

Feature Selection

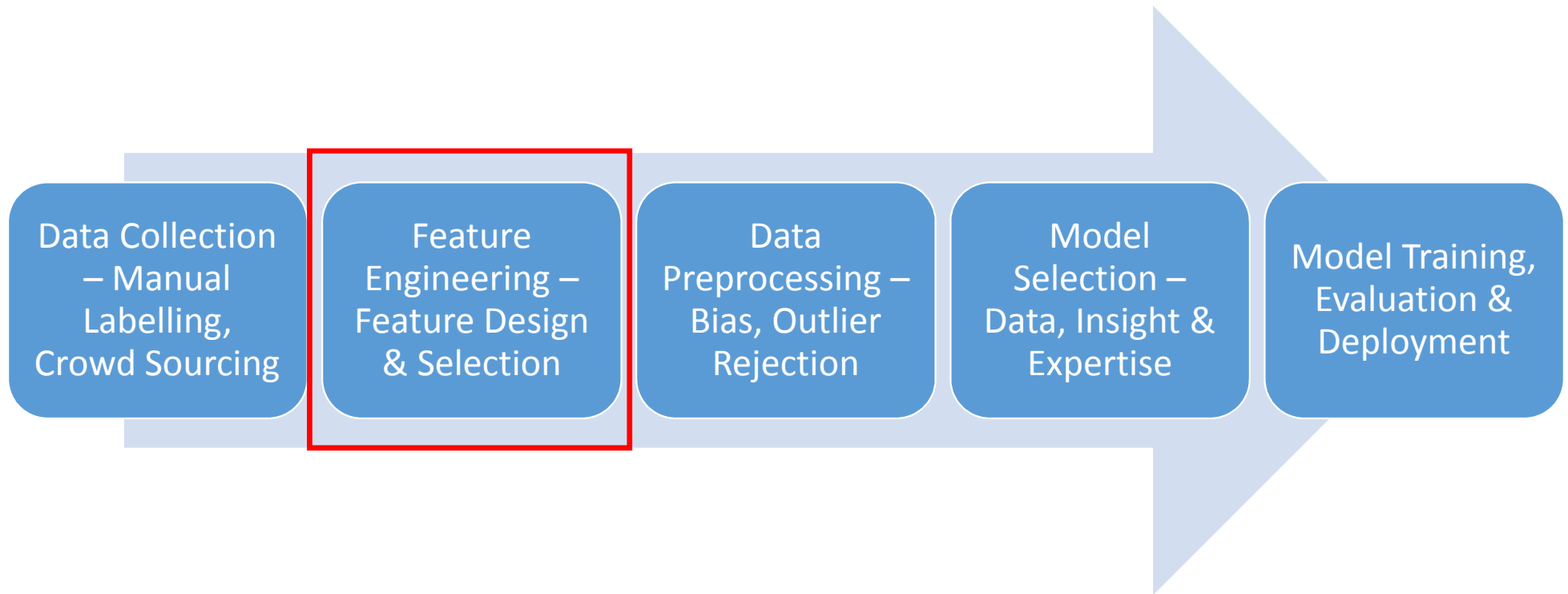


Prithwijit Guha
Dept. of EEE, IIT Guwahati

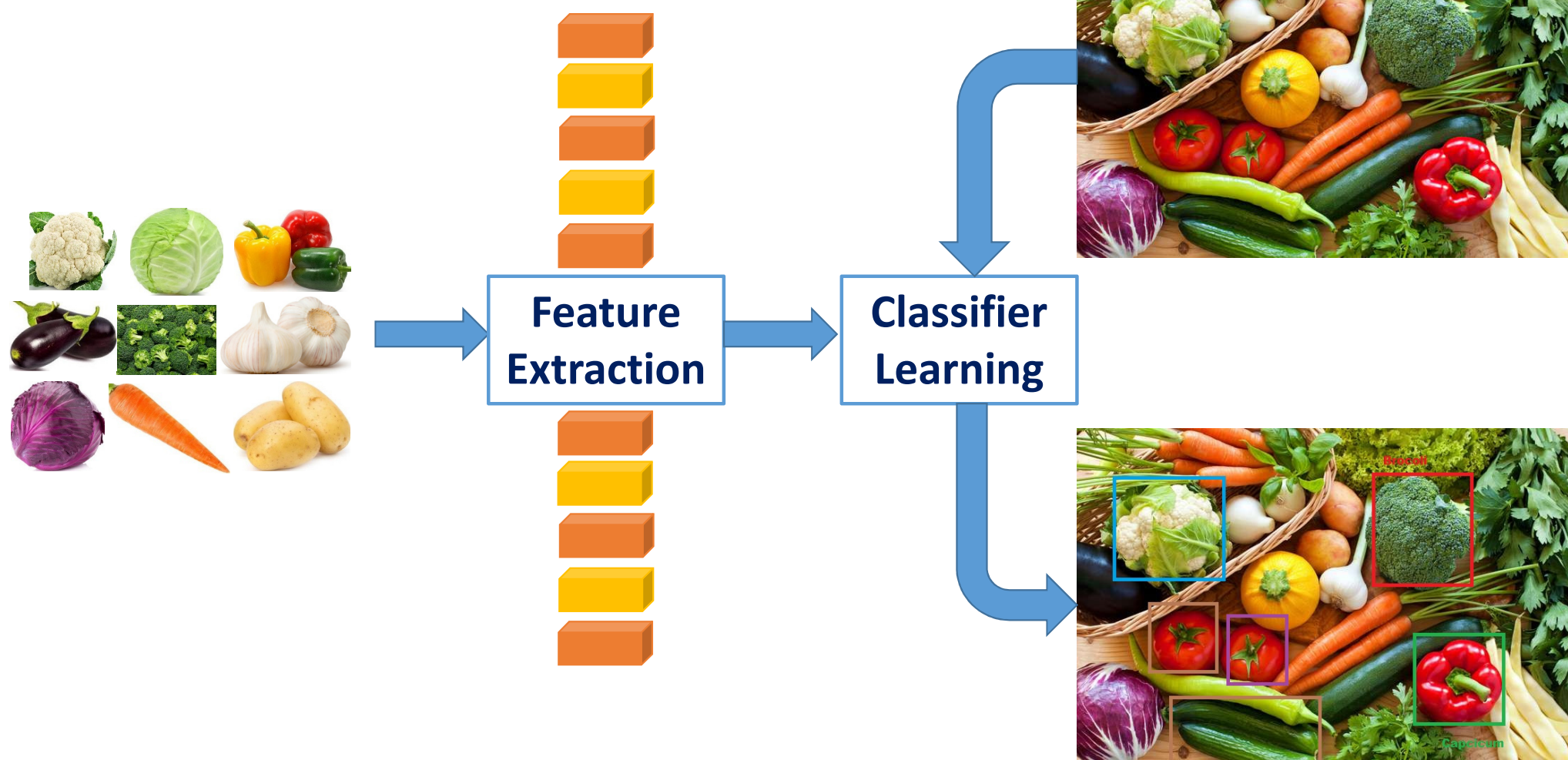
Overview

- Machine Learning Workflow
- Feature Extraction
- Separability & Classification
- The Curse of Dimensionality
- Feature Selection
- Hashing Techniques

