

---

**CSE 519: Data Science**  
**Steven Skiena**  
**Stony Brook University**

---

Lecture 19: Nearest Neighbor Methods

---

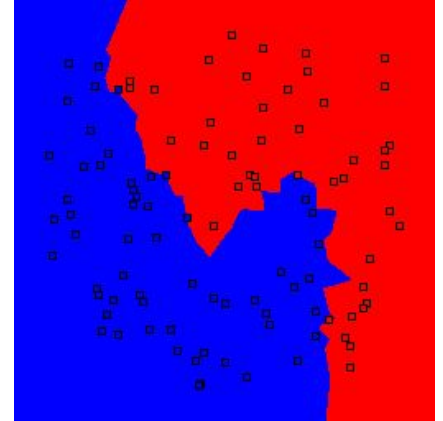
# Nearest Neighbor Classification

---

Identify which training example is most similar to the target, and take the class label from it.

The key issue here is devising the right distance function between rows/points.

**Advantages:** simplicity, interpretability, and non-linearity.



# Distance Metrics

---

Certain mathematical properties are expected of any distance measure, or *metric*:

- $d(x,y) \geq 0$  for all  $x, y$  (positivity)
  - $d(x,y) = 0$  iff  $x=y$  (identity)
  - $d(x,y) = d(y,x)$  (symmetry)
  - $d(x,y) \leq d(x,z) + d(z,y)$  (triangle inequality)
-

# Not a Metric

---

Many natural similarity measures are not distance metrics:

- Correlation coefficient (-1 to 1)
  - Cosine similarity / dot product (-1 to 1)
  - Mutual information measures (often not symmetric)
  - Cheapest airfare (think triangle inequality)
-

# Euclidean Distance Metric

---

The traditional Euclidean distance metric weighs all dimensions equally:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^d |x_i - y_i|^2}$$

We might use a coefficient  $c_i$  to give a different weight to each dimension, but *at least* normalize to make dimensions comparable.

---

# L\_k Distance Norms

---

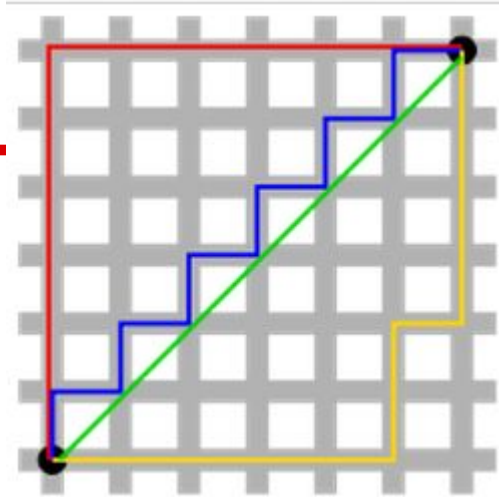
To generalize Euclidean distance:

$$d(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^k \right)^{1/k}$$

- $k=1$  gives the Manhattan distance metric
- $k = \infty$  gives the maximum component

$k$  regulates the trade off between largest and total dimensional difference.

---



# Which Point is Farther from (0,0)?

---

Is  $p_1=(2,0)$  or  $p_2=(1.5,1.5)$  farther from  $(0,0)$ ?

- For  $k=1$ , the distances are 2 and 3, so  $p_2$
- For  $k=2$ , the distances are 2 and 2.12, so  $p_2$
- For  $k=\infty$ , distances are 2 and 1.5, so  $p_1$

The distance metric sets which point is closer.

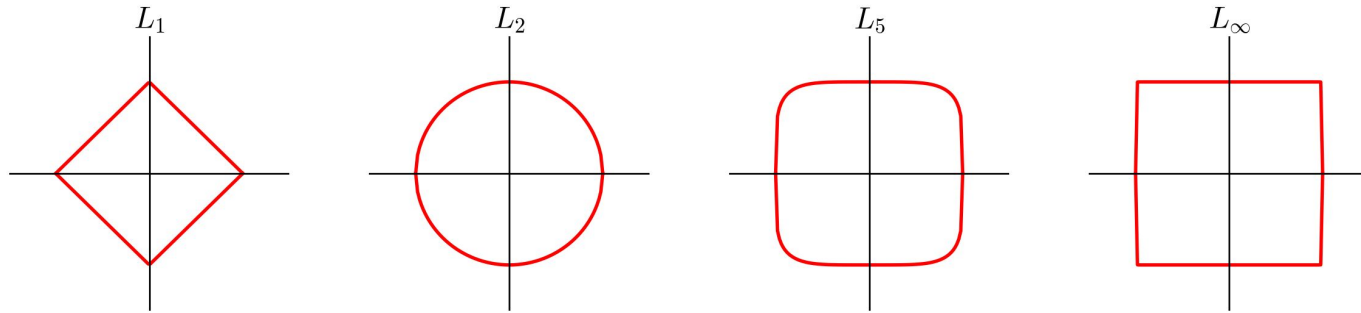
Are we more worried about random noise or dimensional outliers/artifacts?

