

Machine Learning from Data

Lecture 17: Spring 2021

Today's Lecture

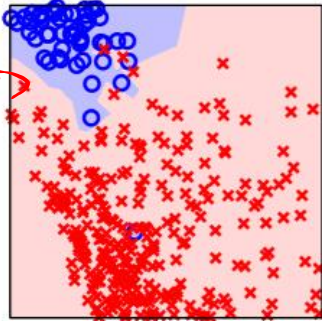
- Memory and Efficiency of Nearest Neighbor

RECAP: Similarity and Nearest Neighbor

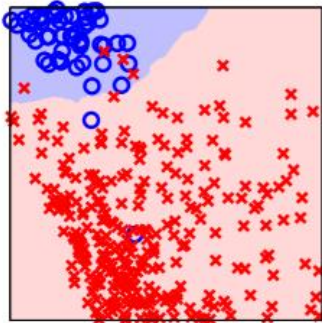
Similarity

$$d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|$$

1-NN rule



21-NN rule



1. Simple. ✓

2. No training. ✓

3. Near optimal E_{out} . ✓

$$k \rightarrow \infty, k/N \rightarrow 0 \Rightarrow E_{\text{out}} \rightarrow E_{\text{out}}^*$$

4. Good ways to choose k :

$$k = 3; k = \lceil \sqrt{N} \rceil; \text{validation/cross validation.}$$

5. Easy to justify classification to customer. ✓

6. Can easily do multi-class.

7. Can easily adapt to regression or logistic regression ✓

$$g(\mathbf{x}) = \frac{1}{k} \sum_{i=1}^k y_{[i]}(\mathbf{x})$$

$$g(\mathbf{x}) = \frac{1}{k} \sum_{i=1}^k \mathbb{I}[y_{[i]}(\mathbf{x}) = +1]$$

8. **Computationally demanding.**

Computational Demands of Nearest Neighbor

①

Memory.

Need to store all the data, $O(Nd)$ memory.

$N = 10^6$, $d = 100$, double precision $\approx 1\text{GB}$

Computational.

②

Finding the nearest neighbor of a test point.

Need to compute distance to every data point, $O(Nd)$.

$N = 10^6$, $d = 100$, 3GHz processor $\approx 3\text{ms}$ (compute $g(\mathbf{x})$)

$\approx 1\text{hr}$ (compute CV error)

$> 1\text{month}$ (choose best k from among 1000 using CV)

Two Basic Approaches

→ Reduce the amount of data.

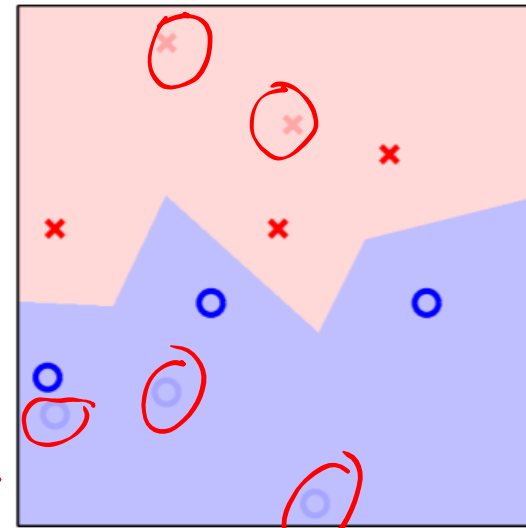
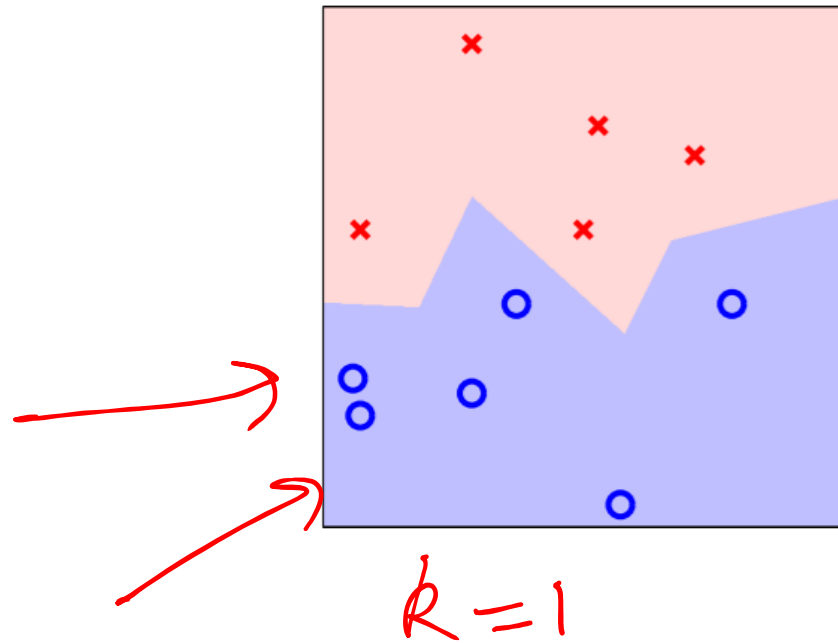
→ Story → Memory