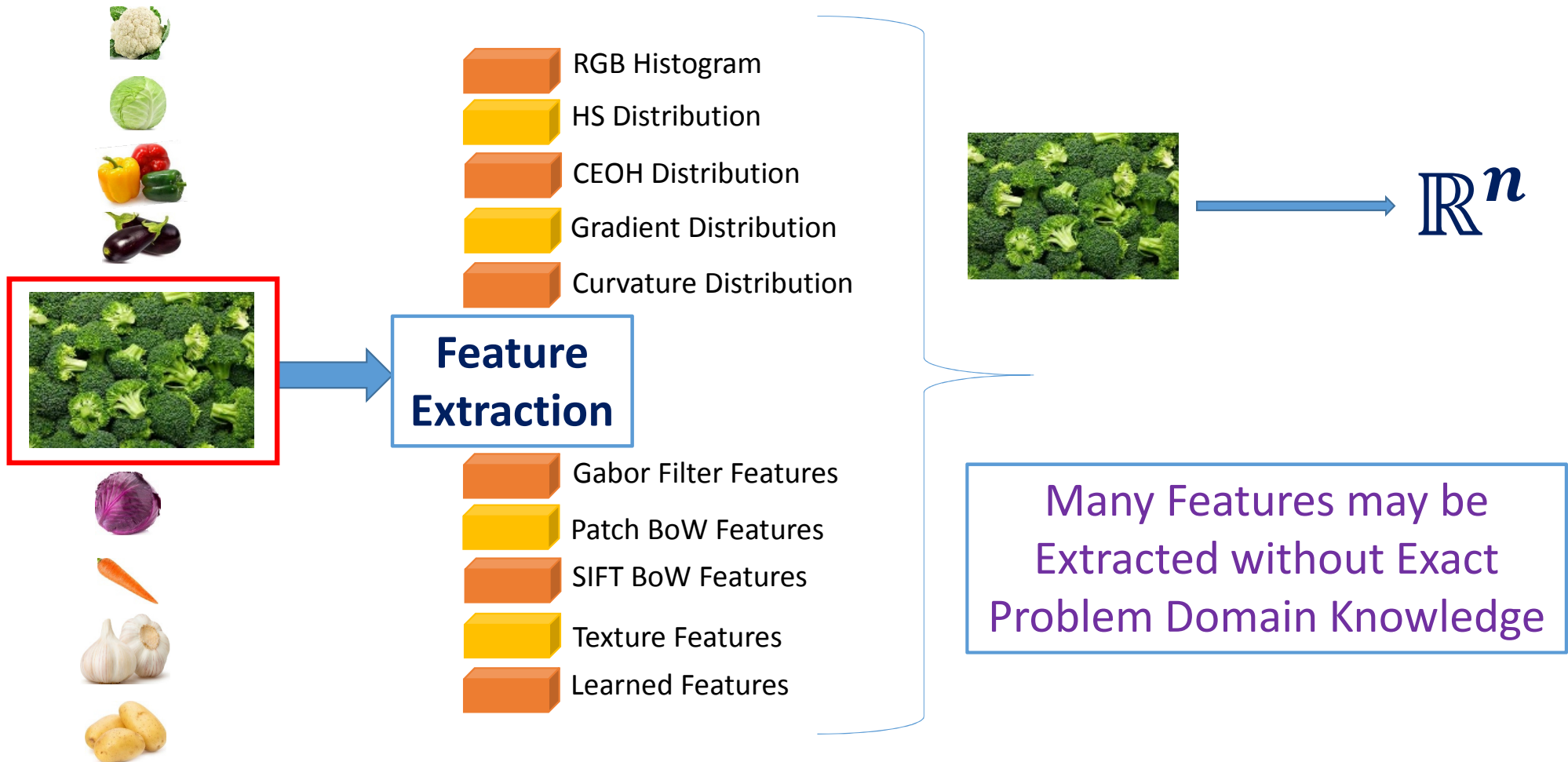
Machine Learning Task Workflow



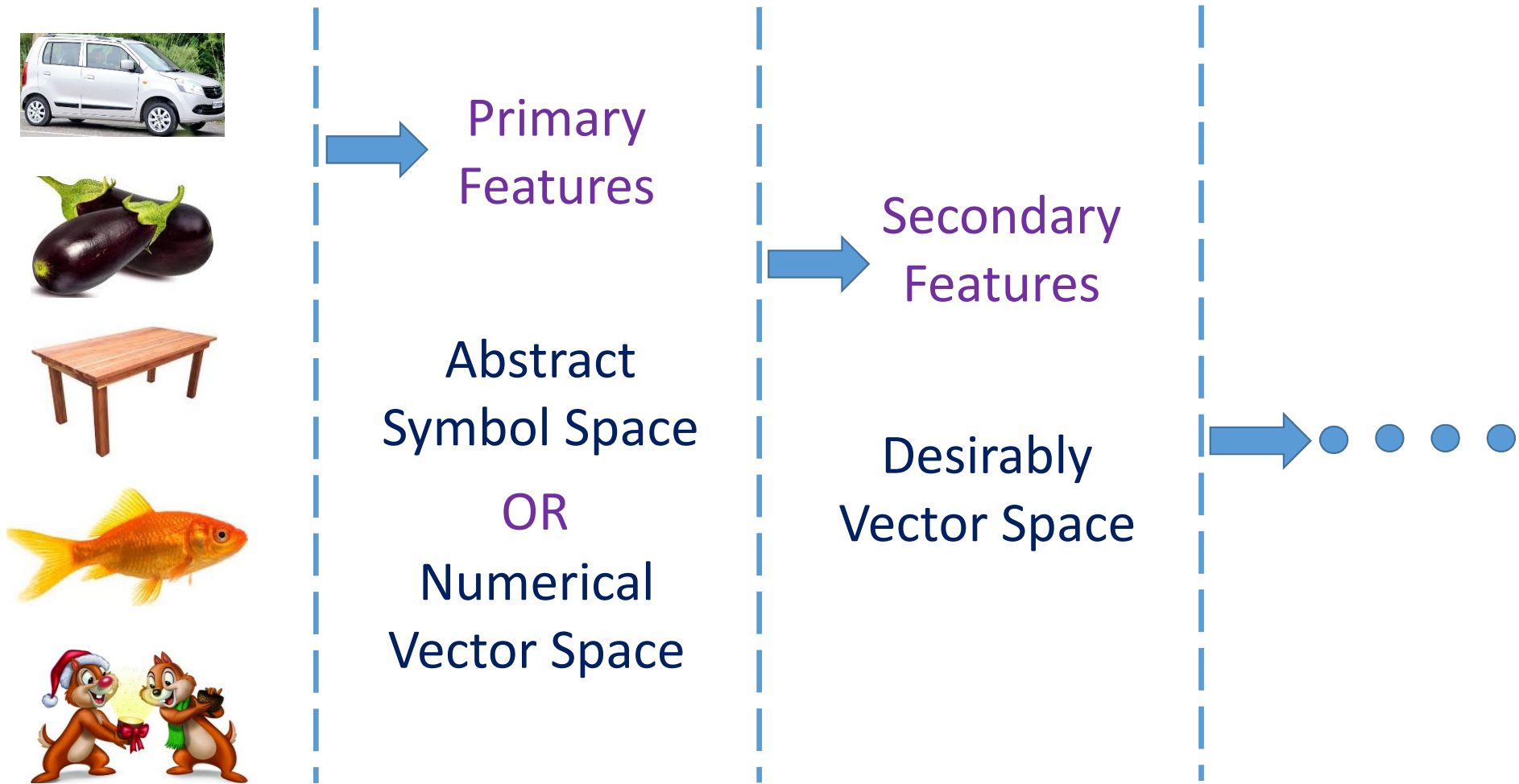
A Classification Task



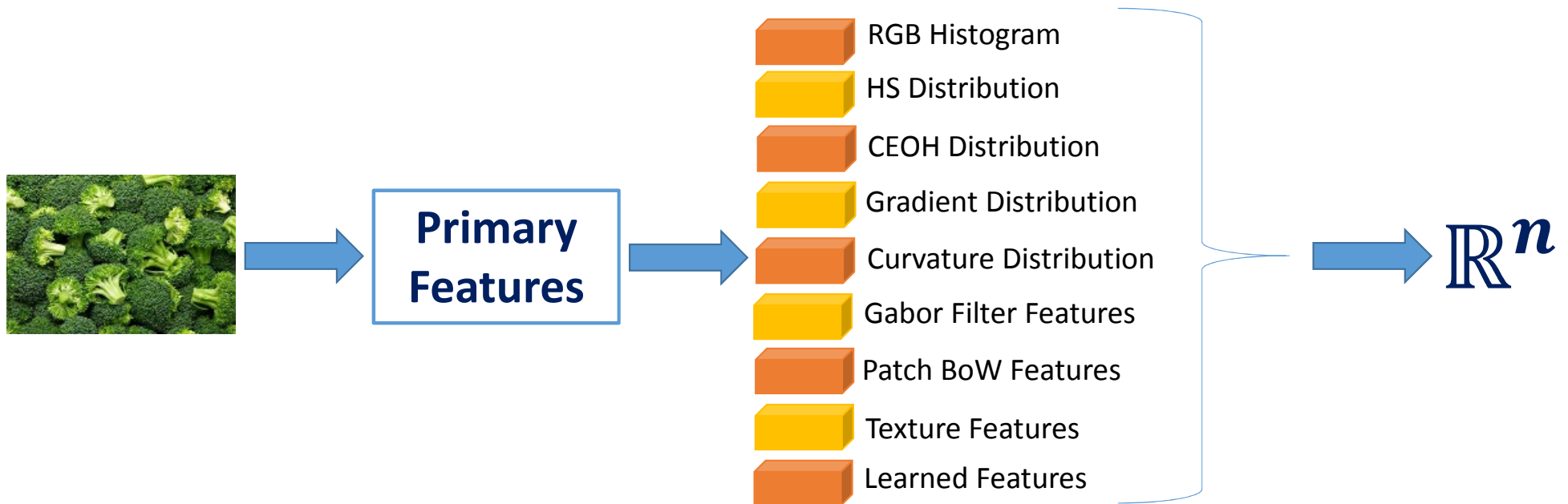
From Real World to Symbol Space



Feature Extraction



Primary Features: Measurement



Secondary Features Onwards: Transformation

$$\begin{pmatrix} x_1 \\ \vdots \\ x_i \\ \vdots \\ x_n \end{pmatrix} \xrightarrow{\quad} \begin{bmatrix} a_{11} & \cdots & a_{nm} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{bmatrix} \xrightarrow{\quad} \begin{pmatrix} y_1 \\ \vdots \\ y_j \\ \vdots \\ y_m \end{pmatrix}$$