---

# Circles for Different $k$

---

The shape of the  $L_k$  “circle” governs which points are equal neighbors about the origin.



The distinction here become particularly important in higher dimensional spaces: do we care about deviations in all dimensions or primarily the biggest?

---



# Projections from Higher Dimensions

---

Projection methods (like SVD) compress or ignore dimensions to reduce representation complexity.

Nearest neighbors in such spaces can be more robust than in the original space.

---

# Dimensional Egalitarianism

---

Although  $L_k$  norms in principle weigh all dimensions equally, the scale matters.

But note that the real impact of height will differ depending upon whether it is measured in centimetres, meters, or kilometers.

**This is why we use Z-scores for normalization!**

---

# Regression / Interpolation by NN

---

The idea of nearest neighbor classification can be generalized to function interpolation, by averaging the values of the  $k$  nearest points.

Weighted averaging schemes can value points differently according to (1) distance rank, (2) actual distances.

Similar ideas work for all classification methods.

---

# K-Nearest Neighbors

---

NN classification produce non-linear classifiers, because each training point changes the separating boundary.

More robust classification or interpolation follow from voting over the  $k$  closest neighbors for  $k > 1$ .

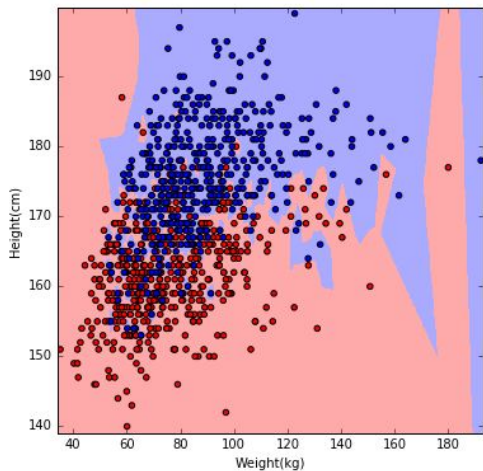
---

# Gender Classification by Height/Weight

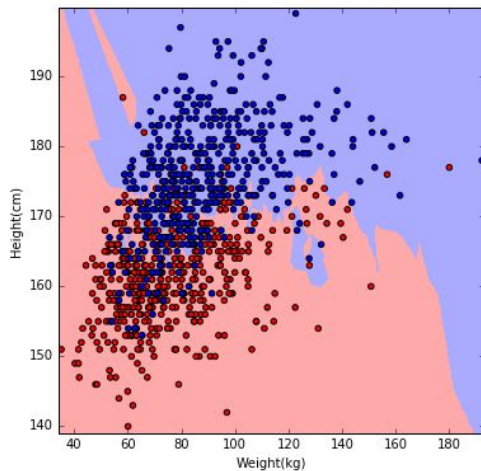
---

Smooother boundaries follow from larger k:

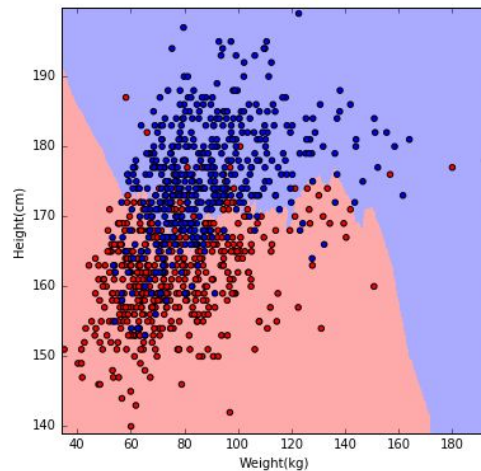
k=1



k=3



k=10



# Seeking Good Analogies

---

Many intellectual disciplines rest on analogies:

- **Law:** which legal precedent was most like this case?
  - **Medicine:** how did I treat patients with similar symptoms, and did they survive?
  - **Real estate:** what price did comparable properties sell for in the neighborhood?
-

# Seeking Good Analogies? (Projects)

---

- Miss Universe?
  - Movie gross?
  - Baby weight?
  - Art auction price?
  - Snow on Christmas?
  - Super Bowl / College Champion?
  - Ghoul Pool?
  - Future Gold / Oil Price?
-

# Finding Nearest Neighbors

---

Given  $n$  points in  $d$ -dimensions, it takes  $O(nd)$  time to find the NN using brute force search.