The 5-year old does not remember every horse he has seen, only a few representative horses.

→ Store the data in a specialized data structure.

Ongoing research field to develop geometric data structures to make finding nearest neighbors fast.

Decision Boundary Consistent

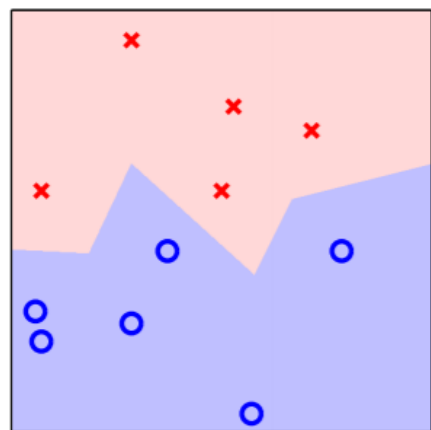


Condensed

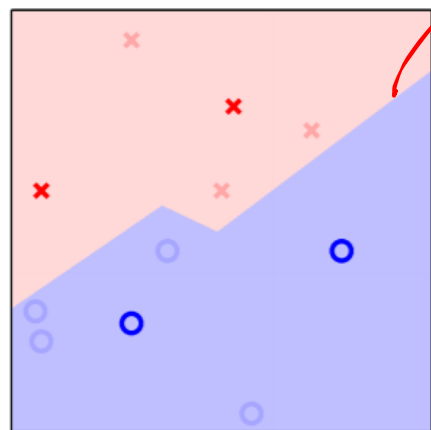