$X \in \mathbb{R}^n$ $A \in \mathbb{R}^{m \times n}$ $Y \in \mathbb{R}^m$

Linear Transformation: $Y = AX$

Secondary Features Onwards: Transformation

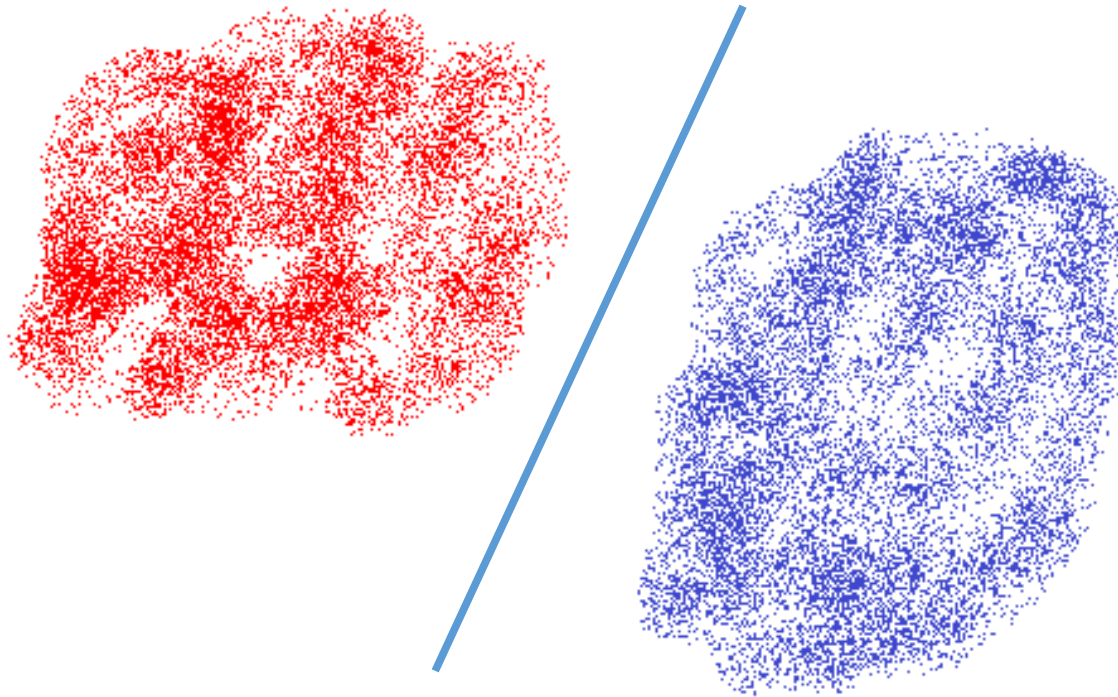
$$\begin{pmatrix} x_1 \\ \vdots \\ x_i \\ \vdots \\ x_n \end{pmatrix} \xrightarrow{\quad} \begin{bmatrix} a_{11} & \cdots & a_{nm} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{bmatrix} \xrightarrow{\quad} \begin{pmatrix} y_1 \\ \vdots \\ y_j \\ \vdots \\ y_m \end{pmatrix} \xrightarrow{\quad} \begin{pmatrix} z_1 = g_1(y_1) \\ \vdots \\ z_j = g_j(y_j) \\ \vdots \\ z_m = g_m(y_m) \end{pmatrix}$$

$X \in \mathbb{R}^n$ $A \in \mathbb{R}^{m \times n}$ $Y \in \mathbb{R}^m$ $Z \in \mathbb{R}^m$

