

# CSE 417T: Introduction to Machine Learning

## Lecture 19: $k$ -Nearest Neighbors

Henry Chai

11/06/18

# Recall

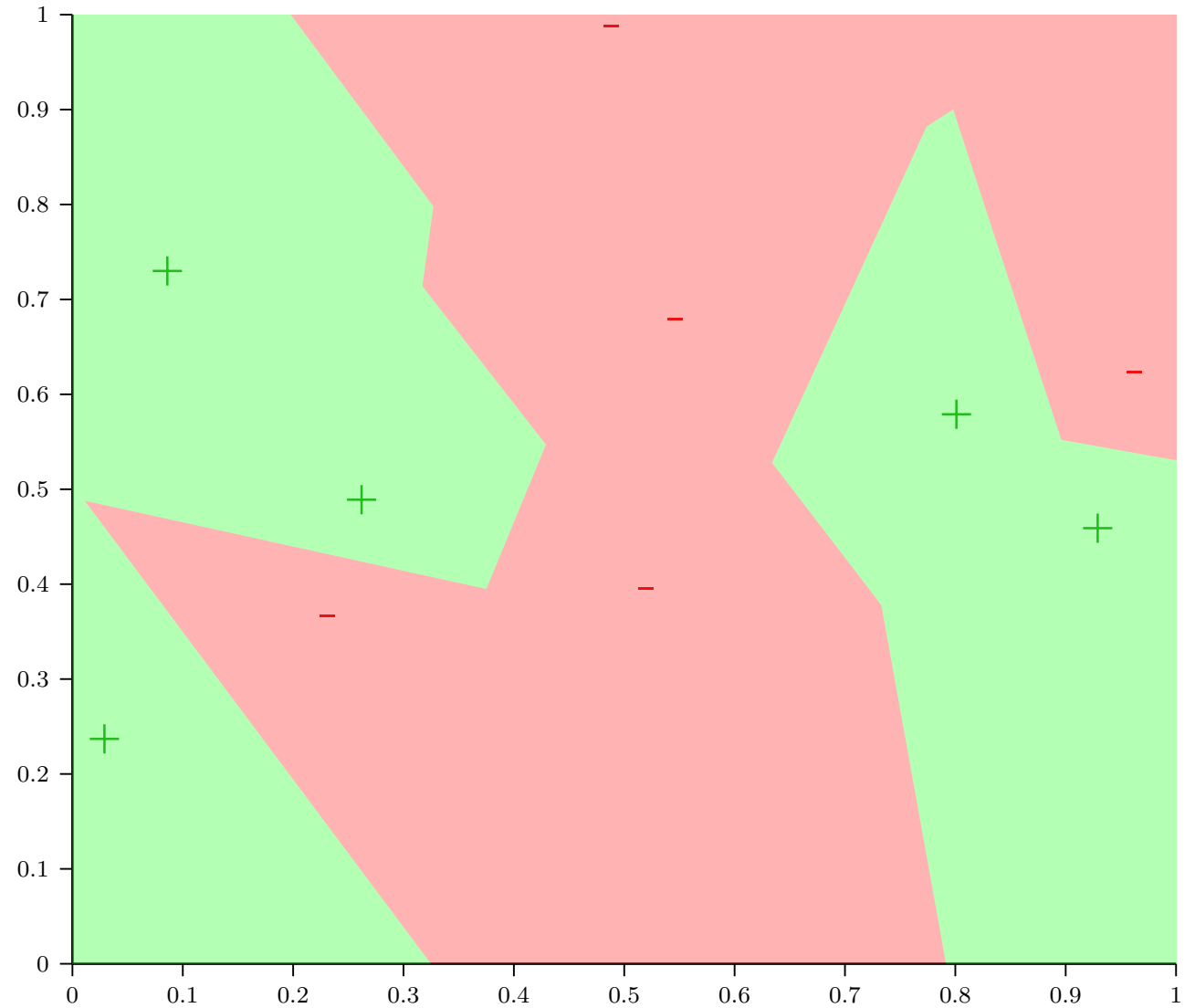
- Classify a point as the label of the “most similar” training point