For large training sets or high dimensionality this becomes very expensive.

This motivates use of more sophisticated data structures: grid indices, kd-trees, Voronoi diagrams, and locality sensitive hashing.

---



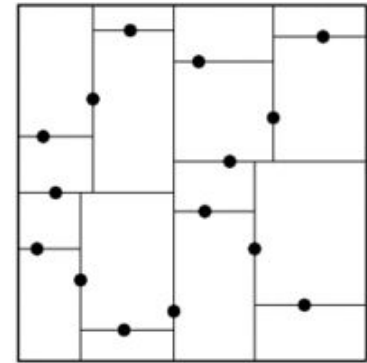
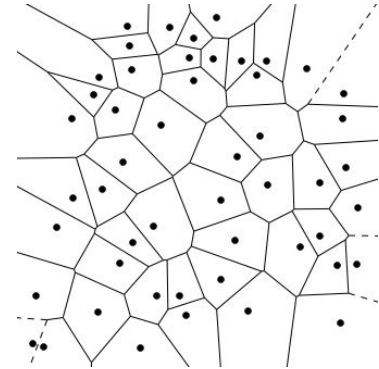
# Voronoi Diagrams / Kd-trees

---

Voronoi diagrams partition space into regions sharing nearest neighbors.

Efficient algorithms exist for finding nearest neighbors in low dimensions, such as kd-trees.

But exact NN search is doomed to reduce to linear search for high-enough dimensionality data.



# Locality Sensitive Hashing

---

Hashing could speed nearest-neighbor search if nearby points got hashed to the same bucket. But normal hashing uses hash functions that spreads similar items to distance buckets.

Locality sensitive hashing (LSH) takes points or vectors  $a$  and  $b$  such that likely  $h(a)=h(b)$  iff  $a$  is near  $b$ .

---

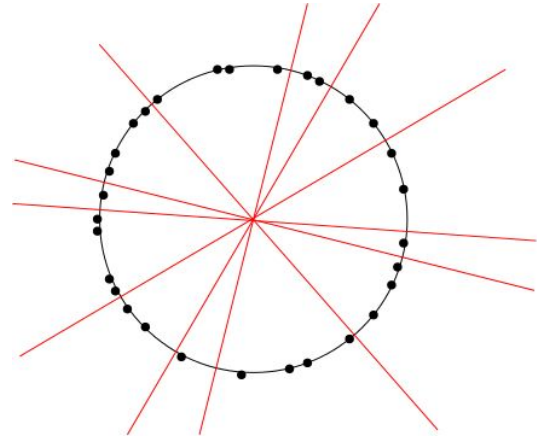
# LSH for Points on a Sphere

---

Pick random planes cutting through the origin.

If near each other, two points are likely on the same side (left or right) of a given random plane.

L/R patterns for  $d$  random planes define a  $d$ -bit LSH hash code.



# Network Data

---

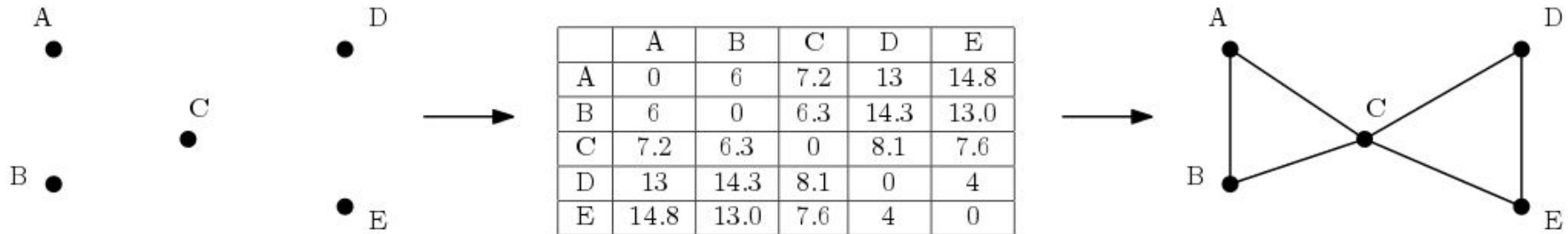
Many datasets have natural interpretations as graphs/networks:

- **Social networks**: vertices are people, edges are friendships.
  - **WWW**: vertices=pages, edges=hyperlinks.
  - **Product/customer networks**: edges=sales.
  - **Genetic networks**: v=genes, e=interactions.
-

# Point Sets and Graphs

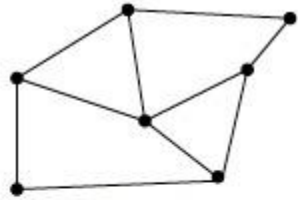
---

- Point sets naturally define graphs: add an edge  $(x,y)$  if  $x$  and  $y$  are close enough.
- Graphs naturally define point sets: perform an SVD of the adjacency matrix.

