

<p>Predictive Models</p> <ul style="list-style-type: none"> Inference: understand the relationship between predictors and an outcome Prediction: modeling for prediction which prioritizes prediction accuracy 	$Y = \beta_0 + \beta_1 X + \epsilon$ <p style="text-align: center;">dependent (outcome) Intercept slope (coefficients) independent (predictor)</p>	<p>Factors effect coefficient</p> <ul style="list-style-type: none"> Scaling: standardizing variables help compare Interaction effects: combined effect differs from individual effects Multicollinearity: misleading if predictors are highly correlated → correlation matrix / VIF 																		
<p>Model Assessment</p> <p>Coefficient of Determination (R^2): the proportion of the variation in Y that is explained by variation in X, measures fit, not correctness</p>	<ul style="list-style-type: none"> Coefficients: magnitude - strength of the effect / sign - direction of the relationship Intercept: the expected value of Y when all predictors are 0 P-value: whether the coefficient is statistically significant 	<p>Variance Inflation Factor (VIF): measures how much the variance of a model coefficient is inflated due to multicollinearity</p> $VIF(\beta_j) = \frac{1}{1 - R^2_{X_j X_{-j}}}$ <p>where $R^2_{X_j X_{-j}}$ is the R^2 of a regression model with X_j as the response variable and predictors $X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_p$</p>																		
<table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="text-align: left; padding: 2px;">Measure</th> <th style="text-align: left; padding: 2px;">Formula</th> </tr> </thead> <tbody> <tr> <td style="padding: 2px;">Total Sum of Squares</td> <td style="padding: 2px;">$\sum (y_i - \bar{y})^2$</td> </tr> <tr> <td style="padding: 2px;">Residual Sum of Squares (RSS) / Sum of Squared Errors (SSE)</td> <td style="padding: 2px;">$\sum_{i=1}^n (y_i - \hat{y}_i)^2$</td> </tr> <tr> <td style="padding: 2px;">Mean Squared Error (MSE)</td> <td style="padding: 2px;">SSE/n</td> </tr> <tr> <td style="padding: 2px;">Root mean squared error (RMSE)</td> <td style="padding: 2px;">\sqrt{MSE}</td> </tr> <tr> <td style="padding: 2px;">Average Squared Error (ASE)</td> <td style="padding: 2px;">SSE/n</td> </tr> <tr> <td style="padding: 2px;">Residual Standard Error (RSE)</td> <td style="padding: 2px;">$\sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-p}}$</td> </tr> <tr> <td style="padding: 2px;">Coefficient of Determination (R^2)</td> <td style="padding: 2px;">$0 \leq R^2 \leq 1$ 0: No changes in Y explained by X 1: All changes in Y explained by X (perfect fit)</td> </tr> <tr> <td style="padding: 2px;"></td> <td style="padding: 2px;">$1 - \frac{RSS}{TSS}$</td> </tr> </tbody> </table>	Measure	Formula	Total Sum of Squares	$\sum (y_i - \bar{y})^2$	Residual Sum of Squares (RSS) / Sum of Squared Errors (SSE)	$\sum_{i=1}^n (y_i - \hat{y}_i)^2$	Mean Squared Error (MSE)	SSE/n	Root mean squared error (RMSE)	\sqrt{MSE}	Average Squared Error (ASE)	SSE/n	Residual Standard Error (RSE)	$\sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-p}}$	Coefficient of Determination (R^2)	$0 \leq R^2 \leq 1$ 0: No changes in Y explained by X 1: All changes in Y explained by X (perfect fit)		$1 - \frac{RSS}{TSS}$	<p>$H_0:$ No relationship between X and Y (coefficient = 0) $H_1:$ No relationship between X and Y (coefficient ≠ 0)</p>	
Measure	Formula																			
Total Sum of Squares	$\sum (y_i - \bar{y})^2$																			
Residual Sum of Squares (RSS) / Sum of Squared Errors (SSE)	$\sum_{i=1}^n (y_i - \hat{y}_i)^2$																			
Mean Squared Error (MSE)	SSE/n																			
Root mean squared error (RMSE)	\sqrt{MSE}																			
Average Squared Error (ASE)	SSE/n																			
Residual Standard Error (RSE)	$\sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-p}}$																			
Coefficient of Determination (R^2)	$0 \leq R^2 \leq 1$ 0: No changes in Y explained by X 1: All changes in Y explained by X (perfect fit)																			
	$1 - \frac{RSS}{TSS}$																			
	<p>Linear Regression Model Assumptions</p> <ul style="list-style-type: none"> Linearity of response-predictor relationship Independence of residuals Normal Distribution of residuals → QQ-plot Equal variance of residuals (Homoscedasticity) 																			
<p>Odds: $p=0.75$, odds=3, success is 3 times more than failure</p> <p>Odds ratios (OR): OR>1, more likely in group 1</p>	$\text{Odds Ratio} = \frac{\text{Odds in Group 1}}{\text{Odds in Group 2}} \quad Odds = \frac{p}{1-p}$																			
<p>Tree-based Model</p> <p>Pros:</p> <ol style="list-style-type: none"> Understand decision path Identify key predictors: top - more important Detect interaction effect: split on different features at different levels, linear model (feature engineering) Threshold effect: keep split at same value forming key decision boundary Classify group <p>Cons:</p> <ol style="list-style-type: none"> Overfitting and noise: deep, big, complex (solution: pruning, max tree depth, min node size) Multicollinearity is not considered Lack of statistical significance Can not make inference beyond observation 	<p>Global Interpretability: how model make decisions across entire data</p> <p>SHAP (Shapley Additive Explanations): How much each feature contributes to model across all samples</p> <p>PDPs (Partial Dependence Plots): Show average effect of a feature on predictions. How predicted values change as one feature varies, averaging over all other</p> <p>ICE (Individual Conditional Expectation Plots): show feature effects at the <i>individual</i> level</p>																			
