

Unsupervised learning

Mostly clustering

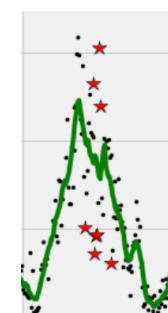
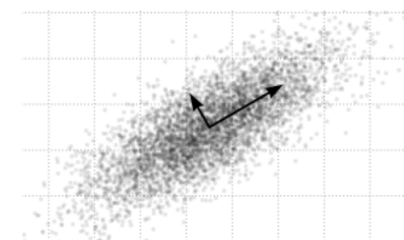
Terence Parr

MSDS program

University of San Francisco

Unsupervised learning techniques...

- In a nutshell, we have just X not $X \rightarrow y$ and would like to know about X , such as density or interesting subregions/vectors of X ($n \times p$ matrix)
- Principle components analysis (orthogonal vectors of most variation)
- Page rank for ranking most important articles/nodes
- Word embeddings like glove (matrix factorization)
- Anomaly detection (fraud or network attack detection)



"Almost all of AI's recent progress is through one type, in which some input data (A) is used to quickly generate some simple response (B)."

Andrew Ng in *What Artificial Intelligence Can and Can't Do Right Now*
Harvard Business Review November 9, 2016

Most common unsupervised learning

- Clustering (unsupervised classifier; i.e., no known classes)
 - k -means / k -medoid / mean-shift
 - hierarchical clustering
 - spectral clustering; graph connecting observations (nodes) by distance-labeled edges
- Recommendation engines (this stuff works great)
 - collaborative filtering (“other people like you bought X”); best done with embeddings; e.g., see [1]
 - market basket analysis / association rules; see *a priori algorithm* (“what do people buy together?”)

[1] https://github.com/fastai/fastbook/blob/master/08_collab.ipynb

The problem with clustering is...

- Clustering sounds awesome but useful only in limited circumstances, such as vector quantization & compression, and usually only when p is small
- Generally doesn't work well with imbalanced data sets such as fraud or network attack classification
- There no clear measure of success, such as the metrics used by supervised learning; e.g., you have bank transactions and no idea which are fraudulent; design algorithm to identify fraud; now, how do you know if your algorithm works?
- You can measure cluster centroid separation, but it still doesn't truly indicate proper clustering; might have too many k etc...

Is clustering really what we want anyway?

- Imagine clustering a customer db into 4 clusters; now what?
- You have 4 groups of, say, 4 million records each; what do you do with it?
- Let's say you can identify marketing-related groups like “technerd”, “shoeshopper”, etc... What do you do with that info?
- Old joke: You know you're wasting half of your marketing money; you just don't know which half!
- Can try to market to those groups but don't we really want to know what kind of ad people click on? Run an ad campaign and track customer->clicks; now you have a supervised problem

Instead can try semi-supervised learning

- Sometimes getting labels is expensive or difficult, such as in medicine
- Still, it's good idea to try to turn unsupervised into supervised learn problem so let's try to start with this “kernel” and gradually broaden the labeled data
- Use a few labeled observations to get things started, such as picking the initial centroids (cluster centers); try to get your client to give you class labels for a few observations
- For a good summary, see **An overview of proxy-label approaches for semi-supervised learning** by S. Ruder
<https://ruder.io/semi-supervised/index.html#selftraining>

One possible self-training procedure

1. Get small initial X^0, y^0 labeled training set from X
2. Train supervised model M^0 on initial X^0, y^0 set
3. Use model M^0 to make predictions for $X \setminus X^0$
4. Combine highest confidence predictions with X^0, y^0 to get, new larger labeled set X^1, y^1
5. Train model M^1 on X^1, y^1
6. Repeat until all X are labeled or no high confidence obs.

Selecting “highest confidence” metric requires experimentation and requires accurate confidence or probabilities from the model

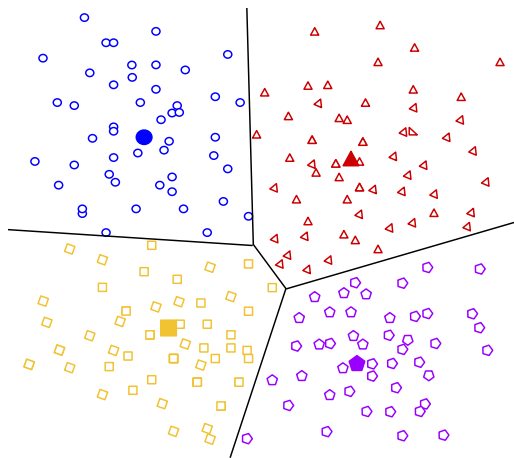
E.g., see **Semi-Supervised Self-Training of Object Detection Models**
<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.68.3602&rep=rep1&type=pdf>

