

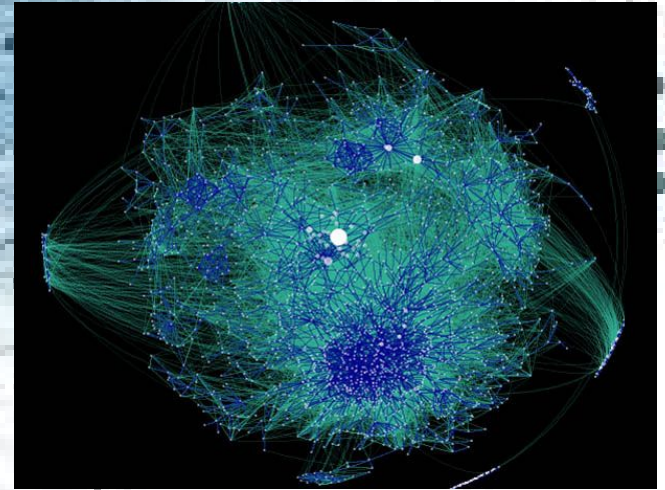


Heidi Matrix Visualization

Data Visualization. Why?

If data is high dimensional, then identifying patterns is difficult task owing to complexity and high dimensionality issues.

As the dimension of dataset increased, data in subspaces are more prominent than original dimensional space.



Heidi Matrix

Heidi matrix gives insight of

1. How the clusters are placed with respect to each other
2. Characteristics of placement of points within a cluster in all the subspaces.
3. Characteristics of overlapping clusters in various subspaces

Heidi matrix Input/output

INPUT:

- d dimensional dataset

OUTPUT:

- 2-D colored matrix called Heidi matrix



Heidi Algorithm:

Given $n \times d$ dataset . (n : no. of rows ; d : no. of columns)

1. compute set of all possible subspaces. ($2^n - 1$)
2. For each subspace
compute knn for every pair of points
3. Merge these ($2^n - 1$) matrices into one (bitwise)
4. Grouping and ordering of points
5. Draw image



Implementation

Heidi approach was implemented in python and the pseudo code is as follows.

1. compute set of all possible subspaces. ($2^n - 1$)

This step was done via binary counting.

2. pick one subspace

compute knn for every pair of points

For computing knn sklearn inbuilt method NearestNeighbors was used,

3. merge these ($2^n - 1$) matrices into one

This step was done using basic binary manipulation.

4. Represent above computed 2-D matrix as image

— Heidi approach was implemented on

Haberman's Survival Dataset.

Number of Instances: 306

Number of Attributes: 4

Attribute Information:

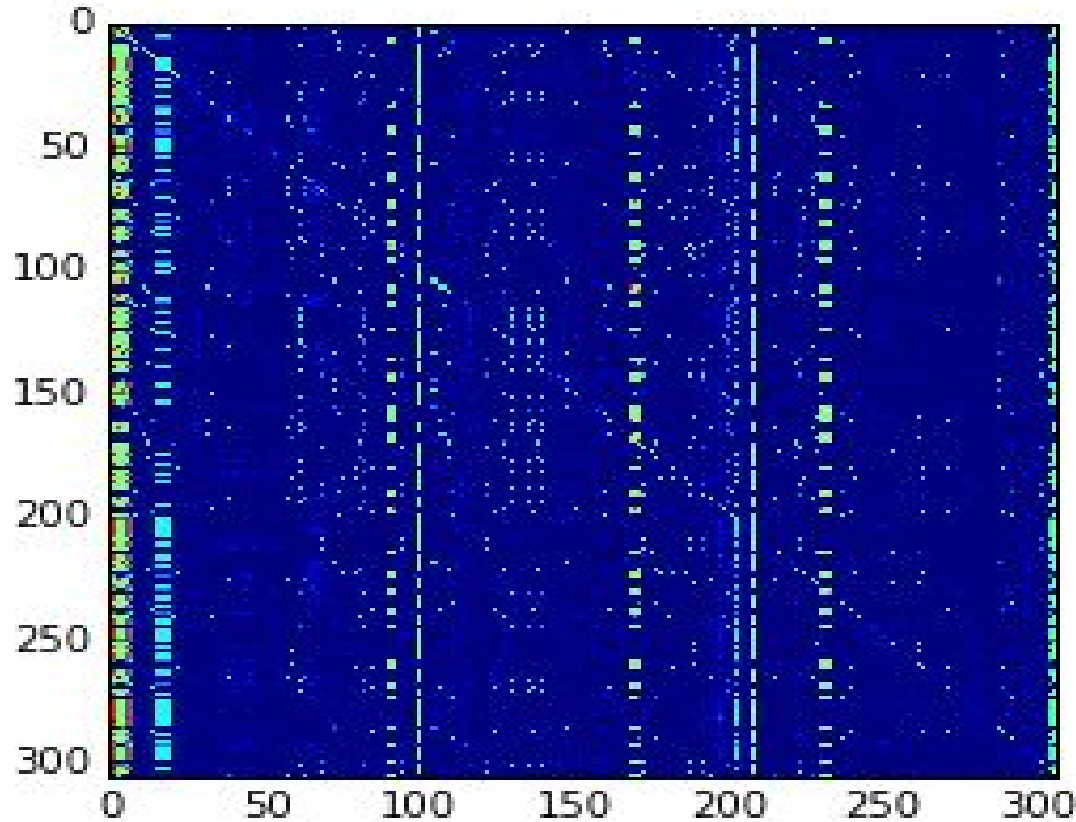
1. Age of patient at time of operation (numerical)
2. Patient's year of operation (year - 1900, numerical)
3. Number of positive axillary nodes detected (numerical)
4. Survival status (class attribute)