Nonlinear Functions: $g_1, \dots, g_j, \dots, g_m$

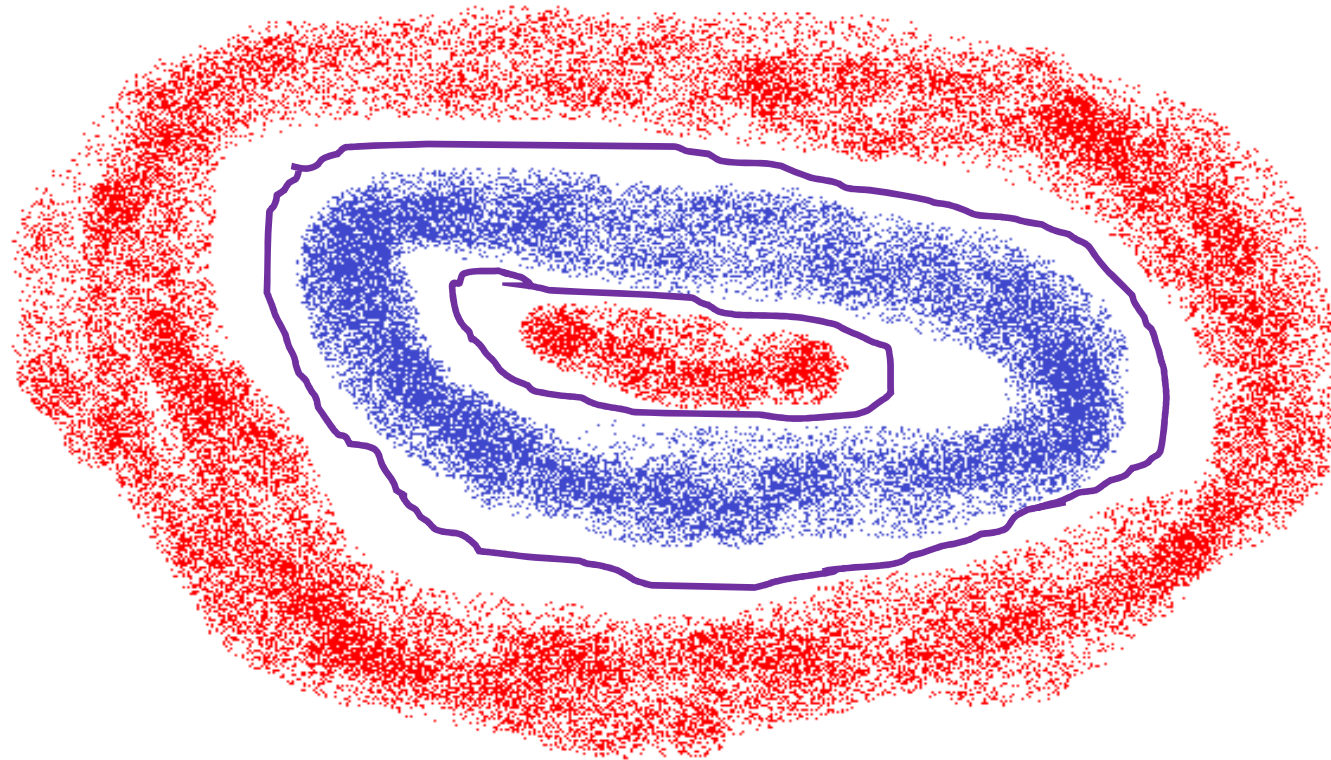
Non-Linear Transformation

The Feature Space: Linearly Separable



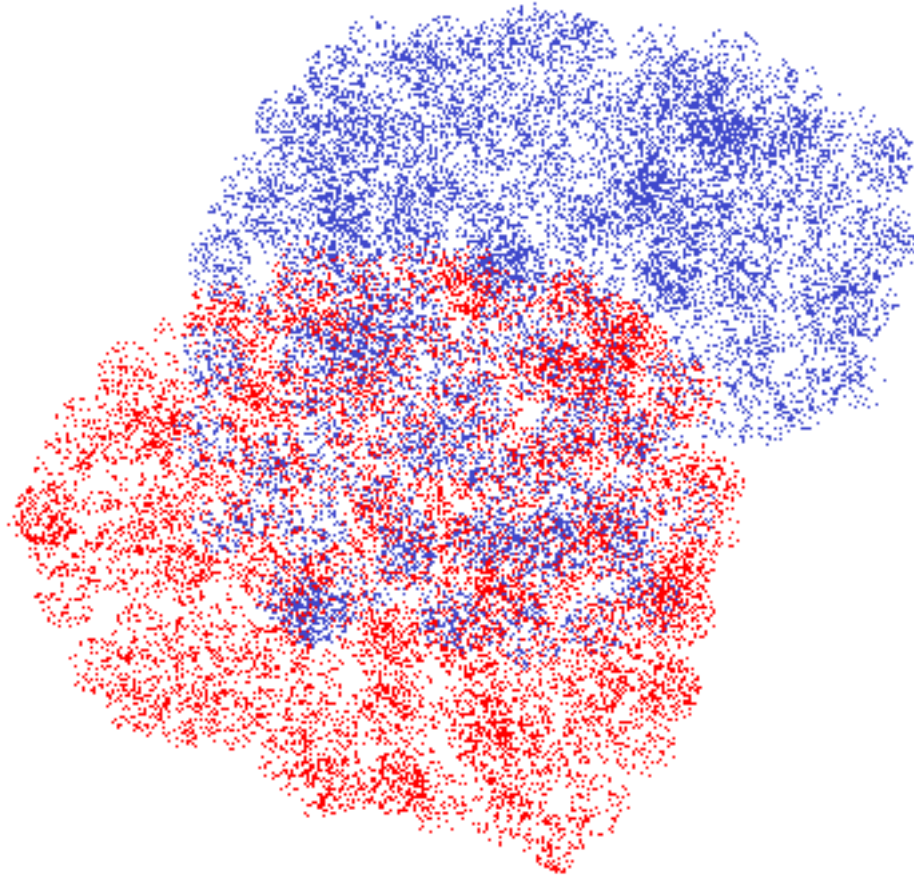
Hyperplane(s) can Separate Features Extracted from
Two (Or More) Classes

The Feature Space: Not Linearly Separable



Nonlinear Hypersurfaces Separate Features Extracted from Two
(Or More) Classes

The Feature Space: Nonseparable Case



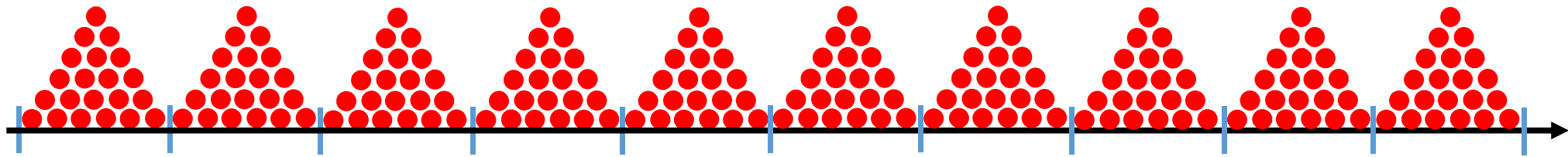
Require Further Transformations for Enhancing Separability

Cover's Theorem of Separability

Non-separable Data can be Transformed to a Higher Dimensional Space through Non-linear Transformation. The Probability of Separability Increases with the Dimension of the Destination Space.

