

Clustering

Hierarchical Agglomerative Clustering (HAC)

Learning Objectives

- ▶ What's the **Hierarchical Agglomerative Clustering (HAC) Algorithm**?
 - ▶ What's hierarchical clustering? What's agglomerative? How does it work?
 - ▶ How does it compare to k -Means?
- ▶ What is a measure of **dissimilarity between points**? How about **dissimilarity between clusters**? What's **linkage**?
- ▶ How does **(high) dimensionality of data** impacts **metrics-based clustering techniques** such as HAC and k -Means?
- ▶ Afternoon Assignment
 - ▶ Topic modeling with k -Means
 - ▶ Implement HAC with *scipy*
 - ▶ Topic modeling with HAC

Overview

Hierarchical Clustering

- Definition

- Example

HAC Algorithm

- Pseudocode

- Step-through

- Linkage

- Choosing k

The curse of dimensionality

Overview

Hierarchical Clustering

- Definition

- Example

HAC Algorithm

- Pseudocode

- Step-through

- Linkage

- Choosing k

The curse of dimensionality

Hierarchical Clustering

- ▶ Type of **agglomerative** clustering
 - ▶ I.e., we will **iteratively group** observations together based on their **distance** from one another
 - ▶ As we continue to group observations together we form a hierarchy of their **similarity** to one another
- ▶ This will force us to answer different questions than we did in *k*-Means
 - ▶ No longer do we have to choose the number of clusters up front
 - ▶ Instead we'll have to define the **nature of successive grouping** of observations

Overview

Hierarchical Clustering

- Definition

- Example

HAC Algorithm

- Pseudocode

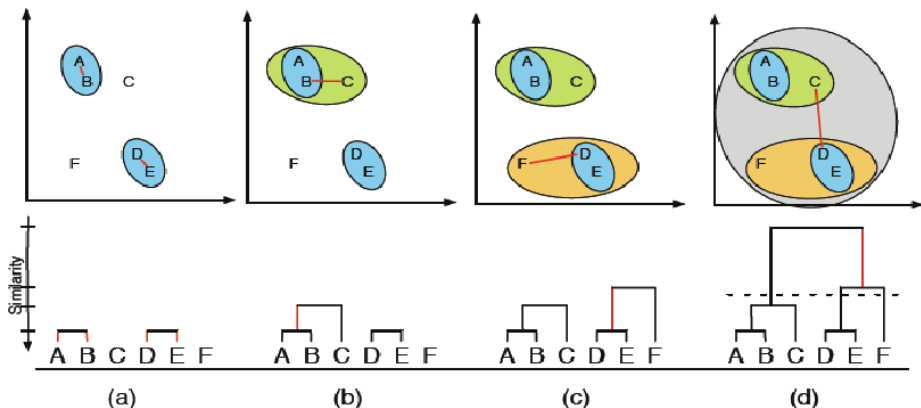
- Step-through

- Linkage

- Choosing k

The curse of dimensionality

Example of Hierarchical Agglomerative Clustering



Overview

Hierarchical Clustering

- Definition

- Example

HAC Algorithm

- Pseudocode

- Step-through

- Linkage

- Choosing k

The curse of dimensionality

HAC

The algorithm in all its glory:

1. Each point as its own cluster
2. Merge "closest" clusters
3. End when all data points are in a single cluster

Like k -Means, this training algorithm may look pretty simple...
and that's because it is

Overview

Hierarchical Clustering

Definition

Example

HAC Algorithm

Pseudocode

Step-through

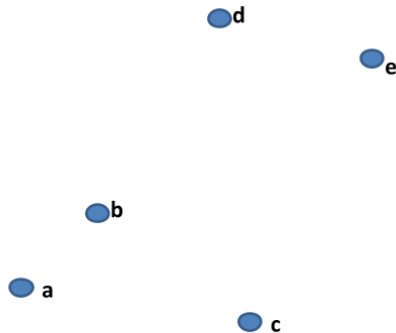
Linkage

Choosing k

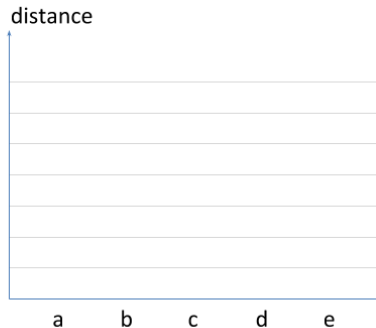
The curse of dimensionality

Step-by-step Execution: DATA!!