E_{in} change? NO
 E_{out} change? NO

$g(\mathbf{x})$ unchanged

Training Set Consistent



g



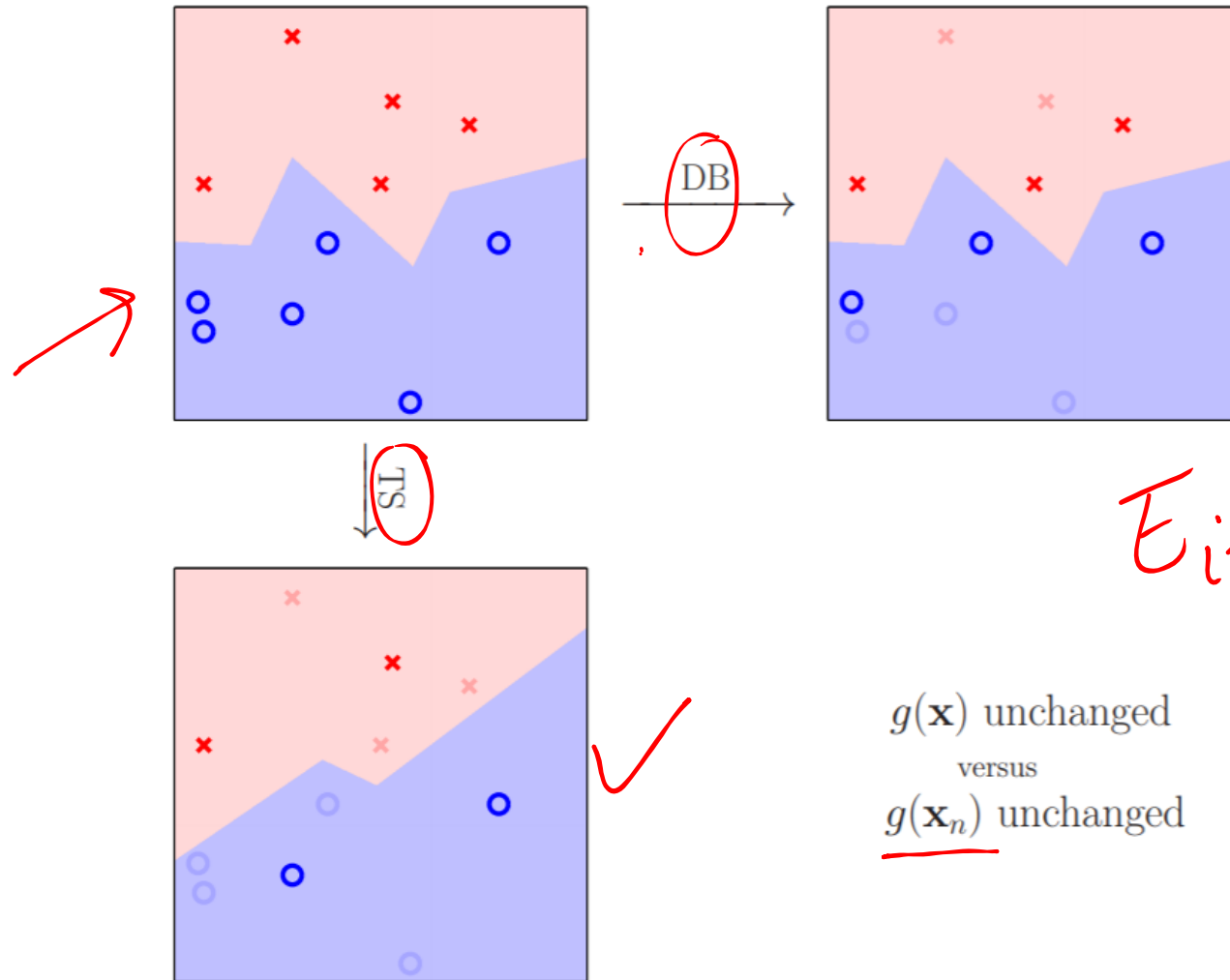
$g_{\hat{}}$ hypothesis

$g(\mathbf{x}_n)$ unchanged

$E_{in} \rightarrow \text{NO.}$

$E_{out} \rightarrow$ hypothesis has changed so E_{out} has changed!

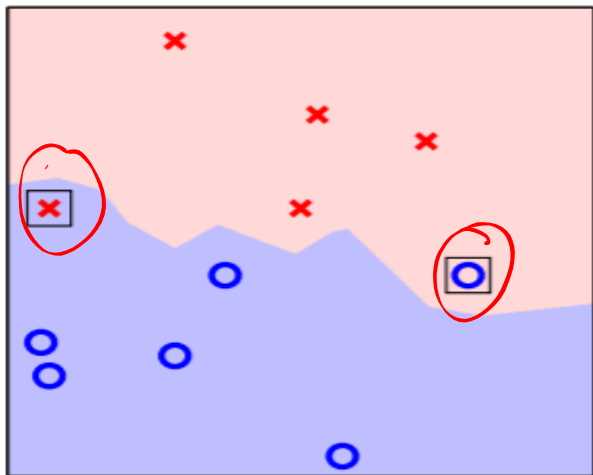
Decision Boundary Vs. Training Set Consistent