Clustering

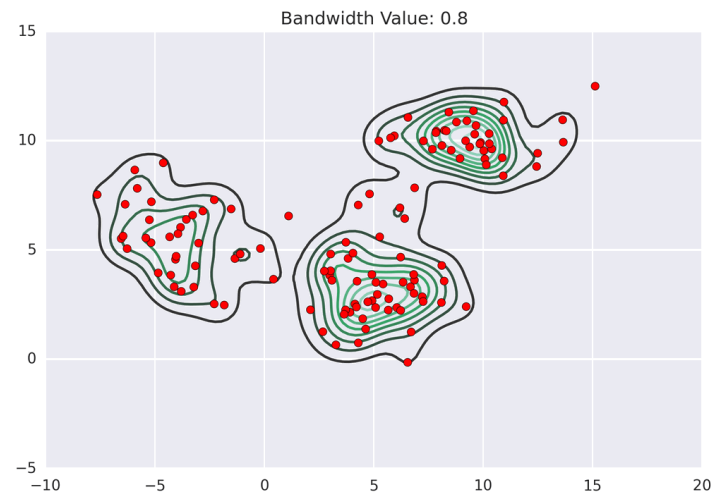
Clustering preliminaries

- Each $x^{(i)}$ in X is a point in p space with p coordinates
- Space can be Euclidean but categorical vars present a challenge
- All such spaces must have distance($x^{(i)}, x^{(j)}$) measure
- Often we need to normalize x values so distance means same thing in all directions
- L_1 and L_2 are common distances for Euclidean space
- L_∞ also useful: max abs difference in any dimension
- For large p and/or binary values, better to use cosine similarity (angle between 2 vectors); $\cos(\theta) = \frac{v \cdot w}{\|v\| \|w\|}$ so $1 - \cos(\theta)$ is distance
- Two flavors: point-assignment and agglomerative

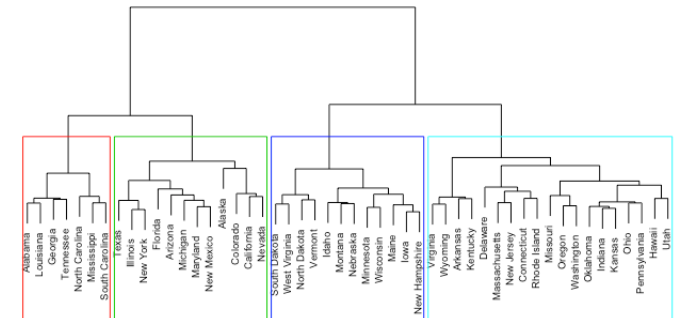
Clustering examples



k-means



mean-shift

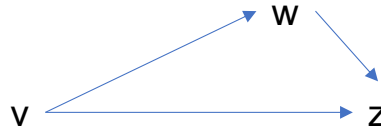


hierarchical clustering

Images: <https://developers.google.com/machine-learning/clustering/clustering-algorithms>,
<https://spin.atomicobject.com/2015/05/26/mean-shift-clustering/>, https://uc-r.github.io/hc_clustering

Distance measure requirements

1. Always nonnegative; only $\text{distance}(v,v)$ is 0
2. Symmetry; $\text{distance}(v,w) = \text{distance}(w,v)$
3. Triangle inequality; $\text{distance}(v,w) + \text{distance}(w,z) \geq \text{distance}(v,z)$

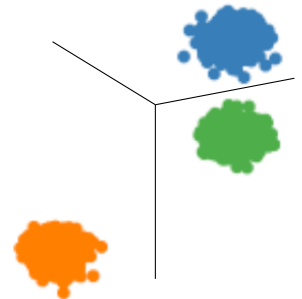


From <http://www.mmds.org/>; see it for more on distance metrics

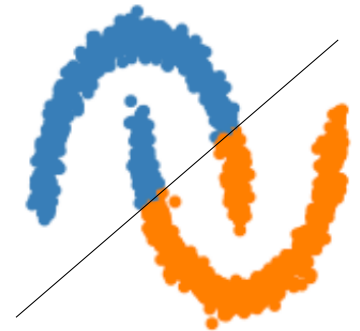
k -means clustering

- Assumes Euclidean space
- Clusters separated by straight lines only
- User provides X and number of clusters to find, k
- Idea is to pick k centroids in p space and assign points to cluster with closest centroid then recompute centroids
- Repeat until the cluster assignments stop changing
- Can select k points as initial centroids:
 - At random (seems to do a crappy job for large p)
 - By picking k distant points (k -means++ is a variation for initial selection)
- Algorithm converges using Euclidean distance
- Not guaranteed to find optimal clusters