- 1 - Compute distances between observations
- 2 - Identify/choose a minimum
- 3 - Fuse observations



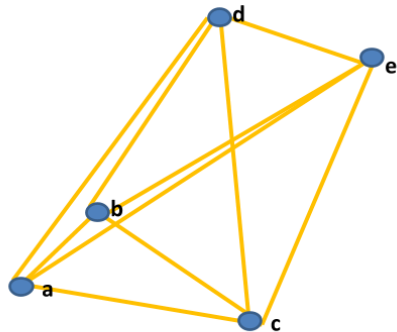
Observations



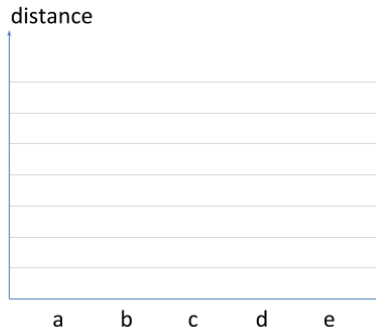
Dendrogram

Step-by-step Execution: Iteration 1 - Compute

- 1 - **Compute** distances between observations
- 2 - Identify/choose a minimum
- 3 - Fuse observations



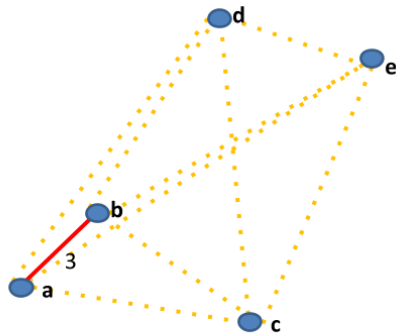
Observations



Dendrogram

Step-by-step Execution: Iteration 1 - Identify

- 1 - Compute distances between observations
- 2 - **Identify**/choose a minimum
- 3 - Fuse observations



