

Regression Trees: Split Optimization Task

$$\min_{d, s} \left\{ \min_{c_1} \sum_{\vec{x} \in R_1(d, s)} (\hat{p}(\vec{x}) - c_1)^2 + \min_{c_2} \sum_{\vec{x} \in R_2} (\hat{p}(\vec{x}) - c_2)^2 \right\}$$

Trees are prone to overfitting - limit depth with regularization:

$$J_{\text{reg}} = \sum_{j=1}^k \int_{\vec{x} \in R_j} (p(\vec{x}) - \hat{p}(\vec{x}))^2 d\vec{x} + \alpha \cdot k$$

k : number of splits

α : weighing parameter

Curse of dimensionality

- statistically all distances become large
- similarity measure (distance) breaks

Manifold Learning

Goal: reduce dimensionality

PCA: Project onto hyperplane with largest variance

MDS: "PCA on distances / dissimilarities"

PCA Concept

C : $p \times p$ covariance matrix

V : matrix of column eigenvectors

L : diagonal matrix with eigenvalue entries

$$C = \frac{X^T X}{n-1} \Rightarrow C = V L V^T \quad (\text{Diagonalized})$$

\Rightarrow Perform SVD of X : $X = U S V^T$ where S is a diagonal matrix

$$\lambda_i = \frac{s_i^2}{n-1}$$

MDS Concept

$$S = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}^d$$

$$X = [x_1, \dots, x_N] \in \mathbb{R}^{d \times N}$$

$$D^2 = [d_{ij}^2]_{i,j \in [1, \dots, N]}, \quad d_{ij}^2 = (x_i - x_j)^T (x_i - x_j)$$

Given: D Compute: X

$$d_{ij}^2 = (x_i - x_j)^T (x_i - x_j) = x_i^T x_i + x_j^T x_j + 2x_i^T x_j$$

$$\Rightarrow D^2 = \text{diag}(X^T X) \cdot \vec{1}^T + \vec{1} \text{diag}(X^T X)^T - 2X^T X$$

C: Centering Matrix

$$C = \left(I - \frac{1}{N} \vec{1} \vec{1}^T \right)$$

$$-\frac{1}{2} C D^2 C = \dots = X^T X \quad (\text{centered distances})$$

$$\Rightarrow \begin{array}{l} X^T X \quad [\text{eigendecomposition}] = U \Sigma U^T \\ \Rightarrow X = \Sigma^{1/2} U^T \end{array} \quad \text{done}$$

1. Determine the m largest eigenvalues & eigenvectors

2. Drop the remaining eigenvectors

 \rightarrow gives us feature vectors (rows) in new, reduced spaceProportion of variance explained by p dimensions:

$$\frac{\sum_{i=1}^p \lambda_i}{\sum_{i=1}^{n-1} \lambda_i}$$

ISOMAP Algorithm

- non-linearity patch for MDS
- close distance: Euclidean distance

others: distance via graph edges (shortest path)

[sparse adjacency matrix]
"geodesic distance"

Locally Linear Embedding

Idea: treat local neighborhood of a point as linear

Algorithm:

- 1.) Define the neighborhood (KNN / distance thresholding)
- 2.) Solve w_{ij} in the high dimensional space

$$\min_i \sum_i \|x_i - \sum_{j \in N(x_i)} w_{ij} x_j\|_2^2 \quad \text{s.t.} \quad \sum_{j \in N(x_i)} w_{ij} = 1$$

[linear constraint]

- 3.) Solve for $x'_i \in \mathbb{R}^{d'}$ ($d' \ll d$)

$$\sum_i \|x'_i - \sum_{j \in N(x_i)} w_{ij} x'_j\|_2^2 \quad \text{s.t.} \quad \frac{1}{N} \sum_i x'_i \cdot x'_i{}^T = I$$

[identity covariance]

$$\sum_i x'_i = \vec{0}$$

[zero mean]

MODIFICATION of step 2

$$2^*) \quad \sum_i \|x_i - \sum_{j \in N(x_i)} w_{ij} x_j\|_2^2 = \sum_i \|(x_i - t) - \sum_{j \in N(x_i)} w_{ij} (x_j - t)\|_2^2$$

translation

$$\text{set } x_i = t \Rightarrow \sum_{i=1}^N \left\| \sum_j w_{ij} (x_j - x_i) \right\|_2^2 = \|M_i \vec{w}_i\|_2^2$$

$$\Rightarrow \text{minimize } (M_i \vec{w}_i)^T (M_i \vec{w}_i) + \lambda (1 - \vec{w}_i^T \vec{w}_i) \quad \dots \Rightarrow M_i^T M_i \vec{w}_i = \lambda \vec{w}_i$$

eigenvector / -value \rightarrow solvable

Laplacian Eigenmaps

- 1.) Build adjacency graph
- 2.) Compute affinities between nodes
- 3.) Perform Eigendecomposition
- 4.) Low-Dim Embedding