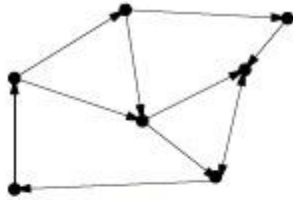


# Taxonomies of Graph Types

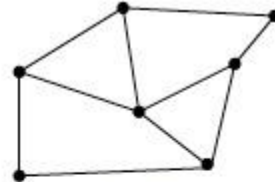
---



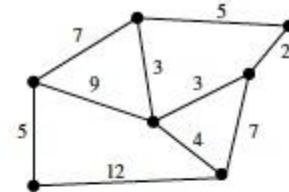
undirected



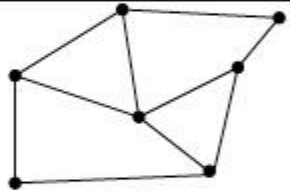
directed



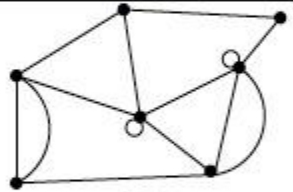
unweighted



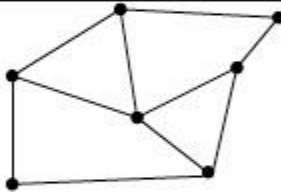
weighted



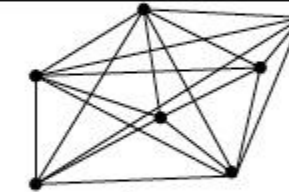
simple



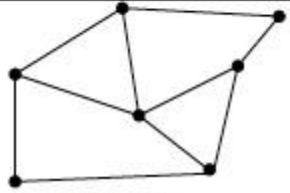
non-simple



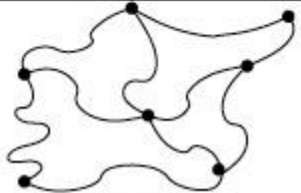
sparse



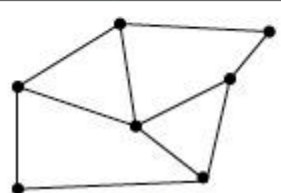
dense



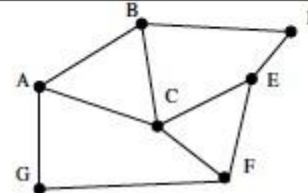
embedded



topological



unlabeled



labeled

Talking  
the talk  
is part  
of  
walking  
the  
walk.

# Classical Graph Algorithms

---

Distances in graphs are naturally defined in terms of the shortest path between vertices.

Classical algorithms for finding shortest paths, connected components, spanning trees, cuts, flows, matchings, topological sorting can be applied to any appropriate network.

---

# PageRank

---

$PageRank(v, G)$  corresponds to the probability that a random walk on  $G$  ends up at vertex  $v$ .

The basic formula is:

$$PR_j(v) = \sum_{(u,v) \in E} \frac{PR_{j-1}(u)}{out-degree(u)}$$



This recursive formula defines an iterative algorithm which quickly converges in practice.

---



# If All Roads Lead to Rome...

---

... then Rome must be an important place.

PageRank is a good feature to measure centrality or importance of vertices, like

Wikipedia pages:

PageRank PR1 (all pages)	
1	Napoleon
2	George W. Bush
3	Carl Linnaeus
4	Jesus
5	Barack Obama
6	Aristotle
7	William Shakespeare
8	Elizabeth II
9	Adolf Hitler
10	Bill Clinton

PageRank PR2 (only people)	
1	George W. Bush
2	Bill Clinton
3	William Shakespeare
4	Ronald Reagan
5	Adolf Hitler
6	Barack Obama
7	Napoleon
8	Richard Nixon
9	Franklin D. Roosevelt
10	Elizabeth II

---

# PageRank in Practice

---

Design decisions in PageRank include:

- Editing the graph to remove irrelevant vertices/edges (spam).
- Removing outdegree zero vertices.
- Allowing random jumps: the damping factor.

PageRank is less important to Google today than popularly supposed.