Observations

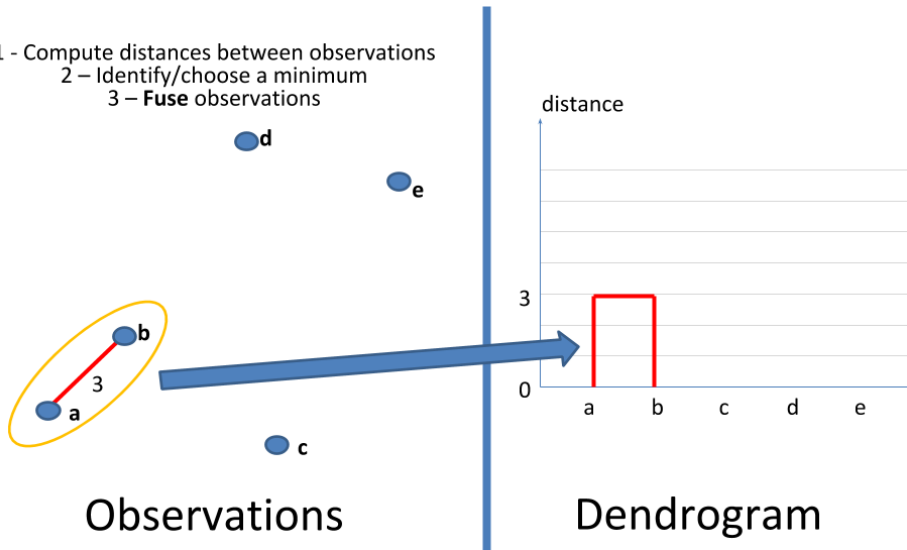
distance



Dendrogram

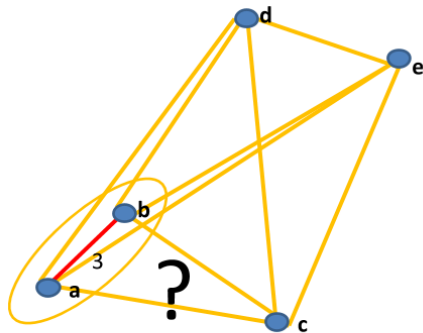
Step-by-step Execution: Iteration 1 - Fuse

- 1 - Compute distances between observations
- 2 - Identify/choose a minimum
- 3 - **Fuse** observations

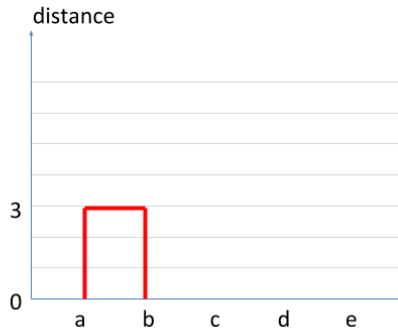


Step-by-step Execution: Iteration 2 - Compute

- 1 - **Compute** distances between observations
- 2 - Identify/choose a minimum
- 3 - Fuse observations



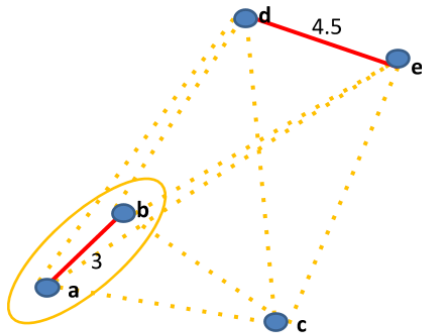
Observations



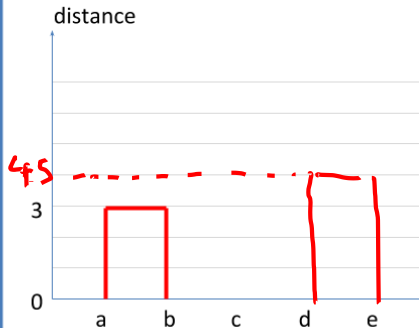
Dendrogram

Step-by-step Execution: Iteration 2 - Identify

- 1 - Compute distances between observations
- 2 - **Identify**/choose a minimum
- 3 - Fuse observations



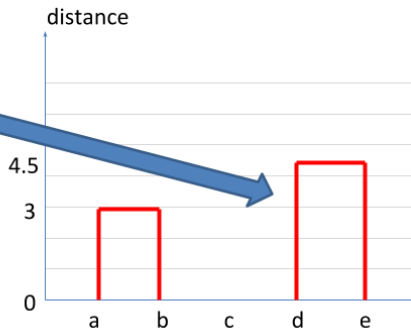
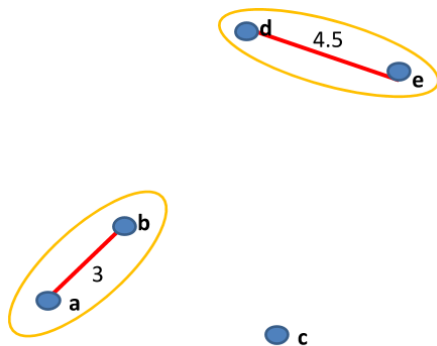
Observations



Dendrogram

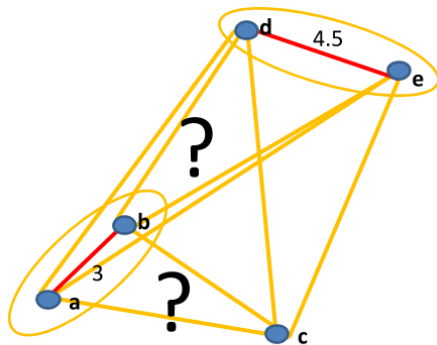
Step-by-step Execution: Iteration 2 - Fuse