} Construct Similarity Matrix

$$L = D - W$$

$$d_{ij} = \begin{cases} \sum_{k=1}^N w_{ik} & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$$

$$w_{ij} = \begin{cases} \text{weight between node } i \text{ and } j, & i \neq j \\ 0 & \text{otherwise} \end{cases}$$

Heat Kernel

$$w_{ij} = e^{-\|x_i - x_j\|_2^2}$$

Binary affinity

$$w_{ij} = \begin{cases} 1 & \text{if } \|x_i - x_j\| \leq \tau \\ 0 & \text{otherwise} \end{cases}$$

Objective Function

$$\min \sum_{i=1}^N \sum_{j=1}^N \|x'_i - x'_j\|_2^2 w_{ij}$$

$$= 2 \vec{x}'^T \underbrace{(D - W)}_{\text{Graph Laplacian}} \vec{x}'$$

$$\Rightarrow \min \vec{x}'^T L \vec{x}' \quad \text{s.t.} \quad \vec{x}'^T D \vec{x}' = 1$$

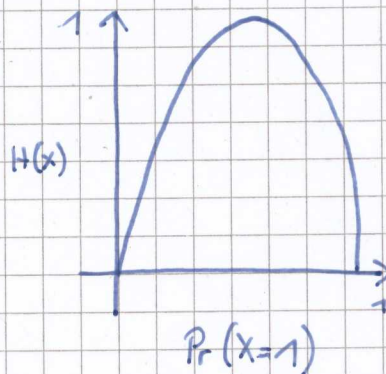
$$\Rightarrow D^{-1} L \vec{x}' = \lambda \vec{x}'$$

Decision TreesInformation Gain (for splitting at S_j)

$$I = H(S_j) - \sum_{i \in \{L, R\}} \frac{|S_j^i|}{|S_j|} \cdot H(S_j^i)$$

Entropy

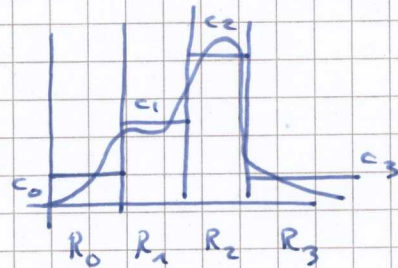
$$H(S_j) = - \sum_{c \in C} p(c) \cdot \log(p(c))$$



Binary Entropy Function

Regression Forests

Goal: Predict continuous label

Regression Tree: Estimation of c_j for each R_j 

Split optimization task:

$$\min_{d, s} \left\{ \min_{c_1} \sum_{\vec{x} \in R_1(d, s)} (\hat{p}(\vec{x}) - c_1)^2 + \min_{c_2} \sum_{\vec{x} \in R_2(d, s)} (\hat{p}(\vec{x}) - c_2)^2 \right\}$$

Prevent overfitting by using a regularizer:

$$J_{\text{reg}} = \sum_{j=1}^K \int_{\vec{x} \in R_j} (p(\vec{x}) - \hat{p}(\vec{x}))^2 d\vec{x} + \alpha \cdot K$$

 K : number of splits α : weighting parameter

Density Forests

- Each leaf node is a multivariate Gaussian

Information Gain - same as Regression

Entropy

$$H(S_j) = \frac{1}{2} \log \left((2\pi e)^d \left| \sum S_j \right| \right)$$

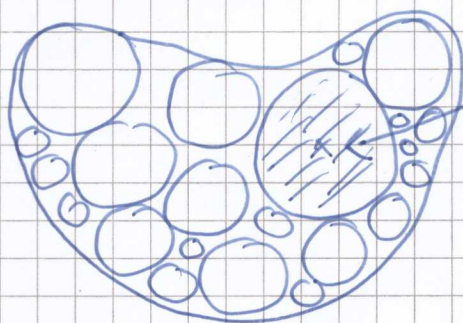
determinant of
the covariance of
the data ("volume of
a cluster")

⇒ Information Gain:

$$I(S_j, \mathcal{Y}) = \log(|\Sigma(S_j)|) - \sum_{i \in \{L, R\}} \frac{|S_j^i|}{|S_j|} \cdot \log(|\Sigma(S_j^i)|)$$

Manifold Forests

1. Partition feature space via Density Forest training
2. Use Regions as local neighborhoods → similarity matrix
3. Apply Laplacian Eigenvectors



Density Tree

partitioning ⇒ sparse similarity matrix

Affinity Model(s)

$$w_{ij} = e^{-Q(x_i, x_j)} \Leftarrow \text{Affinity}$$

$$Q(x_i, x_j) \Leftarrow \text{distance function}$$

Mahalanobis

$$Q(x_i, x_j) = \begin{cases} d_{ij}^T (\Lambda_{\mathcal{L}(x_i)})^{-1} d_{ij} & \text{if in same leaf } \mathcal{L}(x_i) \\ \infty & \text{otherwise} \end{cases}$$

Binary

$$Q(x_i, x_j) = \begin{cases} 0 & \text{if in the same leaf} \\ \infty & \text{otherwise} \end{cases}$$