Data is Sparse in Higher Dimensions

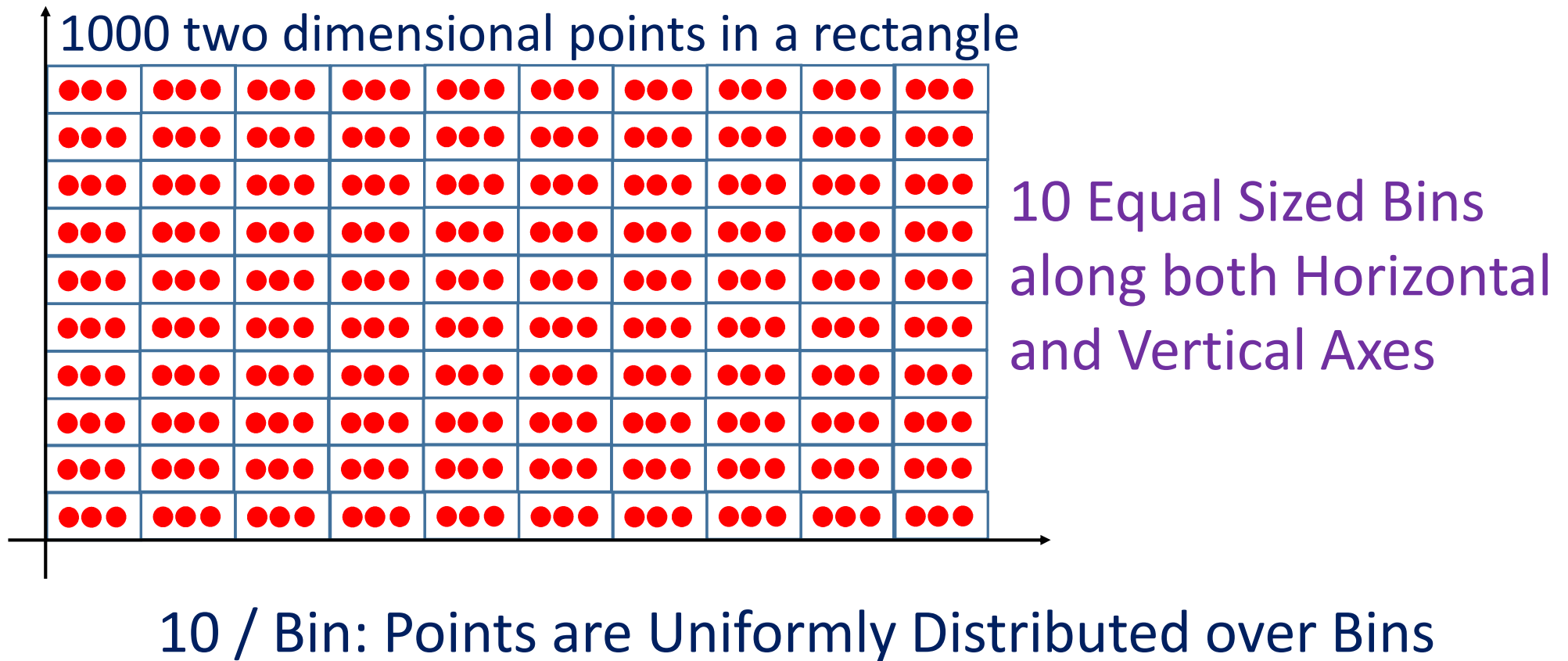
1000 one dimensional points on a line



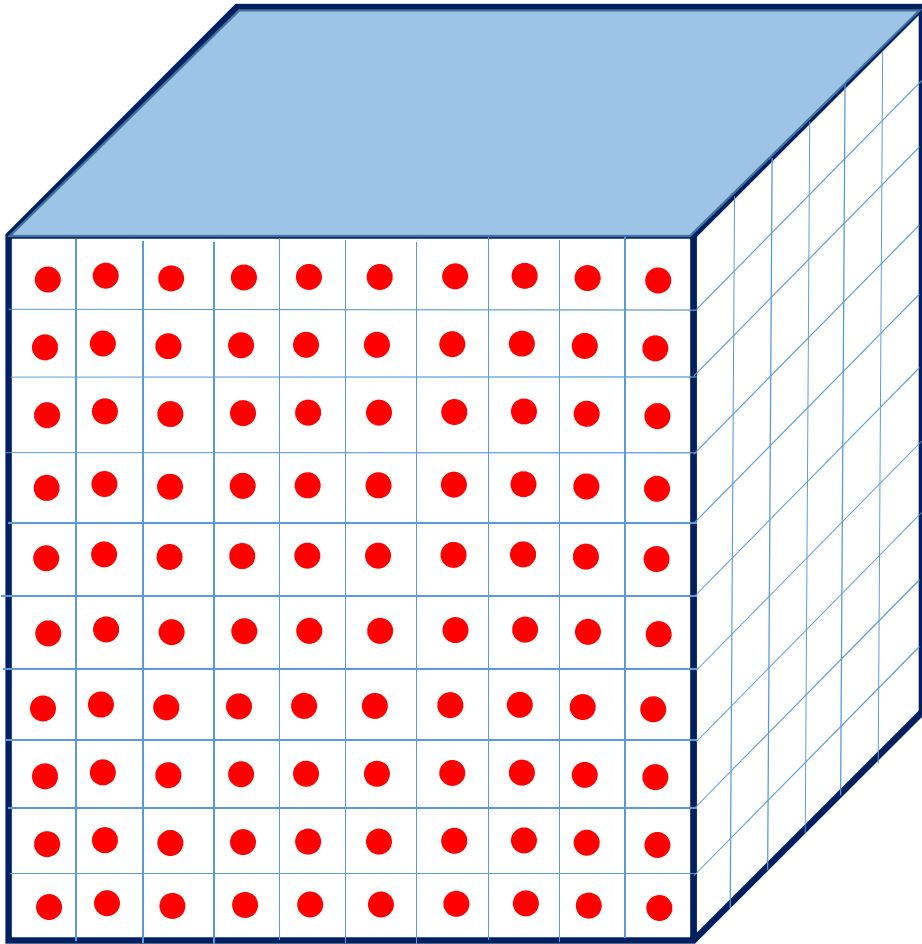
10 Equal Sized Bins on the Line

100 / Bin: Points are Uniformly Distributed over Bins

Data is Sparse in Higher Dimensions



Data is Sparse in Higher Dimensions



1000 Three Dimensional
points in a Cube

10 Equal Sized Bins
along each Axes

1 / Bin: Points are Uniformly
Distributed over Bins

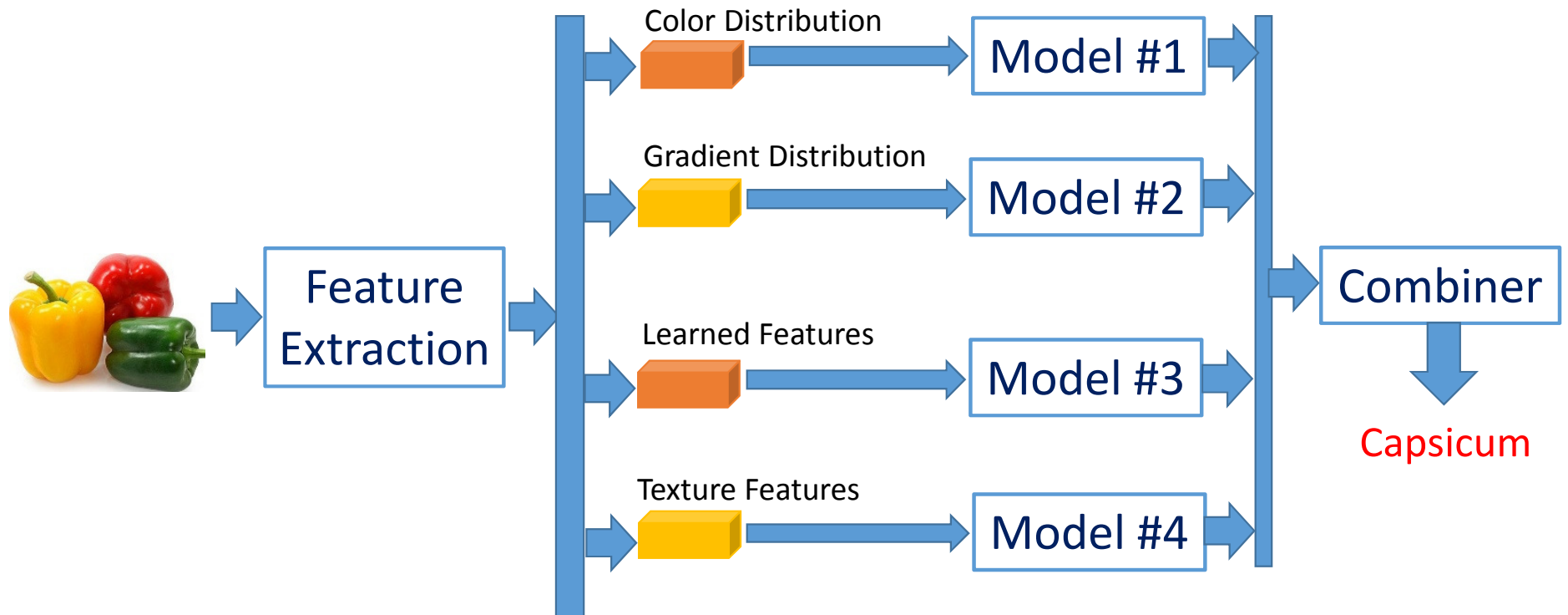
Data is Sparse in Higher Dimensions

- Consider 1000 n-dimensional ($n > 3$) points
- We have 10 Equal Size Bins on each Axis
- There are 10^n (> 1000) bins
- Vacant bins even with 1 point/bin
- This can be generalized...

Curse of Dimensionality

- High Dimension due to Large Number of Features
- High Dimension Leads to More Storage & Computation
 - More Parameters in (Non)Linear Transformations
 - Complex Classification & Regression Models
 - More Operations (e.g. Matrix Inverse in Mahalanobis Distance)
 - Unnecessary Features often lead to Poor Performance
- Working on Feature Subspaces
 - Late Fusion in Ensemble Framework
 - Dimensionality Reduction

Late Fusion in Ensemble Framework



Dimensionality Reduction

- Feature Subset Selection
- Hashing Techniques
- Principal Component Analysis (PCA)
- Linear Discriminant Analysis (LDA)
- Exploratory Factor Analysis (EFA)

Feature Subset Selection

- Feature Selection vs. Extraction/Transformation(s)
- Supervised vs. Unsupervised Approaches
- Feature Subset Selection – 2^n possibilities for $x \in \mathbb{R}^n$
- Forward & Backward Feature Subset Selection
- Useful Tool for Feature Combination Analysis
- Generally used for Handcrafted Features

Feature Subset Selection: Notations

X_j : The j^{th} Feature Subset; $j = 1, \dots, m$

S_i : Feature Set formed in i^{th} Iteration

$$S_i = \{X^{(1)} \cup X^{(2)} \dots \cup X^{(i)}\}$$

C_i : Classifier Trained with Feature Set S_i

$E(S_i; C_i)$: Evaluation Error of Classifier C_i using Feature Set S_i

ϵ : A Threshold on Evaluation Error

Forward Feature Subset Selection

1. $S_0 = \varphi; \text{SizeOf}(\cup_{i=1}^m X_i) = N; E(S_0; C_0) = \text{LARGE_NUMBER}$
2. **WHILE** $\text{SizeOf}(S_{i-1}) < N$ **DO**
3. $J_i = \{k: X_k \notin S_{i-1}\}$
4. $\forall k$ Train $C_{i-1}^{(k)}$ with $S_{i-1} \cup X_k$
5. $j = \text{argmin}_k E(S_{i-1} \cup X_k; C_{i-1}^{(k)}); k \in J_i$
6. **IF** $E(S_{i-1} \cup X_j; C_{i-1}^{(j)}) < E(S_{i-1}; C_{i-1}) - \epsilon$
7. **THEN** $S_i = S_{i-1} \cup X_j; C_i = C_{i-1}^{(j)}$
8. **ELSE** break
9. **END WHILE**