k-means works

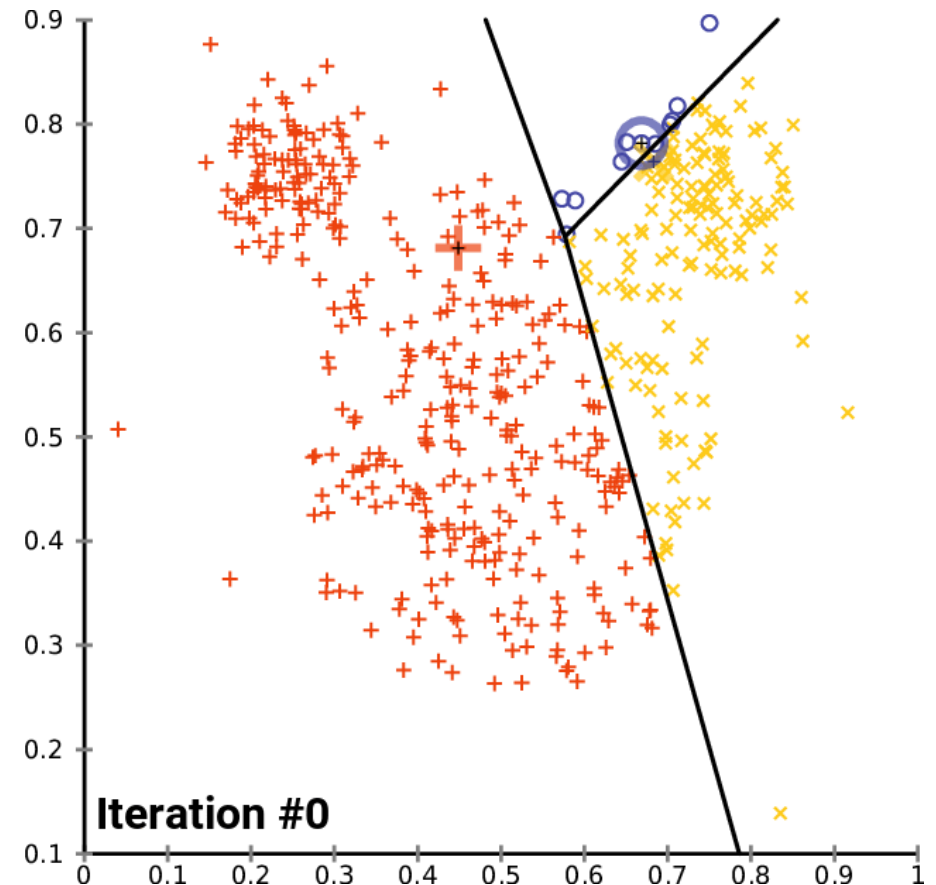


k-means fails



k -mean clustering animation

- k -means gives Voronoi tessellation
- It assumes that all points in the cluster are contiguous; sounds obvious, but for most real problems this assumption doesn't hold
- If true clusters are noncontiguous, k -means will give poor results



Animation from https://en.wikipedia.org/wiki/K-means_clustering

k-means algorithm

Algorithm: $kmeans(X, k)$

Select k unique points from X as initial centroids $m_{1..k}^{(t=0)}$ for clusters $C_{1..k}^{(t=0)}$

repeat

foreach $x \in X$ **do**

$j^* = \arg \min_j \text{distance}(x, m_j^{(t)})$ (*find closest centroid to x*)

 Add x to cluster $C_{j^*}^{(t+1)}$ (*assign x to cluster*)

end

for $j = 1..k$ **do**

$m_j^{(t+1)} = \frac{1}{|C_j^{(t+1)}|} \sum_{x \in C_j^{(t+1)}} x$ (*recompute centroids*)

end

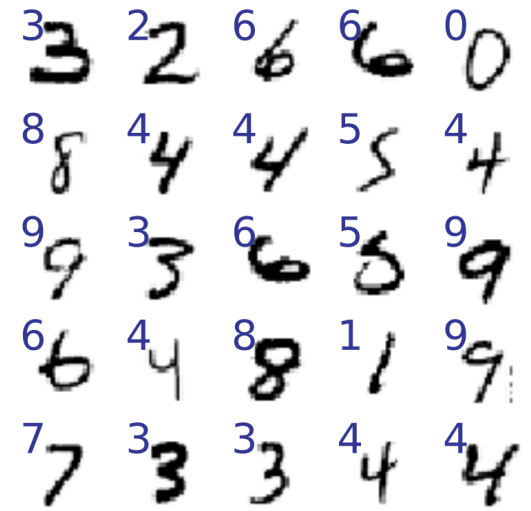
$t = t + 1$

until $C_{1..k}^{(t)} = C_{1..k}^{(t-1)}$ (*until clusters don't change*)

k-means application: MNIST

(Known digits so we can measure error)

- Goal: cluster MNIST digit greyscale images
- $p=28 \times 28=784$ pixels/image
- Results vary a lot and depends a lot on initial cluster centroid selection
- *k*-means doesn't work that well ☹️
- *k*-means finds $k=10$ clusters but cluster number 5 doesn't necessarily correspond to the images of fives
- So, to test quality, let's use known y to group images by digit then ask for impurity of predictions for each image group