- Euclidean distance:  $d(\vec{x}, \vec{x}') = \|\vec{x} - \vec{x}'\| = \sum_{j=1}^D (x_j - x'_j)^2$

- Given  $\mathcal{D} = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)\}$  and a point  $\vec{x}$ , let  $\vec{x}_{[i]}(\vec{x})$  be the  $i^{\text{th}}$  closest point to  $\vec{x}$  in  $\mathcal{D}$

# The Nearest Neighbor Hypothesis

$$g(\vec{x}) = y_{[1]}(\vec{x})$$



# Generalization of Nearest Neighbor

- Claim:  $E_{out}$  for the nearest neighbor hypothesis is not much worse than the best possible  $E_{out}$ !
- Formally: under certain conditions, with high probability,  $E_{out}(g) \leq 2E_{out}(g^*)$  as  $n \rightarrow \infty$
- Interpretation: half of the data's predictive power is in the nearest neighbor!

# $k$ -Nearest Neighbors ( $k$ NN)

- Classify a point as the most common label among the labels of the  $k$  nearest training points
- If we have a binary classification problem and  $k$  is odd:

$$g(\vec{x}) = \text{sign} \left( \sum_{i=1}^k y_{[i]}(\vec{x}) \right)$$

- $k$  controls the complexity of the hypothesis set  $\implies k$  affects how well the learned hypothesis will generalize
  - $k = 3$
  - $k = \lfloor \sqrt{n} \rfloor$
  - Cross-validation

# $k$ NN Pros and Cons

- Pros:
  - Intuitive / explainable
  - No training / retraining
  - Self-regularizes
  - Provably near-optimal in terms of  $E_{out}$
- Cons:
  - Computationally expensive
    - Always needs to store all data:  $O(nD)$
    - Computing  $g(\vec{x})$  requires computing  $d(\vec{x}, \vec{x}') \forall \vec{x}' \in \mathcal{D}$  and finding the  $k$  closest points:  $O(nD + n \log(k))$
  - Suffers from the “curse of dimensionality”

# Curse of Dimensionality

- The fundamental assumption of  $k$ NN is that “similar” points or points close to one another should have the same label
- The closer two points are, the more confident we can be that they will have the same label
- As the number of dimensions the input has grows, the less likely it is that two random points will be close
- As the number of dimensions the input has grows, it takes more points to “cover” the input space

# Curing the Curse of Dimensionality

- More data
- Fewer dimensions
- Blessing of non-uniformity: data from the real world is rarely uniformly distributed across the input space