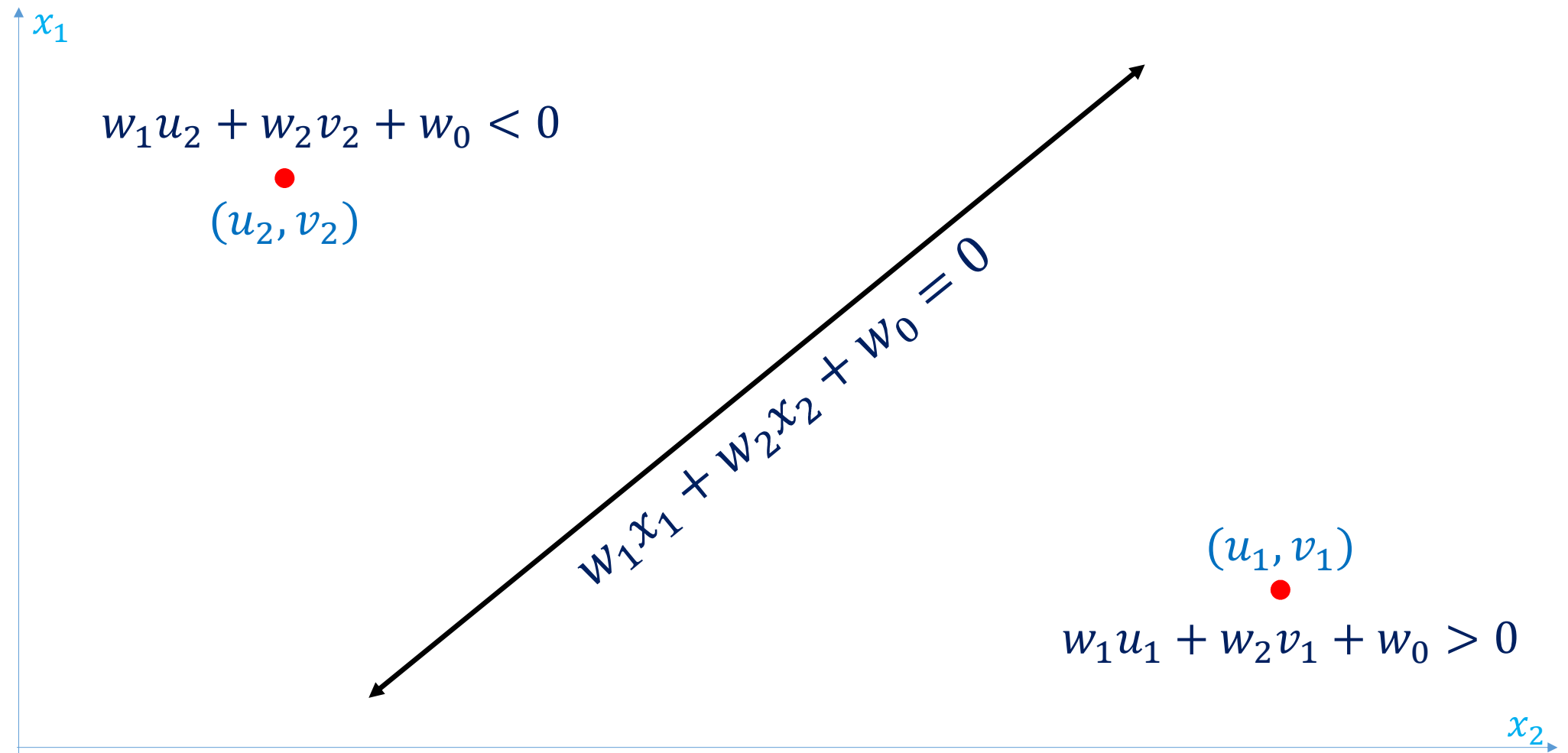
Backward Feature Subset Selection

1. $S_0 = \cup_{i=1}^m X_i$
2. **WHILE** $SizeOf(S_{i-1}) > 0$ **DO**
3. $J_i = \{k: X_k \in S_{i-1}\}$
4. $\forall k$ Train $C_{i-1}^{(k)}$ with $S_{i-1} - X_k$
5. $j = argmin_k E(S_{i-1} - X_k; C_{i-1}^{(k)}); k \in J_i$
6. **IF** $E(S_{i-1} - X_j; C_{i-1}^{(j)}) < E(S_{i-1}; C_{i-1}) - \epsilon$
7. **THEN** $S_i = S_{i-1} - X_j; C_i = C_{i-1}^{(j)}$
8. **ELSE** break
9. **END WHILE**

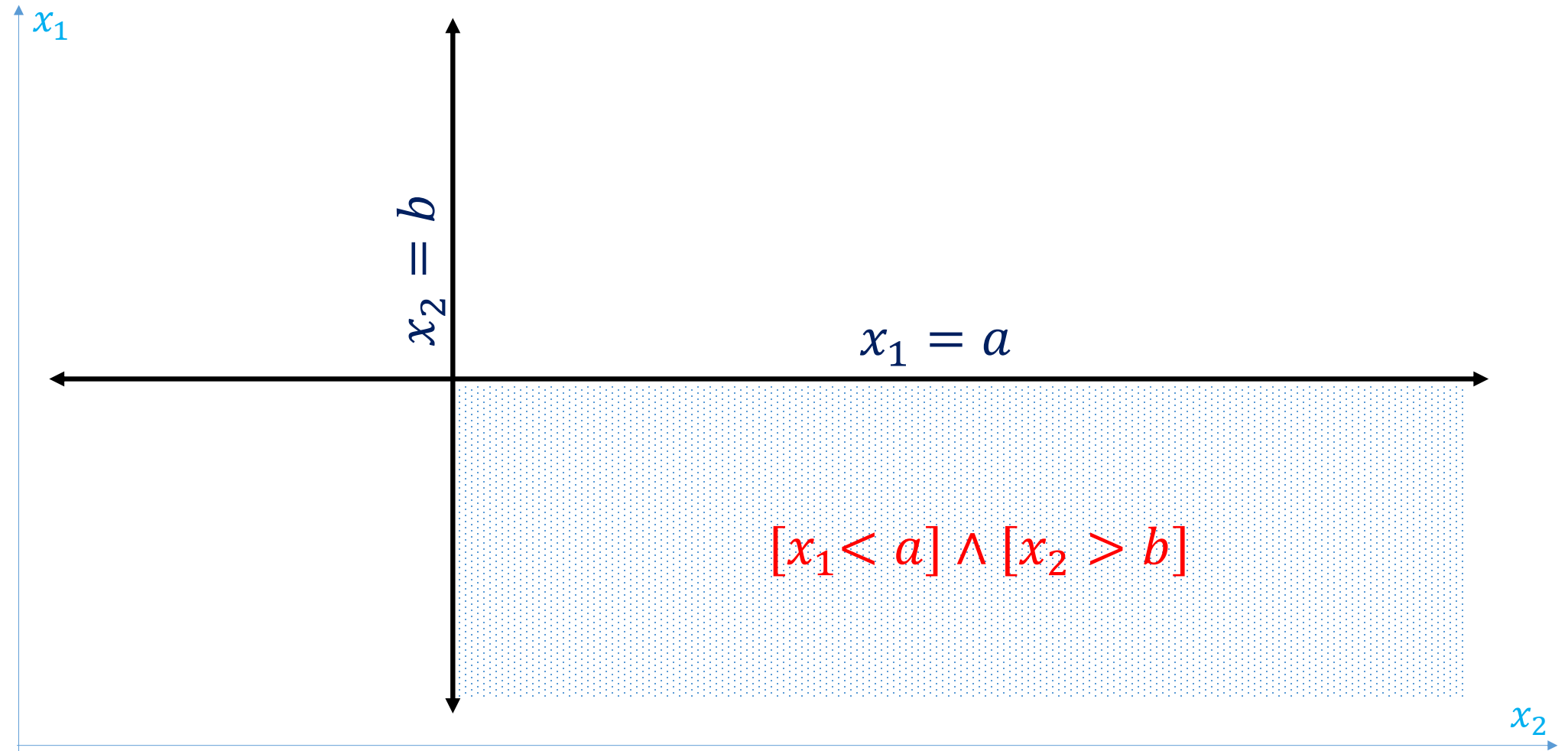
Hashing Techniques

- Dividing Feature Space into Partitions
- Assign Integer IDs (Hashes) to each Partition
- Mechanisms for Partitioning Feature Space
 - Locality Sensitive Hashing
 - Hashing with Hierarchical Structures
- Features in a Partition has the Same Hash
- Features of Similar Entities are Close in Feature Space
- Proximal Points have High Probability of Similar Hash