Compare clustering of multiple image subsets with random initial centroids (doesn't seem to work well)

```
for i in range(10): # do 10 trials
    df_subset = df_digits.sample(n=1000, replace=False)
    X = df_subset.drop('digit', axis=1) # get just pixels
    y = df_subset['digit']
    X = normalize(X)
    kmeans = KMeans(10, init='random')
    kmeans.fit(X)
    y_pred = kmeans.predict(X)
    # digit_idxes[d] tells us which X indexes are for digit d
    digit_idxes = [np.where(y==d)[0] for d in range(10)]
    # we don't know which of k clusters is for which digit but we can
    # ask for the impurity of y_pred associated with known digit d group
    print([round(gini(y_pred[digit_idxes[d]]),2) for d in range(10)])
```

```
[0.27, 0.55, 0.43, 0.34, 0.46, 0.65, 0.3, 0.53, 0.54, 0.66]
[0.27, 0.51, 0.24, 0.41, 0.54, 0.58, 0.48, 0.53, 0.57, 0.71]
[0.36, 0.51, 0.33, 0.52, 0.52, 0.72, 0.24, 0.4, 0.47, 0.61]
[0.49, 0.53, 0.31, 0.56, 0.59, 0.68, 0.44, 0.51, 0.45, 0.57]
```

...

Gini max is $(k - 1)/k = 9/10 = 0.9$

See <https://github.com/parrt/msds621/blob/master/notebooks/clustering/kmeans.ipynb>



UNIVERSITY OF SAN FRANCISCO

Example predictions in cluster

```
print(y_pred[digit_idxes[0]])  
print('gini =', gini(y_pred[digit_idxes[0]]))
```

```
[9 5 5 5 5 5 9 9 5 5 9 5 5 5 5 5 5 5 5 5 5 5 5 5 9 5 5 9 5 5 9 5 5 5 5 9  
 5 5 5 9 5 5 8 5 5 9 5 5 5 9 5 5 9 9 5 5 5 5 5 9 5 9 5 9 9 2 2 0 9 5 5 5  
 8 9 5 9 9 5 5 9 5 5 5 9 9 5 5 9 5 5 9 5 5 5 9 5 9]  
gini = 0.48219569431690645
```

```
print(y_pred[digit_idxes[9]])  
print('gini =', gini(y_pred[digit_idxes[9]]))
```

```
[1 6 8 6 1 1 1 1 8 6 6 8 4 6 6 8 1 6 3 6 6 6 6 6 1 6 6 8 1 1 4 1 6 8 1 8 1  
 6 8 6 6 6 3 1 1 8 6 3 6 8 6 8 8 8 1 5 6 6 1 6 6 6 1 6 1 6 6 6 1 1 1 6 6 6  
 6 8 6 6 6 6 6 1 1 1 8 6 6 1 1 6 6 6 6 6 1 6 6 6 8 6 6 8 6 8 8 6 6 1 6 6 6  
 6 6 1 6 6 6 6 8 6 1 3 6 6]  
gini = 0.6221383975026016
```

k-means application: Breast cancer

(Known cancer/benign target so we can measure error)

- 212 cancer, 357 non-cancer
- *k*-means isn't great; uncertainty is low for one class (0.0056) but uncertain for the other (.4743)
- To the right are the two possible conf matrices, depending on the cluster number chosen by *k*-means for cancer and for non-cancer

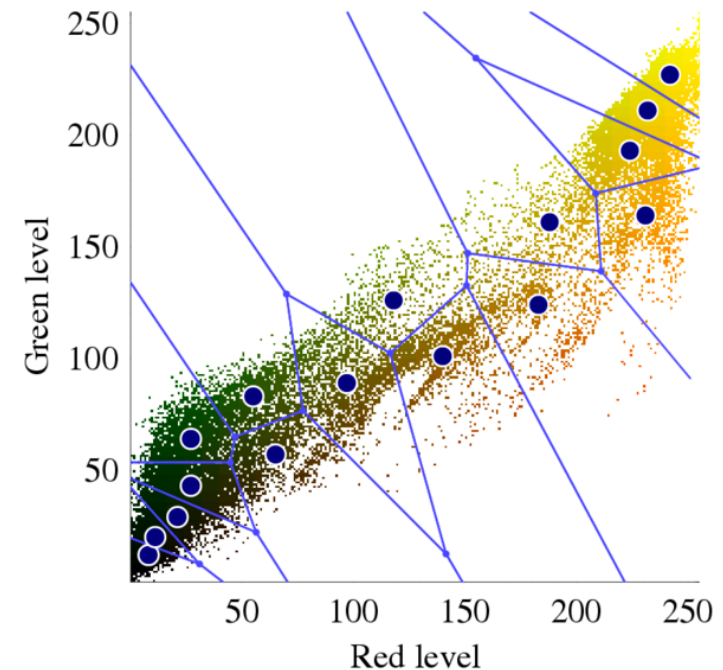
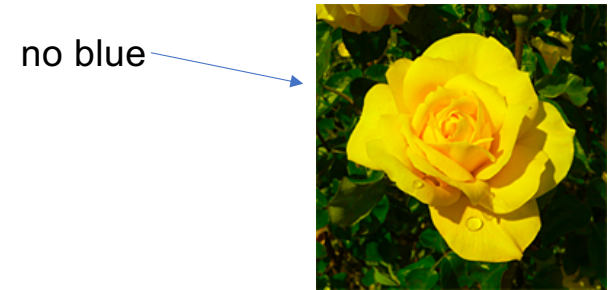
```
kmeans = KMeans(n_clusters=2, init='k-means++')
kmeans.fit(X)
y_pred = kmeans.predict(X)
cancer = np.where(y==0)[0]
benign = np.where(y==1)[0]
print( gini(y_pred[benign]), gini(y_pred[cancer]) )
```