<p>Feature Importance Score: measuring how much a feature reduces impurity across all splits</p> <p>→ calculate decrease in impurity for that feature and normalize to get relative importance score</p> <p>*Random forest: sum the decrease across all trees</p>	<p>Local Interpretability: explain individual prediction</p> <p>SHAP: computes the marginal contribution of each feature across all permutations, can be applied to any predictive model, providing consistent, additive explanations for both linear and non-linear models</p> <p>LIME (Local Interpretable model-agnostic explanations): fits a simple surrogate model (e.g., linear) around a single prediction to approximate the model's local behavior near that point</p>																			
<p>SHAP force plot: the base value is 0.45, and the final predicted probability is 0.75. The biggest SHAP contributors are: Credit Score: +0.18, Debt-to-Income Ratio: +0.12.</p> <p>The prediction started at the base value (0.45) and increased primarily due to Credit Score and Debt-to-Income Ratio. Together, these features pushed the prediction up by +0.30, leading to a final probability of 0.75. This shows how specific features influenced this individual prediction.</p>	<table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="text-align: left; padding: 2px;">Interpretability Type</th> <th style="text-align: left; padding: 2px;">Scope</th> <th style="text-align: left; padding: 2px;">Answers the Question</th> <th style="text-align: left; padding: 2px;">SHAP Tools</th> </tr> </thead> <tbody> <tr> <td style="padding: 2px;">Global</td> <td style="padding: 2px;">Entire dataset</td> <td style="padding: 2px;">"Which features matter most overall?"</td> <td style="padding: 2px;">Summary plot, bar plot</td> </tr> <tr> <td style="padding: 2px;">Local</td> <td style="padding: 2px;">Single prediction</td> <td style="padding: 2px;">"Why was this prediction made?"</td> <td style="padding: 2px;">Force plot, waterfall plot</td> </tr> </tbody> </table>	Interpretability Type	Scope	Answers the Question	SHAP Tools	Global	Entire dataset	"Which features matter most overall?"	Summary plot, bar plot	Local	Single prediction	"Why was this prediction made?"	Force plot, waterfall plot							
Interpretability Type	Scope	Answers the Question	SHAP Tools																	
Global	Entire dataset	"Which features matter most overall?"	Summary plot, bar plot																	
Local	Single prediction	"Why was this prediction made?"	Force plot, waterfall plot																	
<p>SHAP summary plot</p> <p>Right (positive): feature increase prediction</p> <p>Left (negative): feature decrease prediction</p> <p>Widely spread: high variability in impact</p> <p>Same side blue(red) dots: low(high) feature have consistent effect</p> <p>Mixed color: complex/nonlinear relationship</p>	<p>Use Case</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="text-align: left; padding: 2px;">Use Case</th> <th style="text-align: left; padding: 2px;">Preferred Method</th> </tr> </thead> <tbody> <tr> <td style="padding: 2px;">Explaining a single prediction</td> <td style="padding: 2px;">LIME</td> </tr> <tr> <td style="padding: 2px;">Auditing model behavior globally</td> <td style="padding: 2px;">SHAP</td> </tr> <tr> <td style="padding: 2px;">Debugging a "high-stakes" model</td> <td style="padding: 2px;">SHAP</td> </tr> <tr> <td style="padding: 2px;">Fast local explanations</td> <td style="padding: 2px;">LIME</td> </tr> <tr> <td style="padding: 2px;">Simple, interpretable explanations</td> <td style="padding: 2px;">LIME</td> </tr> </tbody> </table>	Use Case	Preferred Method	Explaining a single prediction	LIME	Auditing model behavior globally	SHAP	Debugging a "high-stakes" model	SHAP	Fast local explanations	LIME	Simple, interpretable explanations	LIME							
Use Case	Preferred Method																			
Explaining a single prediction	LIME																			
Auditing model behavior globally	SHAP																			
Debugging a "high-stakes" model	SHAP																			
Fast local explanations	LIME																			
Simple, interpretable explanations	LIME																			
	<p>PDP (age predict health risk): The plot shows a flat line up to Age 50, then a sharp upward trend. The model predicts little to no effect of age on health risk below 50, but increased risk after 50, suggesting a threshold or non-linear effect that the model captures.</p>																			
	<p>ICE plot (for income): most curves slope upward, but 10% of individuals show downward-sloping lines.</p> <p>For most individuals, higher income increases predicted approval probability (positive effect), but for a small group, the effect is negative, possibly due to interactions with other features. Heterogeneous effects and a possible interaction or non-monotonic relationship.</p>																			
	<ol style="list-style-type: none"> ICE lines differ from avg PDP line: heterogeneous effects Around 700, the ICE and PDP both show a sharp increase: threshold highest mean SHAP value: contributes the most to prediction decisions upward trend of avg PDP line: positive relationship 																			

