SVM can be memory efficient! uses only a subset of the training data (support vectors)<br>- Can handle non linear data sets<br>- Can handle high dimensional spaces (even when D>N)<br>- Linear SVM are not very sensitive to overfitting (soft margin; regularization)<br>- Can have high accuracy (even compared to NN) | - Requires feature scaling<br>- No probability outcome!<br>- Does not perform well with noisy data<br>- Limited interpretability (specially for Kernel SVM)<br>- Memory intensive: Long training time when we have large data sets. | - Kernel: linear, rbf, poly, ...<br>- C: Cost of misclassification<br>- Gamma (for rbf): how far the influence reach<br>- d (for poly): degree of polynomial |
| | **Decision Trees** | Tree-based (non-linear)<br>Non-parametric<br><br>Used for both regression and classification | - Progressively divide data sets into smaller data groups based on a descriptive feature , until they reach sets that are small enough to be described by some label | -Easy to interpret and visualize<br>- Can easily handle categorical data without the need to create dummy variables<br>- Can easily capture Non linear patterns<br>- Can handle data in its raw form (no preprocessing needed ).<br>- No assumption (non parametric)<br>- Can handle colinearity efficiently | - Sensitive to noisy data. It can overfit noisy data. Small variations in data can result in the different decision tree<br>- Can lead to overfitting<br>- Poor level of predictive accuracy | - Max tree depth, min samples per leaf (node),  min samples split<br>- Cost complexity alpha<br>- Criterion: gini/entropy/ ... |
| | **Random Forest** | Ensemble method (non-linear)<br>Non-parametric<br><br>Used for both regression and classification | - Many trees are created on bootstrapped data and combined using averaging. | All the advantages of Decision Trees +<br>- Typically more accurate<br>- Avoid overfitting by reducing the model variance.<br>- very flexible and parallelizable!<br>- No data preprocessing (no feature scaling)<br>- Great with high dimensionality | - no interpretability<br>- complexity<br>- many hyper parameters<br>- slow on large data sets | DTs parameters +<br>- m: subset of features<br>- B: number of bootstrapped trees |
| | | Ensemble method (non-linear) | | All the advantages of Random Forests + | - no interpretability | RF parameters + |

| | | | | | | |
|---|---|---|---|---|---|---|
| **Unsupervised** | **Boosting (XGboost)** | Non-parametric<br><br>Used for both regression and classification | - Implements boosting to build decision trees of weak prediction models and generalizes using a loss function. | - Regularization for avoiding overfitting<br>- Efficient handling of missing data<br>- In-built cross validation capability<br>- Cache awareness and out-of-core computing<br>- Tree pruning using depth-first approach<br>- Parallelized tree building | - many hyper parameters | - Regularization terms |
| | **Principle Component Analysis (PCA)** | Non- Parametirc<br><br>Used for dimension reduction | Principal components are vectors that define a new coordinate system in which the first axis goes in the direction of the highest variance in the data. The second PC is orthogonal to PC1 and etc. | - Reducing the number of features to the most relevant predictors is very useful in general.<br>- Dimension reduction facilitates the data visualization in two or three dimensions.<br>- Before training another supervised or unsupervised learning model, it can be performed as part of EDA to identify patterns and detect correlations .<br>- Machine learning models are quicker to train , tend to reduce overfitting (by avoiding the curse of dimensionality), and are easier to interpret if provided with lower dimensional datasets. | - Hard to interpret<br><br>- Requires feature scaling | None |
| | **K-Means** | Cab be both Parametric and Non-Parametirc<br><br>Used for clustering the data | - Uses a measure of similarity to detect groups within data set: K means is an algorithm that repeatedly partitions observations into a fixed, pre-specified number of clusters | - Simple to understand<br>- The k means algorithm is fast and works well on very large datasets<br>- Can help visualize the data and facilitate detecting trends or outliers. | - Need to choose k before running the algorithm<br>- Requires feature scaling<br>- Poor performance with clusters of irregular shapes<br>- Not applicable for categorical data<br>- Unable to handle noisy data | K: Number of clusters<br>- Distance metrics |
| | **Hierarchical clustering** | Cab be both Parametric and Non-Parametirc<br><br>Used for clustering the data | - Uses a measure of similarity to detect groups within data set: Hierarchical clustering is an iterative procedure used to build a hierarchy of clusters | - The optimal number of clusters can be obrained by the model itself, | - The choice of distance metrics and linkage methods can be tricky<br>- Requires feature scaling | - Distance metrics<br>- Linkage methods |