		predicted	
		0	1
actual	0	356	1
	1	82	130

		predicted	
		0	1
actual	0	1	356
	1	130	82

k-means application: color quantization

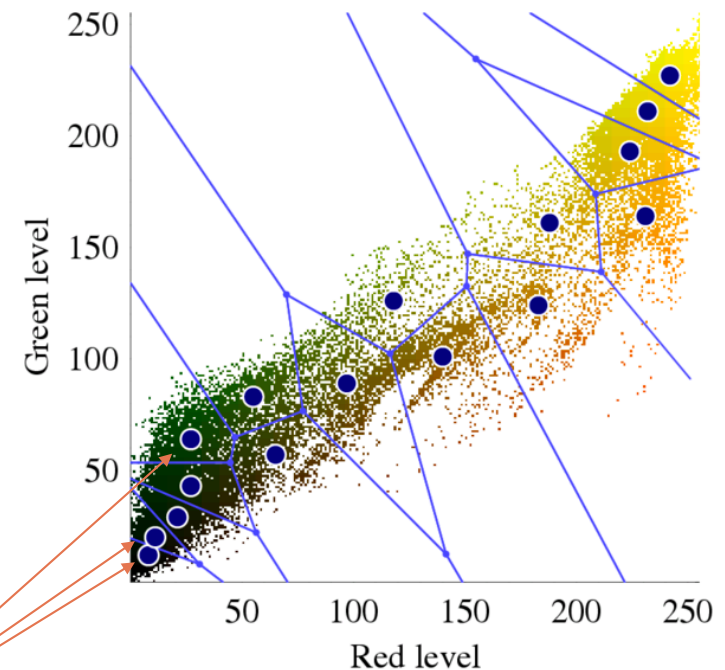
- Color pictures typically use lots of unique colors, possibly 10s of thousands
- Each pixel in the image has Red/Green/Blue colors, 1 byte per RGB = 3 bytes (24 bits)
- Each RGB is a 3D coordinate of color: (R, G, B), with possibly tens of thousands of unique combinations
- Example with just red/green (omit blue):



See https://scikit-learn.org/stable/auto_examples/cluster/plot_color_quantization.html
Image from https://en.wikipedia.org/wiki/Color_quantization

Color quantization cont'd

- (R,G,B) takes 3 bytes per pixel which makes images really big
- If a picture only has 256 unique colors we can map all (R,G,B) vectors to a single byte; the color “index” 0..255 points into a color palette with the full (R,G,B) vectors; 3x compression for each pixel so a massive compression
- If picture has more than 256 colors, we can cluster in RGB space with $k=256$ and it will group similar colors together; then we pick the centroid as the colors in the palette



Color quantization example, $k=10$

Original image
(96,615 colors)



Quantized image
(10 colors, k-Means)



Quantized image
(10 colors, at random)



See <https://github.com/parr/msds621/blob/master/notebooks/clustering/kmeans.ipynb>

Color quantization example, $k=4$

Original image
(96,615 colors)



Quantized image
(4 colors, k-Means)



Quantized image
(4 colors, at random)



Confusion point

- *k*-means' centroids don't have to be points in X , usually aren't
- *k-medians* uses median not mean for centroids (minimizes w.r.t. L1 not L2 distance); median even for single dimension doesn't have to be point in $x^{(i)}$ space
- *k-medoids* (not spelled *k-medoids*) requires medoids to be points in X ; works with any distance measure; sounds like *k*-means but algorithm is pretty different; gotta pick "centrally located point"