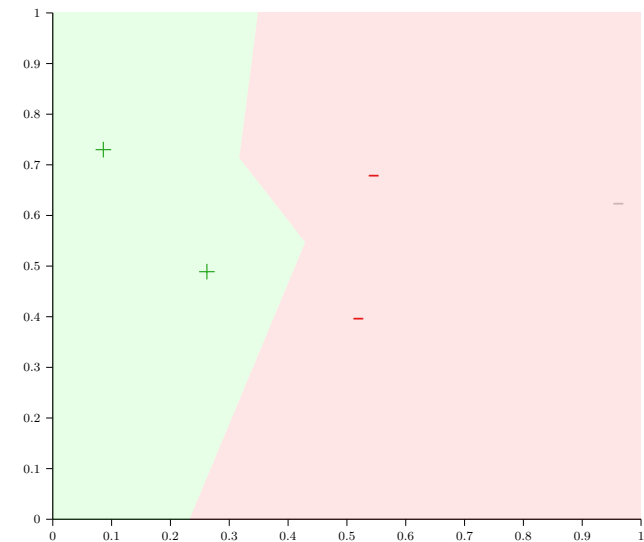
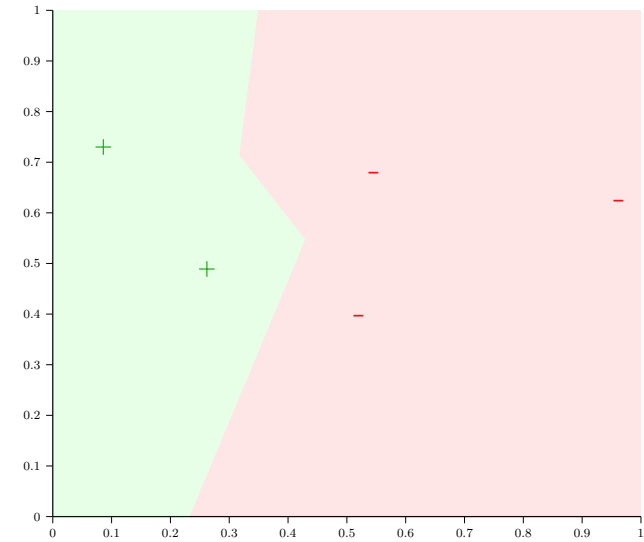


## Computational Cost of $k$ NN

- No training required!
- Memory:  $O(nD)$
- Computing  $g(\vec{x})$ :  $O(nD + n \log(k))$
- Idea: preprocess inputs in order to speed up predictions
  - Reduce the number of inputs held in memory by eliminating redundancies
  - Organize inputs in data structures that make searching for nearest neighbors more efficient

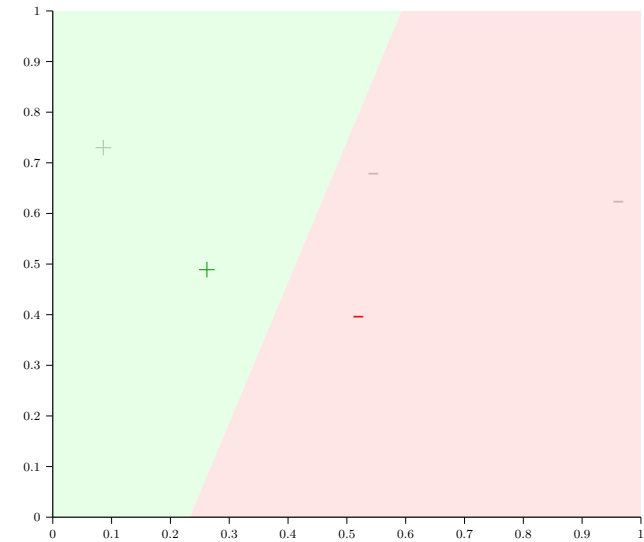
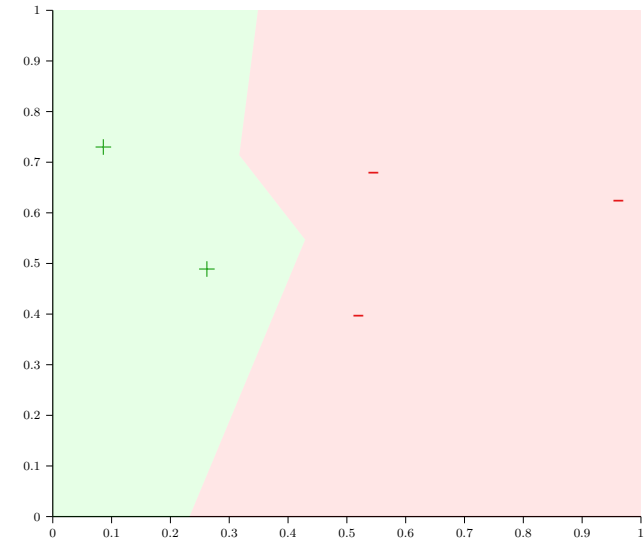
# Data Condensing