- 1 - Compute distances between observations
- 2 - Identify/choose a minimum
- 3 - **Fuse** observations



Step-by-step Execution: Iteration 3 - Compute

- 1 - **Compute** distances between observations
- 2 - Identify/choose a minimum
- 3 - Fuse observations



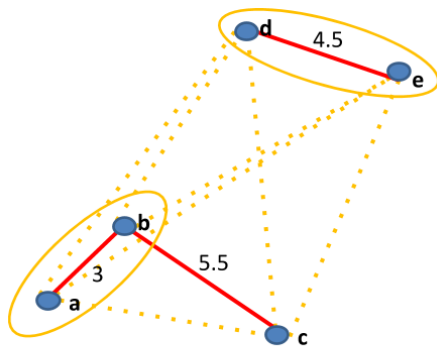
Observations



Dendrogram

Step-by-step Execution: Iteration 3 - Identify

- 1 - Compute distances between observations
- 2 - **Identify**/choose a minimum
- 3 - Fuse observations



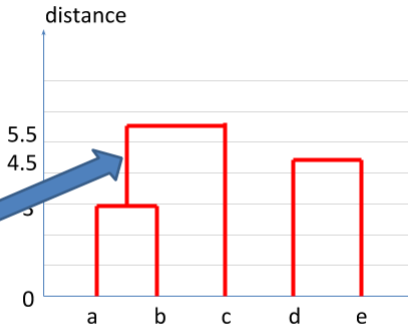
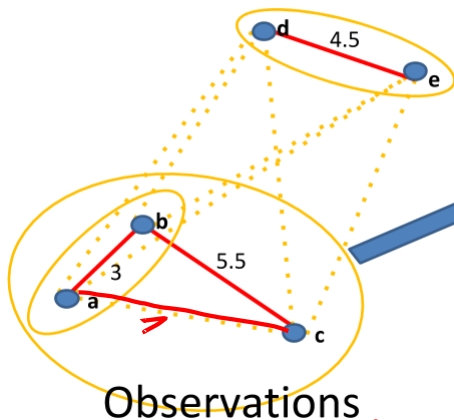
Observations



Dendrogram

Step-by-step Execution: Iteration 3 - Fuse

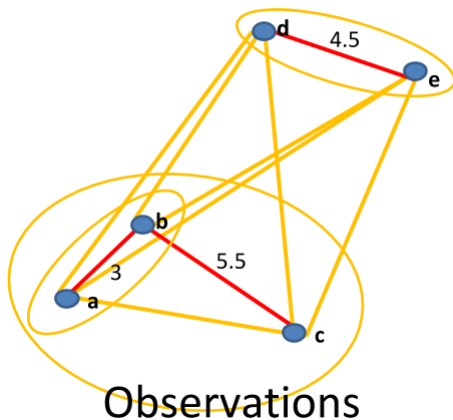
- 1 - Compute distances between observations
- 2 - Identify/choose a minimum
- 3 - **Fuse** observations



$$D(ab, c) = \frac{1}{2 \times 1} (d(a, c) + d(b, c)) = \dots$$

Step-by-step Execution: Iteration 4 - Compute

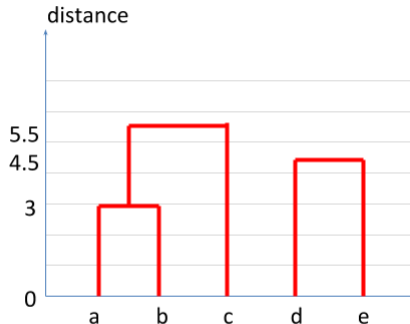
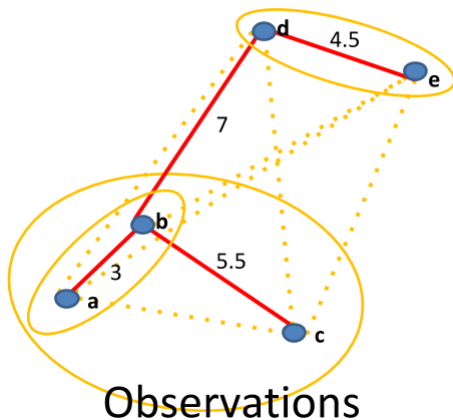
- 1 - **Compute** distances between observations
- 2 - Identify/choose a minimum
- 3 - Fuse observations