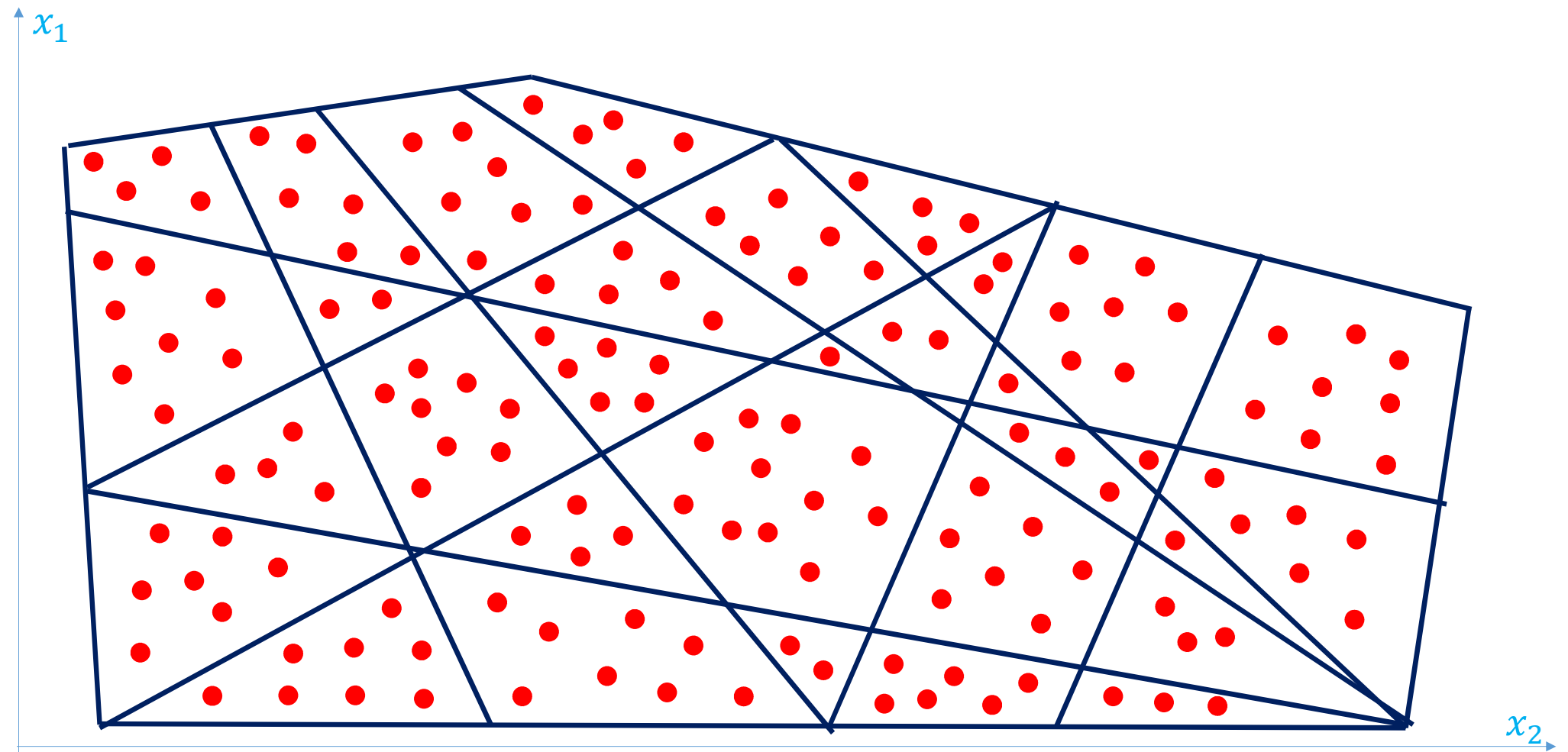
Feature Space Partitioning: Oblique



Feature Space Partitioning: Axis Aligned

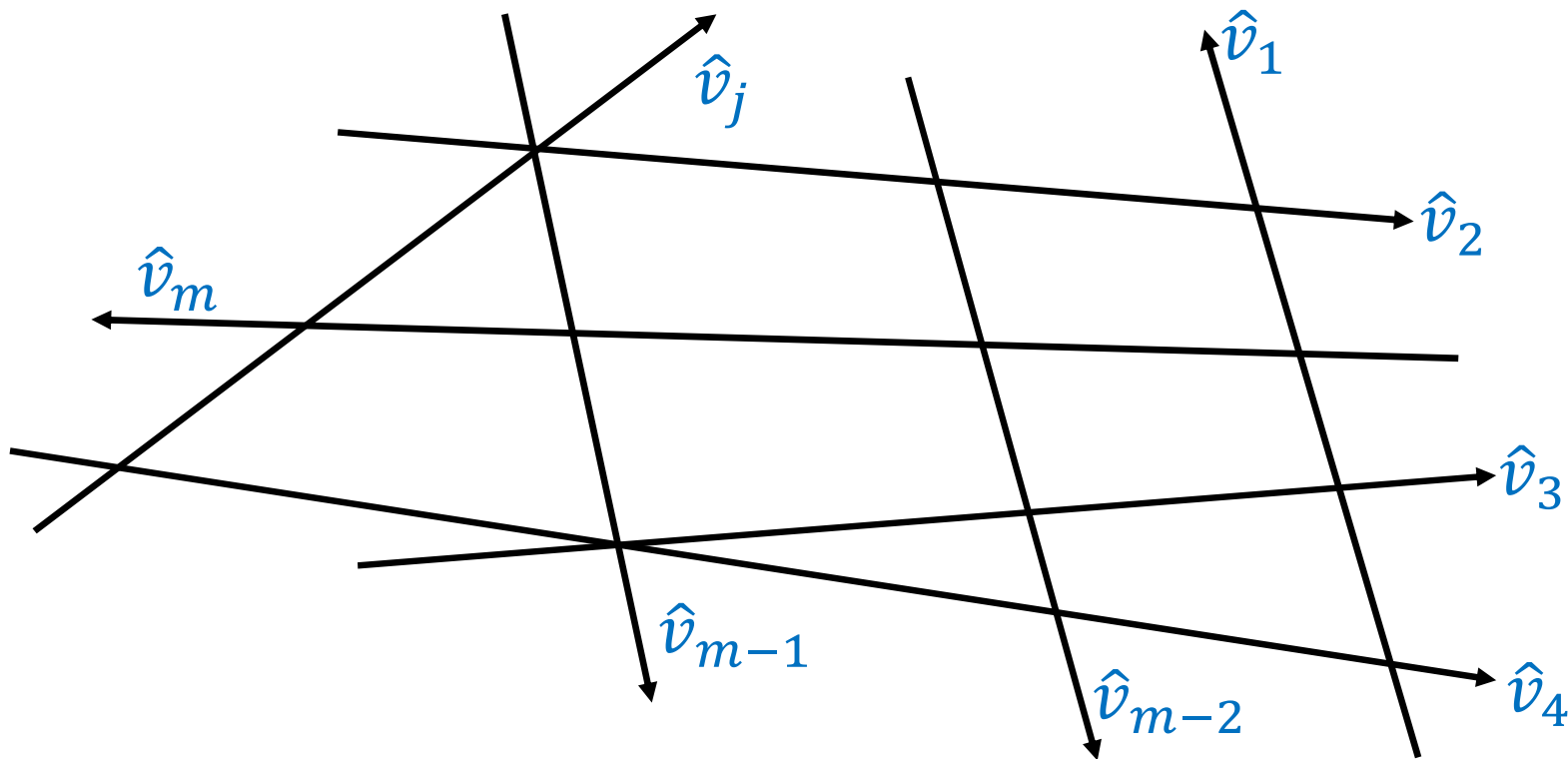


Feature Space Partitioning

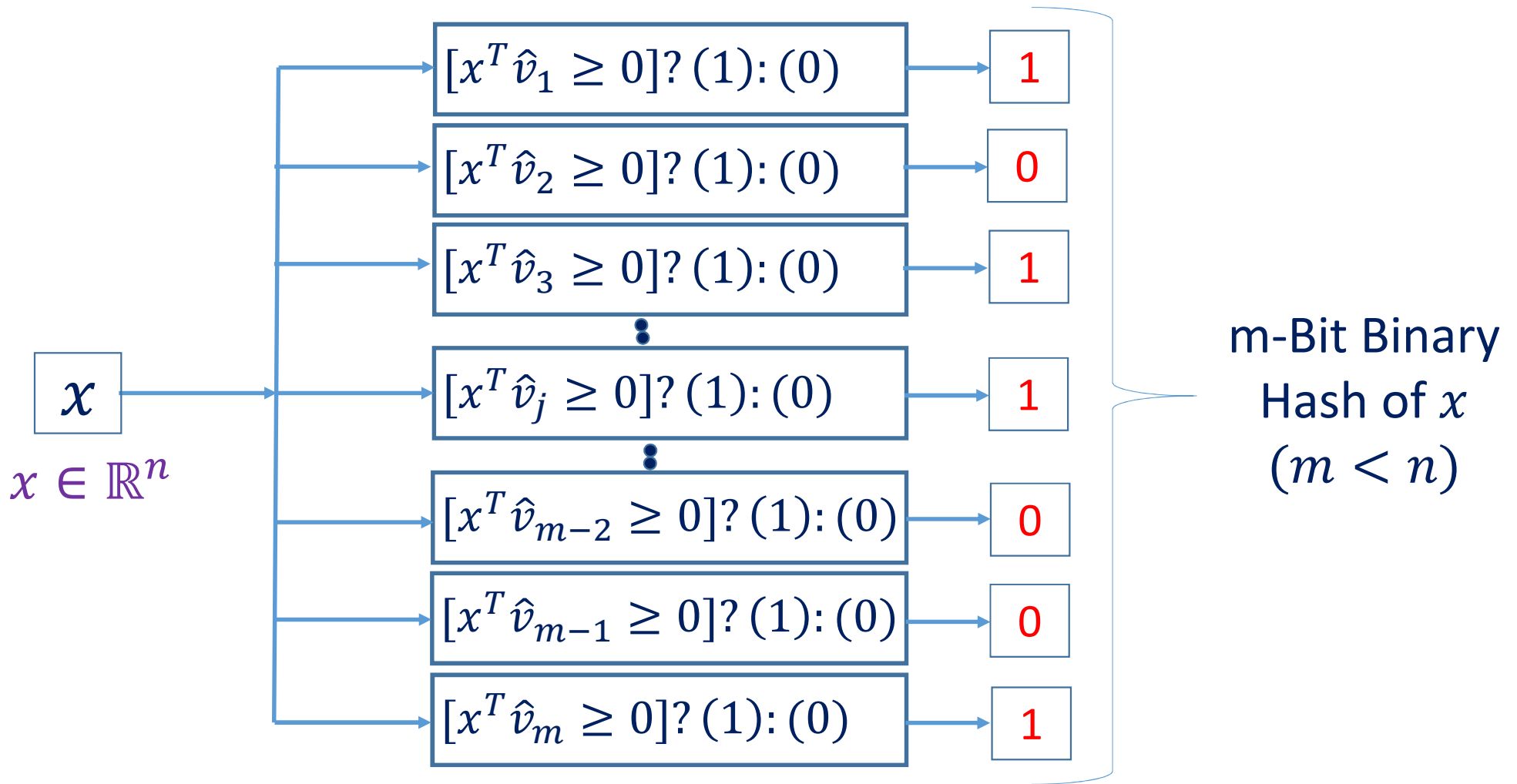


Locality Sensitive Hashing – Random Projections

Generate Random Unit Vectors in Feature Space: $\hat{v}_1, \dots \hat{v}_i \dots \hat{v}_m$