- Reduce the number of inputs while maintaining the same decision boundary
- Let  $g_{\mathcal{D}}$  be the  $k$ NN hypothesis when trained on  $\mathcal{D}$
- $S \subseteq \mathcal{D}$  is decision-boundary consistent if:
$$g_S(\vec{x}) = g_{\mathcal{D}}(\vec{x}) \quad \forall \vec{x} \in \mathcal{X}$$
- Decision-boundary consistent subsets are computationally expensive to find



# Data Condensing

- Reduce the number of inputs while maintaining the same predictions on all inputs
- Let  $g_{\mathcal{D}}$  be the  $k$ NN hypothesis when trained on  $\mathcal{D}$
- $S \subseteq \mathcal{D}$  is training-set consistent if:
$$g_S(\vec{x}_i) = g_{\mathcal{D}}(\vec{x}_i) \quad \forall \vec{x}_i \in \mathcal{D}$$
- Training-set consistent is a much weaker constraint than decision-boundary consistent

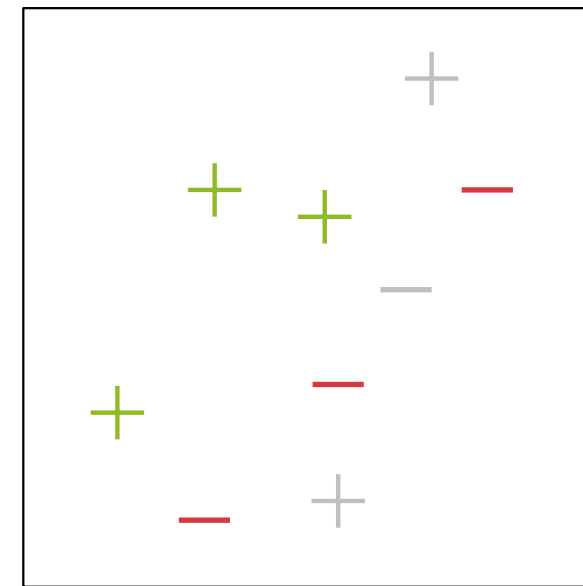
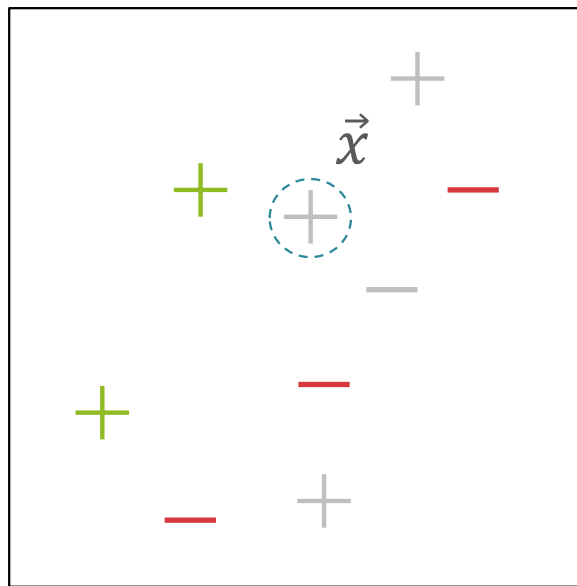