Dendrogram

Step-by-step Execution: Iteration 4 - Identify

- 1 - Compute distances between observations
- 2 - **Identify**/choose a minimum
- 3 - Fuse observations



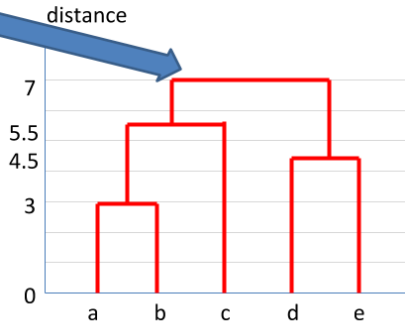
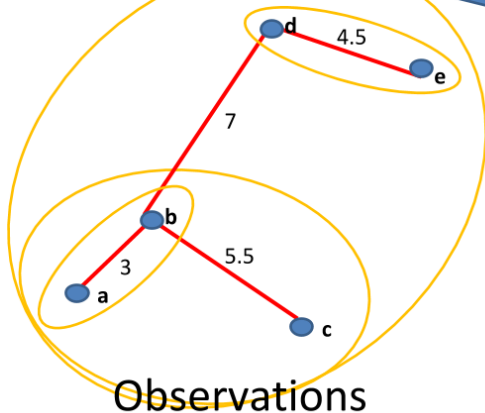
Dendrogram

Step-by-step Execution: Iteration 4 - Fuse

1 - Compute distances between observations

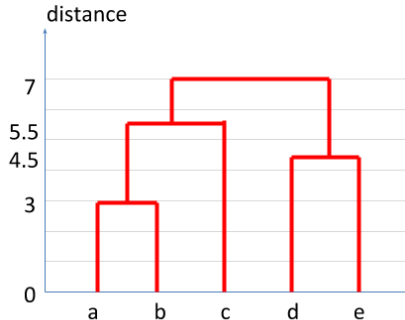
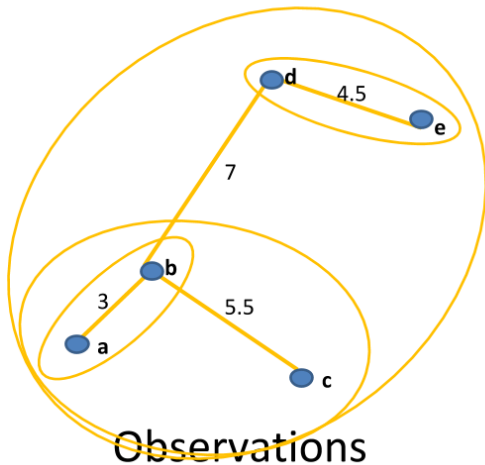
2 - Identify/choose a minimum

3 - Fuse observations



STOP !
Dendrogram

HAC: Final Dendrogram



Dendrogram

Overview

Hierarchical Clustering

Definition

Example

HAC Algorithm

Pseudocode

Step-through

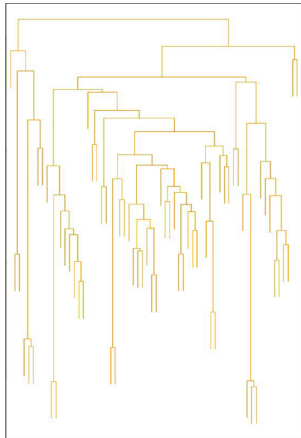
Linkage

Choosing k

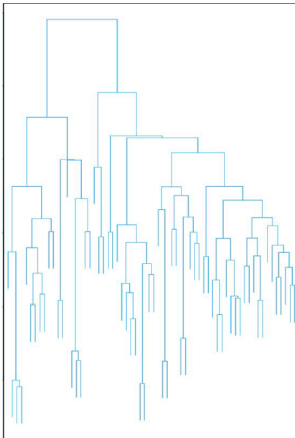
The curse of dimensionality

Distance between clusters?