Clustering Quality

Hopkins statistic: To find whether the data has meaningful clusters, measures the probability that the data is generated by a uniform random distribution

- Random sample n points and find distance from each point to its nearest neighbor (x_i)
- Random sample n points (not datapoints) and find distance from each point to its nearest neighbor (y_i)

$0 < \text{no clusters} < 0.5$ (uniformly distributed) < cluster < 1

Cohesion: How close the point is to other points in the same cluster (min intra-cluster)

Separation: How far the point is from points in the nearest different cluster (max inter-cluster)

$$H = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i + \sum_{i=1}^n y_i}$$

Point	X	Y	Dist_to_A1	Dist_to_B1	Dist_to_C1	Assigned Cluster
P1	2	10	0	9.2	9	A1
P2	4	8	4	5.2	9	A1
P3	5	5	8	1.2	7	B1
P4	6	4	10	0.8	7	B1
P5	8	4	12	2.8	9	B1
P6	1	2	9	7.8	0	C1
P7	7	3	12	2.8	7	B1

Cluster	X	Y
A	2	10
B	6	4.8
C	1	2

Aspect	Silhouette Score	Elbow Method (inertia)
What it measures	How well-separated and cohesive the clusters are	How compact (tight) the clusters are
Based on	Distance between each point and: <ul style="list-style-type: none">its own cluster center (cohesion)the nearest other cluster center (separation) $s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$	Sum of squared distances between each point and its own cluster center
Output range	-1 to +1 (higher is better) <small>weak < 0.2 < acceptable < 0.5 < good</small>	Always positive; lower inertia means tighter clusters
Goal	Maximize the average silhouette score	Finding a good k value where adding more clusters doesn't help much
Limitations	Can be misleading when clusters vary in shape	Always decreases with more clusters – needs interpretation

K-means (sensitive to outliers)

1. Select k cluster centers
2. Assign each object to closest cluster center
3. Update cluster centers (**Average of new cluster**)
4. Repeat 2 and 3 until convergence

K-medoids

3. For each cluster, attempt to swap a medoid with another point in cluster to see if it minimizes the total distance within the cluster

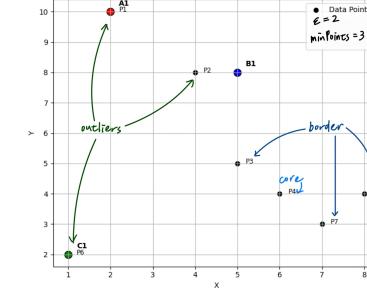
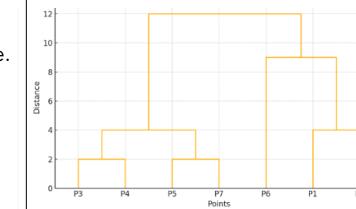
Hierarchical Clustering

- Agglomerative: individual cluster and merge each step
- Divisive: all-inclusive cluster and split each step

Calculating distance between clusters (Linkage)

- Single: closest entities between clusters
- Complete: furthest entities between clusters
- Group average: average of all pairs
- Centroid: centroid between clusters
- Manhattan distance (given matrix)

	K-means (centroids)	K-medoids
Center	Average, not real point	Center, real data point
Outliers	Sensitive	Robust
Distance-based	Without scaling, variables with larger numeric ranges will dominate the distance calculation and bias the clustering.	



Model-based Clustering

Data is generated from a mixture of underlying probability distributions, typically Gaussian. Each cluster corresponds to 1 component in the mixture.