# Condensed Nearest Neighbor (CNN)

- Input:  $\mathcal{D} = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)\}, k$
- Compute  $g_{\mathcal{D}}(\vec{x}_i) \forall \vec{x}_i \in \mathcal{D}$
- Initialize  $S$  to  $k$  random points in  $\mathcal{D}$  and compute  $g_S(\vec{x}_i) \forall \vec{x}_i \in \mathcal{D}$
- While  $\exists \vec{x}_j \in \mathcal{D}$  s.t.  $g_S(\vec{x}_j) \neq g_{\mathcal{D}}(\vec{x}_j)$ 
  - Randomly pick a point  $\vec{x}_j \in \mathcal{D}$  s.t.  $g_S(\vec{x}_j) \neq g_{\mathcal{D}}(\vec{x}_j)$
  - Let  $\vec{x}^*$  be the point closest to  $\vec{x}_j$  that is not already in  $S$  and has label  $y^* = g_{\mathcal{D}}(\vec{x}_j)$
  - Add  $\vec{x}^*$  to  $S$  and recompute  $g_S(\vec{x}_i) \forall \vec{x}_i \in \mathcal{D}$
- Output:  $S$ , a training-set consistent subset of  $\mathcal{D}$

# CNN Example

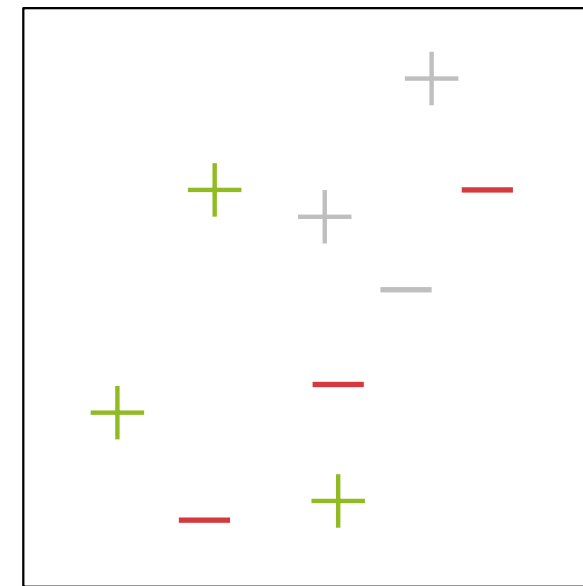
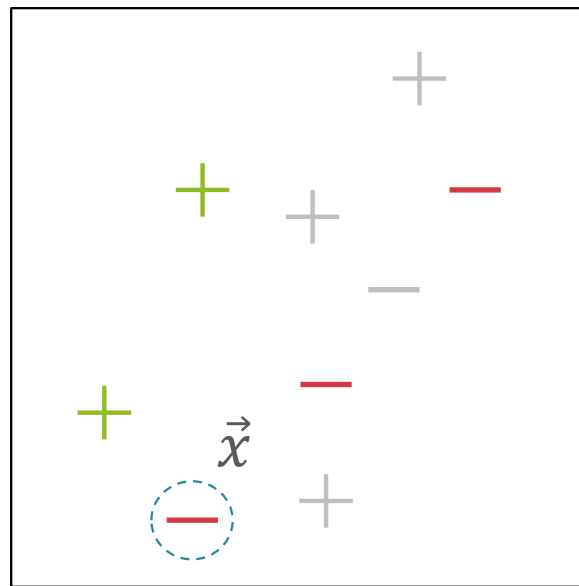
$$k = 3$$