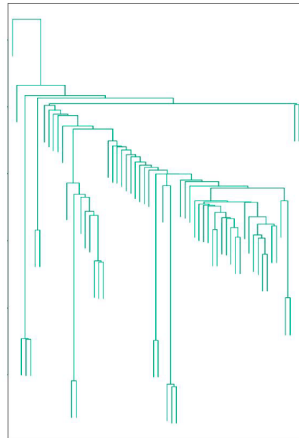
Average Linkage



Complete Linkage



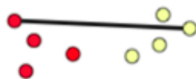
Single Linkage



Complete Linkage

Maximal intercluster dissimilarity

- ▶ Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the **largest** of these dissimilarities



$$D(A, B) = \max_{a \in A, b \in B} D(a, b)$$

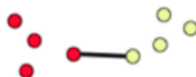
- ▶ Pro: Balanced clusters
- ▶ Cons: More sensitive to outliers; may violate "closeness"
 - ▶ Forces "spherical" clusters with consistent "diameter"

→ More commonly used

Single Linkage

Minimal intercluster dissimilarity

- ▶ Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the smallest of these dissimilarities



$$D(A, B) = \min_{a \in A, b \in B} D(a, b)$$

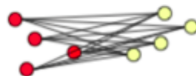
- ▶ Pro: Less sensitive to outliers; handles irregular shapes fairly naturally
- ▶ Con: Extended, trailing clusters
 - ▶ Can fuse single observations one-at-a-time, producing long chains $a \rightarrow b \rightarrow \dots \rightarrow z$

→ Less commonly used

Average Linkage

Mean intercluster dissimilarity

- ▶ Compute all pairwise dissimilarities between the observations in cluster A and the the observations in cluster B, and record the **average** of these dissimilarities



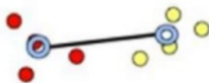
$$D(A, B) = \frac{1}{|A||B|} \sum_{a \in A} \sum_{b \in B} D(a, b)$$

- ▶ Compromise between complete and single linkages
- ▶ Pro: Balanced clusters
 - ▶ Less affected by outliers

→ More commonly used

Centroid Linkage

Dissimilarity between the centroid for cluster A and the centroid for cluster B



$$D(A, B) = D\left(\frac{1}{|A|} \sum_{a \in A} \vec{a}, \frac{1}{|B|} \sum_{b \in B} \vec{b}\right)$$

- ▶ Centroid linkage can result in undesirable inversions
- Not as commonly used, though popular in Genomics