$E_{in} = 0$
if $k=3$

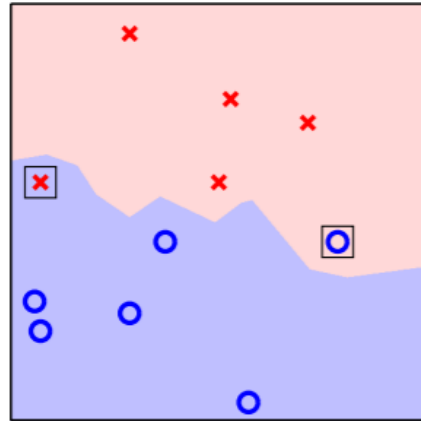
$g(\mathbf{x})$ unchanged
versus
 $g(\mathbf{x}_n)$ unchanged

Consistent Does Not Mean $g(\mathbf{x}_n) = y_n$

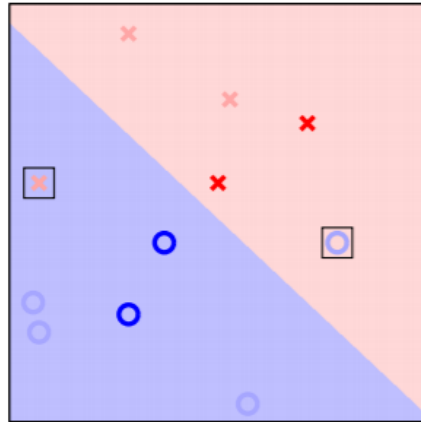


$$k = 3$$

Training Set Consistent $(k = 3)$



$E_{in} \neq 0$



$g(\mathbf{x}_n)$ unchanged

Condensed Nearest Neighbor (CNN)

Take data \rightarrow condense it \rightarrow NN, $k=3$

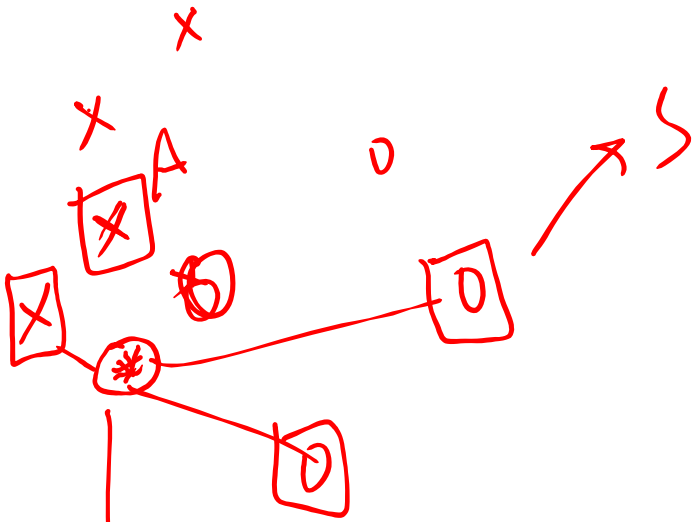
$$\underline{D} = \begin{Bmatrix} x_1 & x_2 & \dots & x_N \\ y_1 & y_2 & \dots & y_N \end{Bmatrix}$$