$$g_S(\vec{x}) = -1, g_D(\vec{x}) = +1$$

# CNN Example

$$k = 3$$

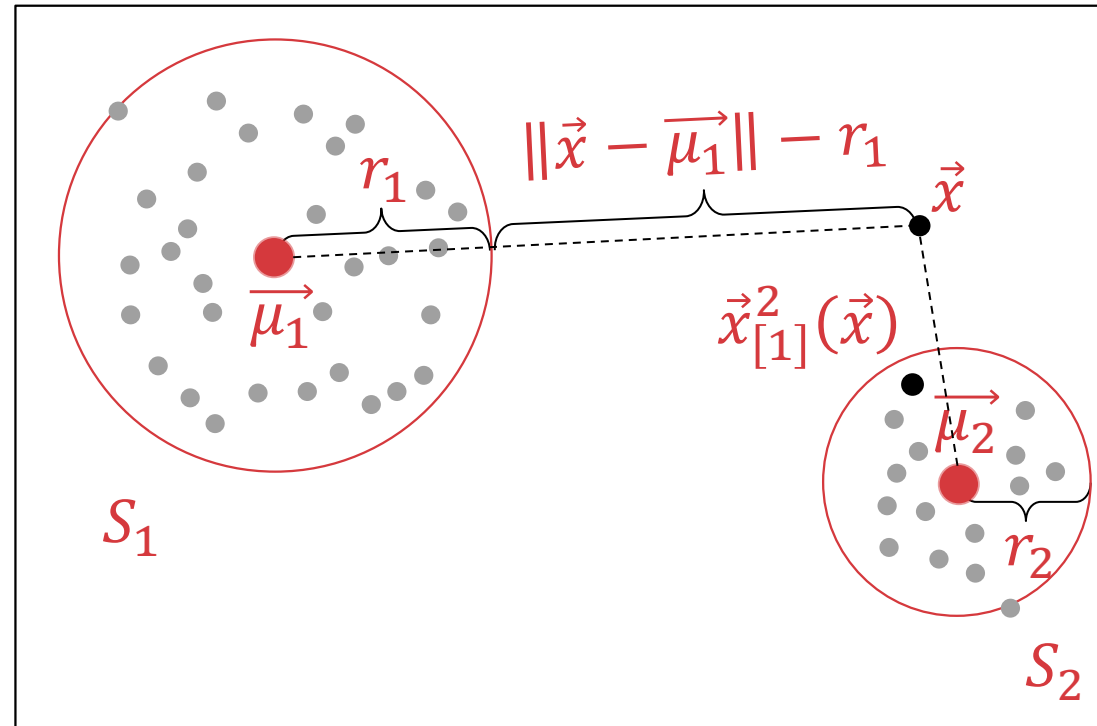


$$g_S(\vec{x}) = -1, g_D(\vec{x}) = +1$$

# Organizing the Inputs

- Intuition: split the inputs into clusters, groups of points that are close to one another but far from other groups.
- If an input point is really close to one group of points and really far from the other groups ...
- ... then we can skip searching through the other groups and just look for nearest neighbors in the close group!
- Questions:
  - What does it mean for a point to be close to a group?
  - How can we split the input into clusters?