Overview

Hierarchical Clustering

- Definition

- Example

HAC Algorithm

- Pseudocode

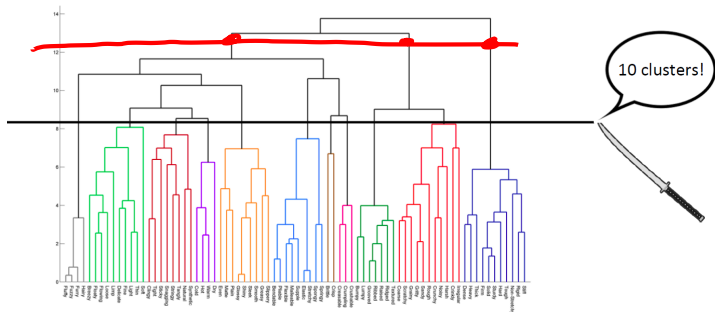
- Step-through

- Linkage

- Choosing k

The curse of dimensionality

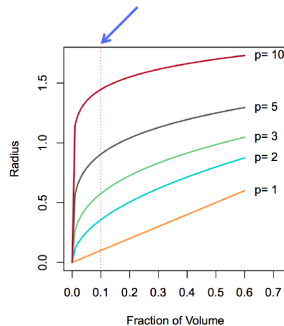
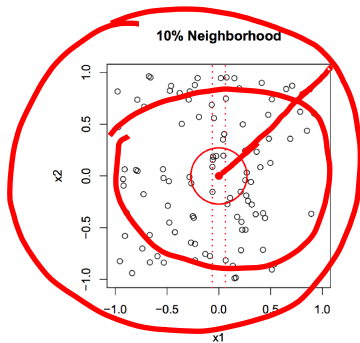
Choosing k



- ▶ In contrast to k -Means, we don't have to choose k from the start
- ▶ Depending on where precisely we cut, we have anywhere from 1 to n clusters

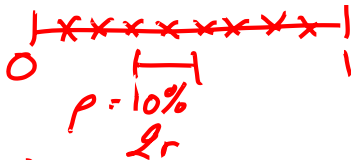
The curse of dimensionality

Just as nearest neighbors breaks down in high dimensional space...
Distance based clustering breaks down in high dimensional space...



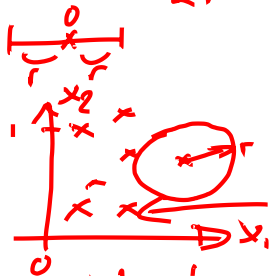
The curse of dimensionality

1D



$$\frac{2r}{1} = p \quad r = \frac{p}{2} = 0.05$$

2D

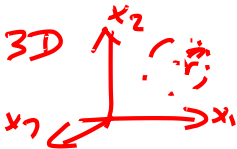


$$\frac{\pi r^2}{1} = p$$

$$r = \sqrt{\frac{p}{\pi}} \\ = \sqrt{\frac{0.1}{3.14}} \\ = 0.178$$

dataset

$$k=5 \rightarrow 10\% \text{ density} \\ n=50$$

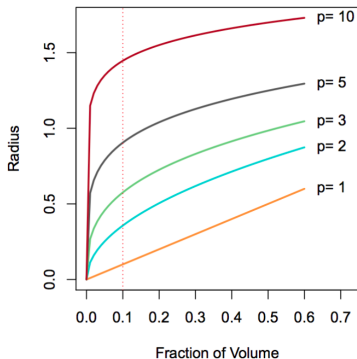


$$\frac{\frac{4}{3}\pi r^3}{1} = p$$

$$r = \left(\frac{p}{\frac{4}{3}\pi} \right)^{1/3} = .287$$

The curse of dimensionality

Can you work out some of the points on the plot?



Dimension	Volume of a ball of radius R	Radius of a ball of volume V
0	1	All balls have volume 1
1	$2R$	$V/2$
2	πR^2	$\frac{V^{1/2}}{\sqrt{\pi}}$
3	$\frac{4}{3}\pi R^3$	$\left(\frac{3V}{4\pi}\right)^{1/3}$
4	$\frac{\pi^2}{2}R^4$	$\frac{(2V)^{1/4}}{\sqrt{\pi}}$
5	$\frac{8\pi^2}{15}R^5$	$\left(\frac{15V}{8\pi^2}\right)^{1/5}$
6	$\frac{\pi^3}{6}R^6$	$\frac{(6V)^{1/6}}{\sqrt{\pi}}$
7	$\frac{16\pi^3}{105}R^7$	$\left(\frac{105V}{16\pi^3}\right)^{1/7}$
8	$\frac{\pi^4}{24}R^8$	$\frac{(24V)^{1/8}}{\sqrt{\pi}}$
9	$\frac{32\pi^4}{945}R^9$	$\left(\frac{945V}{32\pi^4}\right)^{1/9}$
10	$\frac{\pi^5}{120}R^{10}$	$\frac{(120V)^{1/10}}{\sqrt{\pi}}$