Heidi Matrix

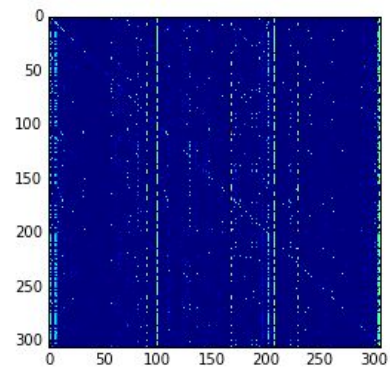
305 X 305 Heidi Matrix

(305 : no of rows/
instances in test
dataset,.

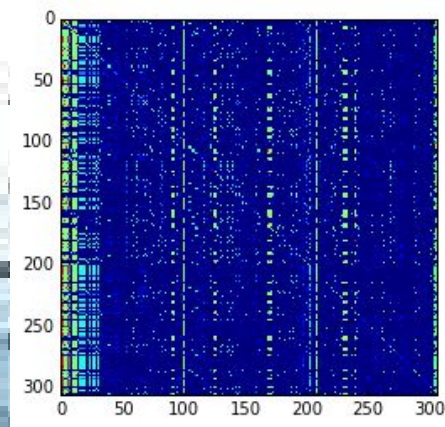


Visual representation of data via Heidi Matrix

This image is generated
output when heidi
matrix algorithm was
implemented in Python



Heidi Matrix with
varying value of k
 $k=[5 \ 10 \ 15 \ 20]$

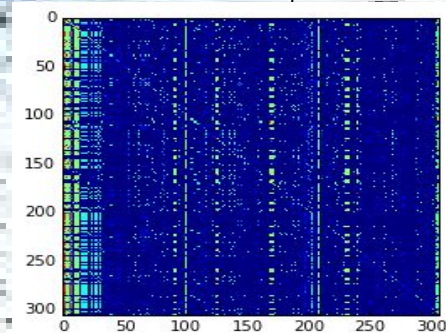
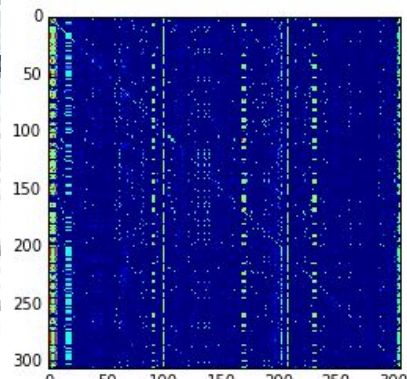


k=5

k=10

k=15

k=20



Summary

Heidi generated
2-d representation
of higher
dimensional data

Presents spatial
overlap among
clusters in
various
subspaces.

Presents nearest
neighbour
proximity
information
among
datapoints.