- **Expectation (E):** computes the probability that each data point belongs to each cluster, based on current estimates of the parameters.
- **Maximization (M):** updates the parameters of each distribution (mean, covariance, and mixing proportion) using the probabilities in the E-step
- **Bayesian Information Criterion (BIC)/Akaike Information Criterion (AIC):** comparing model and choose lower values that indicating better fit with fewer parameters

Association Rule Mining

- **Support:** the proportion of transactions that contain a particular itemset.
- **Confidence:** the conditional probability that a transaction containing the antecedent also contains the consequent.
- **Lift:** how much the variables deviate from being independent, adjusted for the baseline frequency of the consequent and measures how much more likely the consequent is when the antecedent is present, relative to chance.
- A rule may have high confidence simply because the consequent is common, not because there's a real association.
- Rare item problem: Rare but important items (e.g., expensive products) may be filtered out by the minimum support threshold.

$$\text{sup}(X) = \frac{|\{T_j \in T; X \subseteq T_j\}|}{|T|} \quad \text{conf}(X \Rightarrow Y) = \frac{\text{sup}(X \cup Y)}{\text{sup}(X)} \quad \text{lift}(X \Rightarrow Y) = \frac{\text{conf}(X \Rightarrow Y)}{\text{sup}(Y)}$$

Apriori algorithm:

1. Generate all frequent itemsets by pruning infrequent ones based on the minimum support threshold.

***Support Monotonicity:** The support of an itemset cannot be greater than the support of any of its subsets → pruning of candidate itemsets, if a subset is infrequent, all its supersets can be ignored, improving efficiency

If $X_1 \subset X_2$, $\text{support}(X_1) \geq \text{support}(X_2)$

2. From those frequent itemsets, generate association rules that meet the minimum confidence threshold.

***Confidence Monotonicity:** $\text{conf}(\{A, B\} \rightarrow \{C\}) \geq \text{conf}(\{A\} \rightarrow \{B, C\})$

Density-based Clustering (DBSCAN)

Clusters are defined as areas of high point density separated by regions of low density, rather than by distance to centroids or hierarchical structure.

- **Core points:** have at least minPts within ϵ
 - **Border points:** have fewer than minPts within ϵ but within ϵ of core point
 - **Noise points:** are neither core nor border (i.e., outliers)
- * ϵ (epsilon): the radius defining a neighborhood around a point.
* minPts : the minimum number of points required within ϵ for a point to be considered a core point.

Pros: explicitly labels low-density points as noise, while k-means forces all points into clusters, making it sensitive to outliers. Not require to define number of clusters

Cons: A single ϵ value cannot simultaneously capture dense and sparse clusters, it may over-cluster sparse areas or merge dense clusters incorrectly.

T1: {Milk, Bread, Butter}
T2: {Bread, Butter}
T3: {Milk, Bread}
T4: {Milk, Butter}
T5: {Bread, Butter}

- $\text{sup}(\{\text{Milk, Bread}\}) = 2/5$
- $\text{conf}(\{\text{Milk}\} \rightarrow \{\text{Bread}\}) = 2/3$
- Lift = $\text{Confidence}(\{\text{Milk}\} \rightarrow \{\text{Bread}\}) / \text{Support}(\{\text{Bread}\}) = 10/12$

Assume the minimum support threshold is 0.4. Given the following itemsets and their supports: {A}=0.6 {B}=0.5 {C}=0.3 {A,B}=0.4 {A,C}=0.25 {B,C}=0.2

Which itemsets are frequent? {A}, {B}, {A,B}

Which itemsets will be pruned by the Apriori algorithm? {C}, {A,C}, {B,C}

{Green Tea Cup} → {Pink Tea Cup}

Support = 0.03: 3% of all transactions contain both items.

Confidence = 0.62: When the Green Tea Cup is purchased, there's a 62% chance the Pink one is also purchased.

Lift = 16.4: The co-occurrence of these items is 16.4 times more likely than if they were independent—a strong association.

	Movie A	Movie B	Movie C	Movie D
User1	5	3	4	
User2	3	1	2	3
User3	4	3	4	

If the cosine similarity between User A and User B is 0.92, and between User A and User C is 0.60, and both have rated the same movie with ratings:

$$\text{User B: } \frac{0.92 \times 4 + 0.6 \times 5}{\sqrt{4^2 + 5^2}} = 4.3$$

$$\text{User C: } \frac{0.92 \times 4 + 0.6 \times 5}{\sqrt{4^2 + 5^2}} = 4.3$$

Collaborative Filtering

Normalization before computing similarity: different rating scales, removing individual user bias and make comparison meaningful

Item X	Item Y	Item Z
U1	1	0
U2	1	1

Compute the Jaccard similarity between U1 and U2 = 1/3

$$\text{Jaccard Similarity} = \frac{|I_A \cap I_B|}{|I_A \cup I_B|}$$

$$|I_A \cap I_B|$$

$$|I_A \cup I_B|$$