# Organizing the Inputs

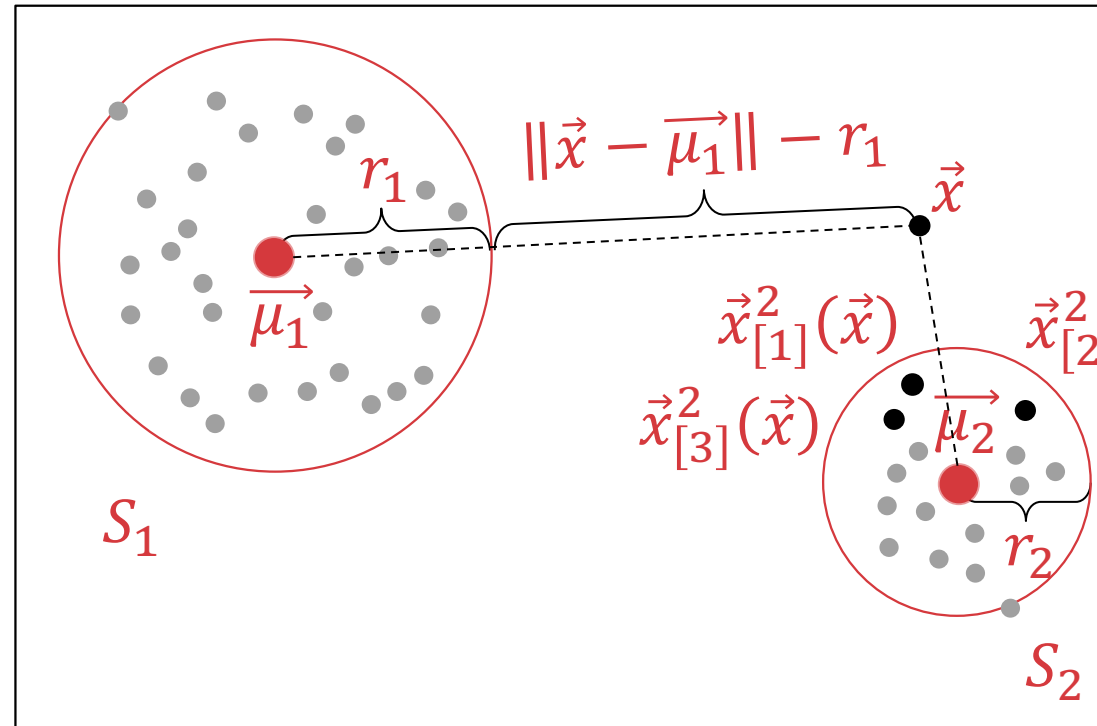


Let  $\vec{x}_{[i]}^j(\vec{x})$  be the  $i^{\text{th}}$  closest point to  $\vec{x}$  in the cluster  $S_j$

- $\|\vec{x} - \vec{\mu}_j\| - r_j \leq \|\vec{x} - \vec{x}_{[1]}^j(\vec{x})\|$
- If  $\exists S_i$  s.t.  $\|\vec{x} - \vec{x}_{[1]}^i(\vec{x})\| \leq \|\vec{x} - \vec{\mu}_j\| - r_j$  ...
- ... then we don't need to search  $S_j$  for the nearest neighbor



# Organizing the Inputs

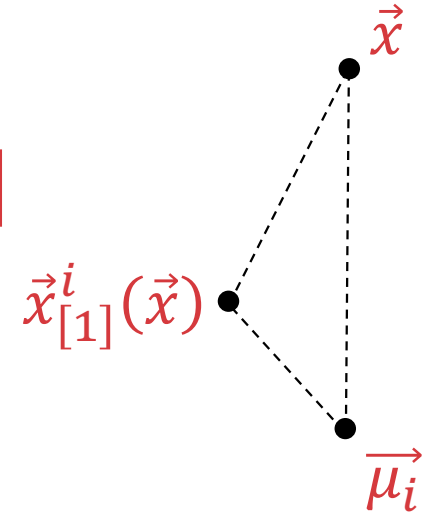


Let  $\vec{x}_{[i]}^j(\vec{x})$  be the  $i^{\text{th}}$  closest point to  $\vec{x}$  in the cluster  $S_j$

- $\|\vec{x} - \vec{\mu}_j\| - r_j \leq \|\vec{x} - \vec{x}_{[1]}^j(\vec{x})\| \leq \dots \leq \|\vec{x} - \vec{x}_{[k]}^j(\vec{x})\|$
- If  $\exists S_i$  s.t.  $\|\vec{x} - \vec{x}_{[k]}^i(\vec{x})\| \leq \|\vec{x} - \vec{\mu}_j\| - r_j \dots$
- ... then we don't need to search  $S_j$  for the  $k$  nearest neighbors

# Clustering Inputs

- We want  $\|\vec{x} - \vec{x}_{[1]}^i(\vec{x})\| \leq \|\vec{x} - \vec{\mu}_j\| - r_j$
- $\|\vec{x} - \vec{x}_{[1]}^i(\vec{x})\| \leq \|\vec{x} - \vec{\mu}_i\| + \|\vec{\mu}_i - \vec{x}_{[1]}^i(\vec{x})\|$   
 $\leq \|\vec{x} - \vec{\mu}_i\| + r_i$