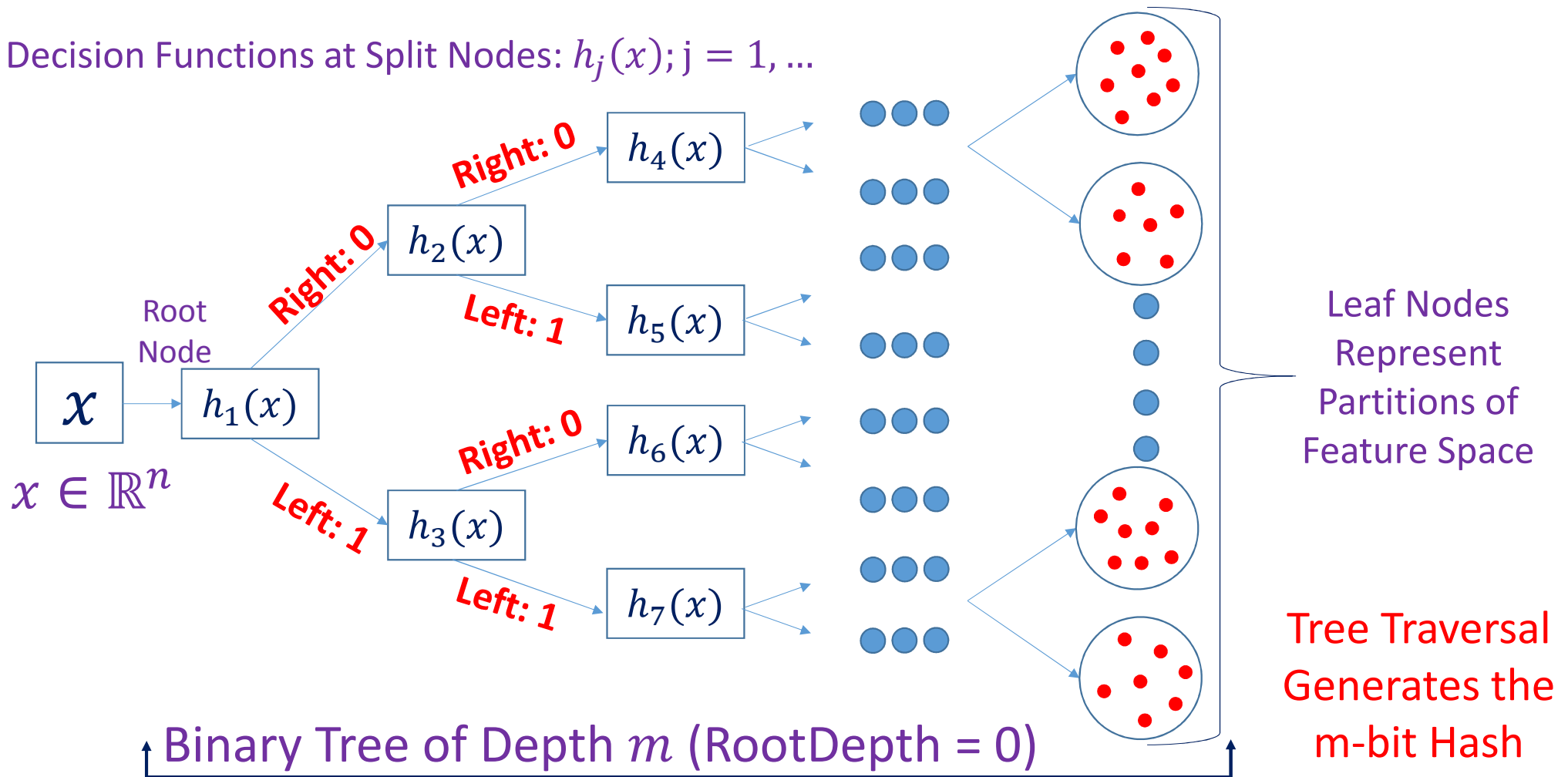


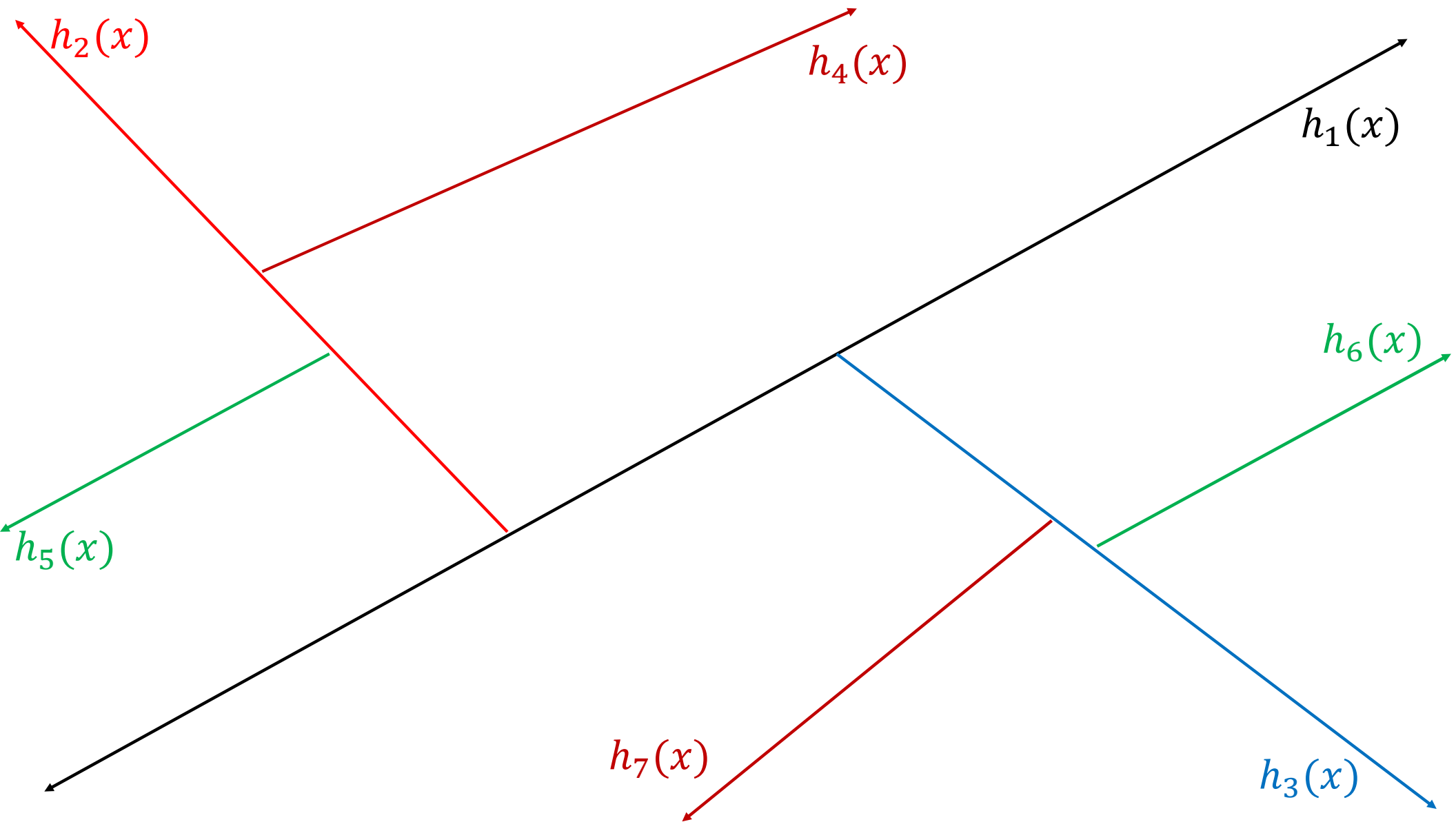
Locality Sensitive Hashing – Random Projections



Hashing with Hierarchical Structures

Decision Functions at Split Nodes: $h_j(x); j = 1, \dots$





Applications of Hashing Techniques

- Near Duplicate Detection (in Archives)
- Hierarchical Grouping
- Search in Image Databases
- Audio Similarity Identification
- Digital Audio-Video Fingerprinting

Summary

- Significance of Feature Engineering
- Feature Extraction & Transformations
- Separability in Feature Space
- Issues with Dimensionality
- Dimensionality Reduction Techniques
- Forward & Backward Feature Selection
- Binary Hashing Techniques



Thank You