Trouble with k -means

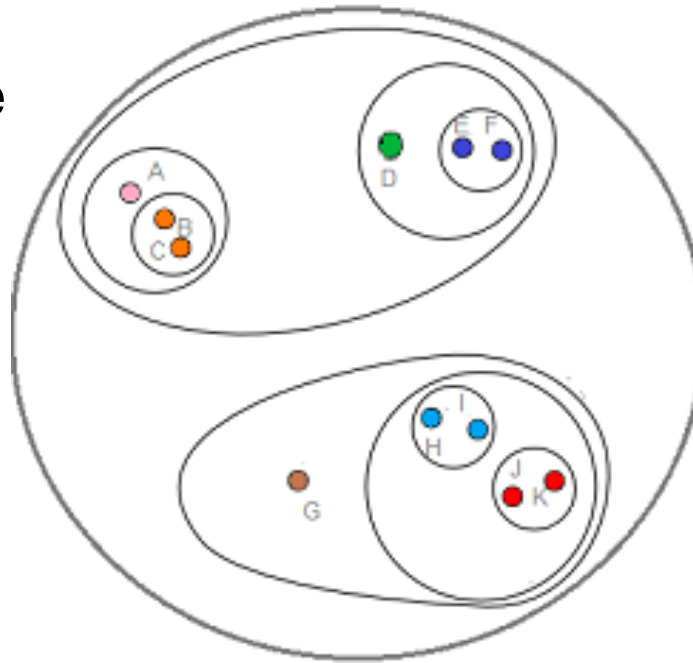
- k -means requires that we specify number of clusters k
- Picking k is usually a problem
- Color quantization and MNIST digits have known k , but few do
- Different starting centroids can lead to very different results
- Each observation is forced into one of the k clusters, but probabilities might be nice; we could use distance to centroid I guess but a density estimate would be better

Hierarchical (agglomerative) clustering

- **Idea:** put every point into its own singleton cluster; repeatedly group two closest clusters into a meta-cluster until there is a single cluster left
- Can also stop when distance between clusters are sufficiently large
- We need a cluster distance metric; called the *linkage criterion*
- Simplest linkage is just the distance between cluster centroids
- Result is a tree of clusters, one cluster per level
- We get all possible clusters

Dendograms

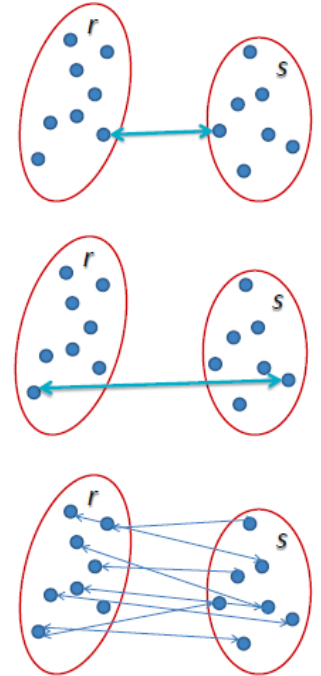
- Dendogram is a tree of clusters, one cluster per level
- Distance from node to children reflects between-cluster-metric



Between-cluster distance metrics (*linkage*)

Warning: statisticians coining terms again!

1. Minimum distance between any two points in the cluster (*single linkage*); tends not to get compact clusters and can get chains of points
2. Max distance between point pairs (*complete linkage*); tends to get compact clusters but points can be closer to other clusters than those within their cluster
3. Average distance of all point pairs from two clusters (*group average linkage*); tries to get compact clusters that are far apart
4. *Ward's method* minimizes within-cluster variance; merge pair with smallest variance at each step