- We want  $\|\vec{x} - \vec{\mu}_i\| + r_i \leq \|\vec{x} - \vec{\mu}_j\| - r_j$
- Suppose  $\vec{x} \approx \vec{\mu}_i$ 
  - $\|\vec{x} - \vec{\mu}_i\| \approx 0$  and  $\|\vec{x} - \vec{\mu}_j\| \approx \|\vec{\mu}_i - \vec{\mu}_j\|$
- We want  $r_i \leq \|\vec{\mu}_i - \vec{\mu}_j\| - r_j \Rightarrow r_i + r_j \leq \|\vec{\mu}_i - \vec{\mu}_j\|$
- We want cluster centers to be far apart and cluster radii to be small

# Clustering Inputs: Example

- Input:  $\mathcal{D} = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\}, C$
- Initialize  $M$  to a random point  $\vec{x} \in \mathcal{D}$  and set  $S_1 = \vec{x}$
- For  $c = 2 \dots C$ 
  - Find the point  $\vec{x}^* \in \mathcal{D} \setminus M$  that is farthest from  $M$ , add  $\vec{x}^*$  to  $M$  and set  $S_c = \vec{x}^*$ :

$$\vec{x}^* = \operatorname{argmax}_{\vec{x}^* \in \mathcal{D} \setminus M} \left( \min_{\vec{x} \in M} \|\vec{x}^* - \vec{x}\| \right)$$

- For  $i = 1 \dots n$ 
  - Find the cluster center  $\vec{x}_c \in M$  closest to  $\vec{x}_i$  and assign  $\vec{x}_i$  to  $S_c$ :

$$\vec{x}_c = \operatorname{argmin}_{\vec{x} \in M} (\|\vec{x}_i - \vec{x}\|)$$

- For  $c = 1 \dots C$ 
  - Update  $S_c$ 's cluster center and compute the radius:

$$\vec{\mu}_c = \frac{1}{|S_c|} \sum_{\vec{x} \in S_c} \vec{x} \quad \text{and} \quad r_c = \max_{\vec{x} \in S_c} \|\vec{\mu}_c - \vec{x}\|$$

- Output: cluster centers  $\mu = \{\vec{\mu}_1, \vec{\mu}_2, \dots, \vec{\mu}_C\}$ , radii  $\vec{r} = \{r_1, r_2, \dots, r_C\}$  and clusters  $S = \{S_1, S_2, \dots, S_C\}$

# Making Predictions: Example

- Input:  $\mathcal{D} = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)\}, C, k, \vec{x}$
- Cluster the input into  $C$  clusters ( $O(CnD)$ )
- Find  $S_c$ , the cluster whose center is closest to  $\vec{x}$
- Find  $\vec{x}_{[k]}^c(\vec{x})$ , the  $k^{\text{th}}$  closest point to  $\vec{x}$  in  $S_c$
- Find  $\{\vec{x}_{[1]}(\vec{x}), \dots, \vec{x}_{[k]}(\vec{x})\}$  the  $k$  closest points to  $\vec{x}$ , ignoring all clusters  $S_{c'}$  s.t.  $\|\vec{x} - \vec{x}_{[k]}^c(\vec{x})\| \leq \|\vec{x} - \vec{\mu}_{c'}\| - r_{c'}$  (hopefully less than  $O(nD + n \log(k))$ )
- Output:  $\{\vec{x}_{[1]}(\vec{x}), \dots, \vec{x}_{[k]}(\vec{x})\}$

# Parametric vs. Non- parametric

- Parametric learning models
  - Hypotheses have a parametrized form
    - Example: linear regression  $g(\vec{x}) = \vec{w}^T \vec{x}$  has  $D + 1$  parameters,  $\{w_0, w_1, \dots, w_D\}$
  - Parameters learned from training data; can discard the training data afterwards
  - Cannot exactly model every target function
- Non-parametric learning models
  - Hypotheses cannot be expressed using a finite number of parameters
    - Example:  $k$ NN, decision trees
  - Training data generally needs to be stored in order to make predictions
  - Can recover any target function given enough data