Effect of between-cluster-metric

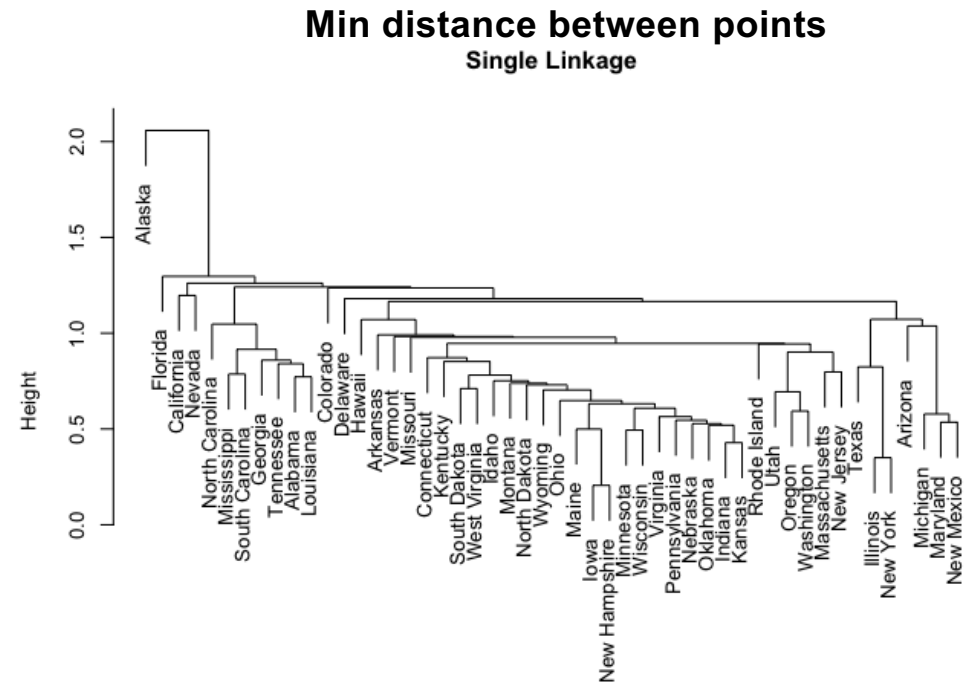
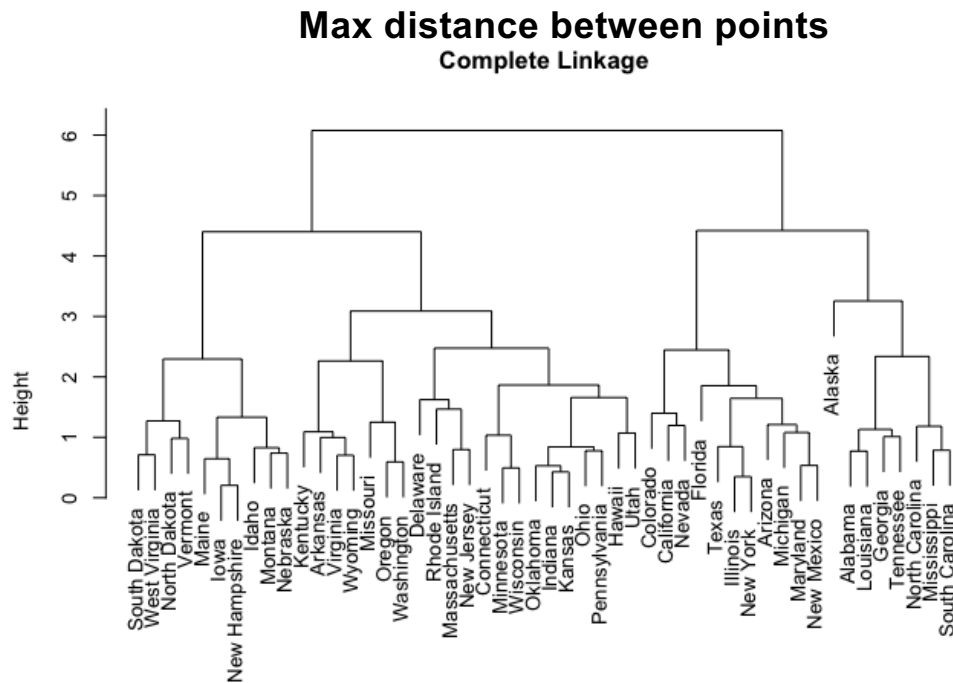


Image from https://uc-r.github.io/hc_clustering

Ward's method

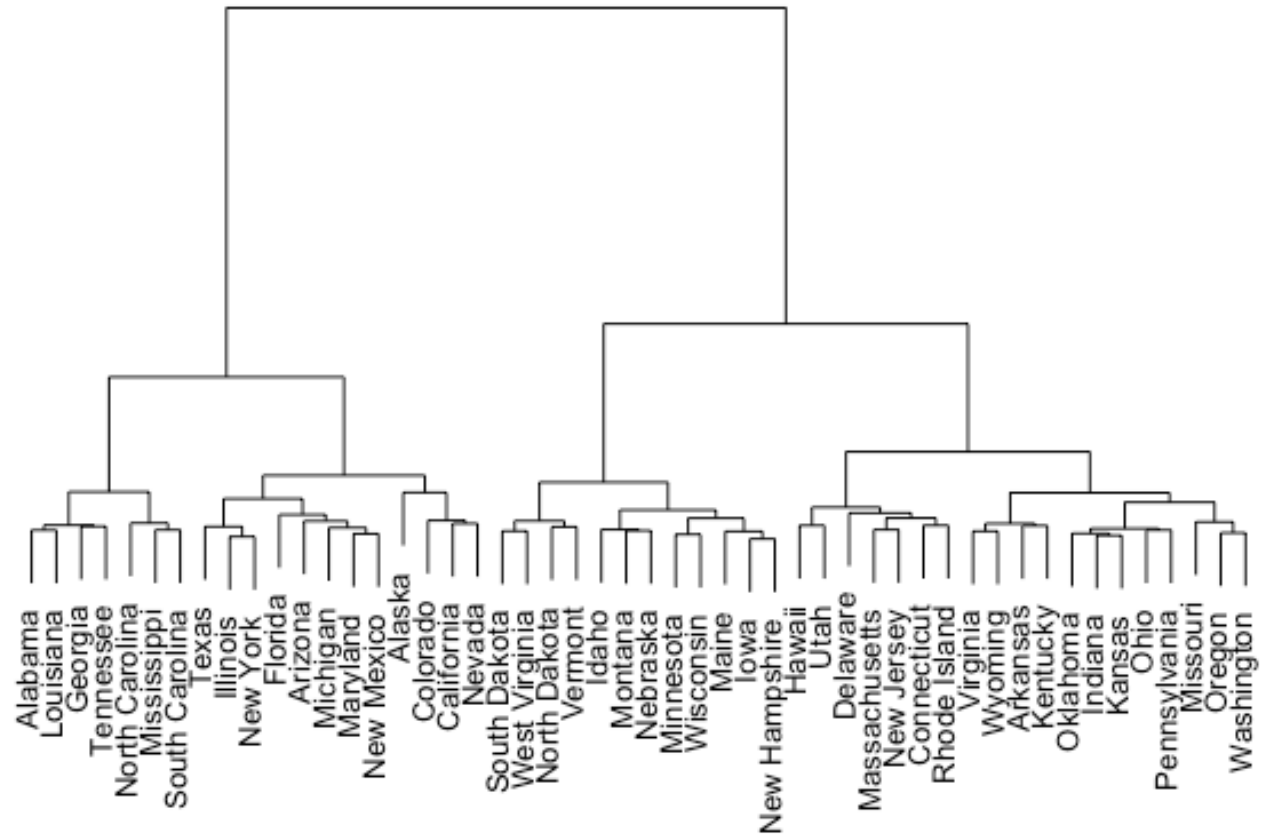


Image from https://uc-r.github.io/hc_clustering

Non-numeric clustering

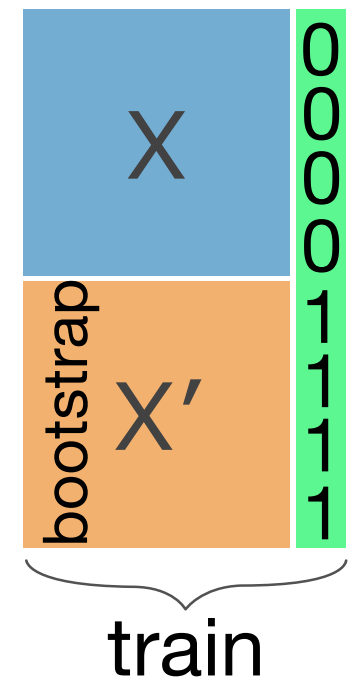
- How do you cluster documents?
- Can use edit distance or Jaccard similarity between text docs
- Maybe convert words to Glove word vectors
- Try your own word embeddings from corpus
- But what about tabular data with nominal categorical variables?
- In non-numeric space, what is a centroid vector?
- There are similarity measures for categoricals, but I'm not a big fan, particularly with mixed numeric and categorical data

Breiman's RF clustering

- Goal: similarity($x^{(i)}, x^{(j)}$) or distance($x^{(i)}, x^{(j)}$) for any two feature vectors in X , even in the presence of mixed categorical and numeric data
- Random Forests to the rescue again with clever trick that turns unsupervised into supervised problem
- Then derive similarity matrix between all $x^{(i)}, x^{(j)}$ pairs
- Proximity matrix: count how often $x^{(i)}, x^{(j)}$ appear in same leaf in all trees of forest; normalize by number of leaves
- Use 1 minus proximity to get distance, then can use any clustering algorithm we want like k -means, ...

Random Forest distance metric

1. Consider all X records as class 0
2. Duplicate and bootstrap columns of X to get X' : class 1
 - Breiman: X' created by “...*sampling at random from the univariate distributions...*” of X
 - X' destroys relationships between columns of X
3. Create y to label/distinguish X vs X'
4. Train RF on stacked $[X, X'] \rightarrow y$
5. Walk all leaves of all trees, bumping proximity $[i, j]$ for all $x^{(i)}, x^{(j)}$ pairs in leaf; divide proximities by num of leaves
6. Cluster using 1-proximity for distance matrix



Breiman's RF gets X' from X

Here's how to create X' from X

```
def df_scramble(X : pd.DataFrame) -> pd.DataFrame:
    X_rand = X.copy()
    for colname in X:
        X_rand[colname] = \
            np.random.choice(X[colname], len(X), replace=True)
    return X_rand
```

Breiman's RF conjures up supervised from unsupervised

```
def conjure_twoclass(X : pd.DataFrame)\
    -> (pd.DataFrame, pd.Series):
    X_rand = df_scramble(X)
    X_synth = pd.concat([X, X_rand], axis=0)
    y_synth = np.concatenate([np.zeros(len(X)),
                              np.ones(len(X_rand))],
                              axis=0)
    return X_synth, pd.Series(y_synth)
```

Train an RF model to recognize structure between variables,
but goal is simply co-existence in leaves

Computing RF similarity matrix

For each tree in RF

For each leaf in tree

Increment similarity for all $x^{(i)}$ and $x^{(j)}$ in leaf (ignoring X' obs.)

Divide each similarity[i,j] by number of leaves to normalize similarities

More on RFs for similarity

- **Similarity Forests**

<http://biorxiv.org/cgi/reprint/258699v1>

- **Unsupervised Learning With Random Forest Predictors**

<https://horvath.genetics.ucla.edu/html/RFclustering/RFclustering/RandomForestHorvath.pdf>