S

(working set)

$$\begin{matrix} g_D(x) \\ \downarrow \\ \hat{y}_1 \quad \hat{y}_2 \quad \hat{y}_3 \quad \dots \quad \hat{y}_N \\ \hat{y}_i = g_D(x_i) \end{matrix}$$

$$S \rightarrow g_S(x) \\ g_D(x_i) = g_S(x_i)$$



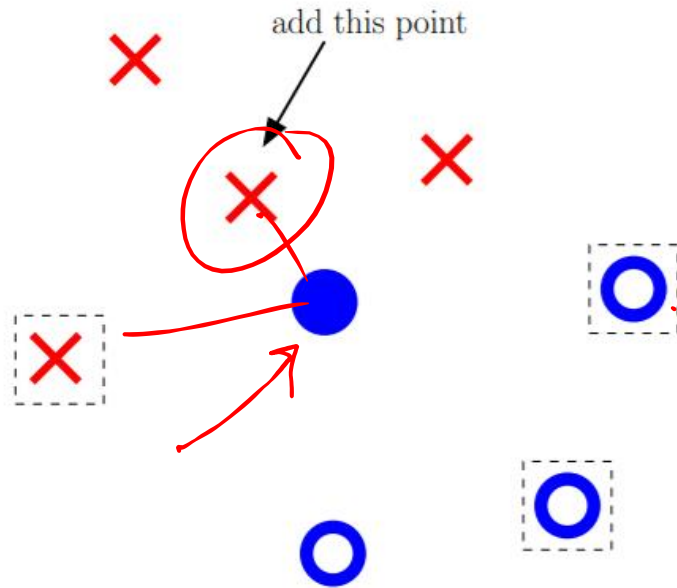
x_* reveals a training set in consistency

$$\begin{matrix} 0 \rightarrow +1 \\ g_S(*) = +1 \\ g_D(*) = -1 \end{matrix}$$

$$S = \emptyset, k=3$$

- ① If S is TS consistent \rightarrow we are done
- ② $x_*(y_*) \rightarrow g_D(x_*) \neq g_S(x_*)$
 \hookrightarrow Add a point to your S
- ③ Add this point
- ④ Repeat these steps \rightarrow TS consistent.

CNN: Condensed Nearest Neighbor



1. Randomly select k data points into \mathcal{S} .
2. Classify all data according to \mathcal{S} .
3. Let \mathbf{x}_* be an inconsistent point and y_* its class w.r.t. \mathcal{D} .
4. Add the closest point to \mathbf{x}_* not in \mathcal{S} that has class y_* .
5. Iterate until \mathcal{S} classifies all points consistently with \mathcal{D} .

Consider the solid blue point:

- i. blue w.r.t. selected points
- ii. red w.r.t. \mathcal{D}

Q1. How do we know such a point exists?
Q2. " " " " that this point is in D?

→ Yes

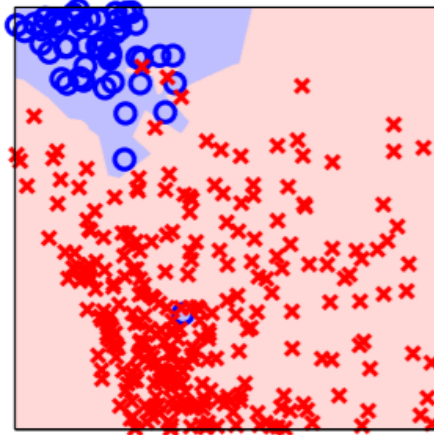
Q3. Algorithm stops or not?
N steps.

Theorem (1) Algorithm is well defined.
2) At most N steps.
3) On termination S must be TS consistent.

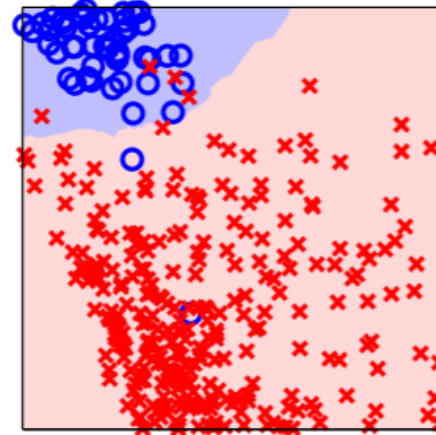
Condensing the Digits Data

Original →

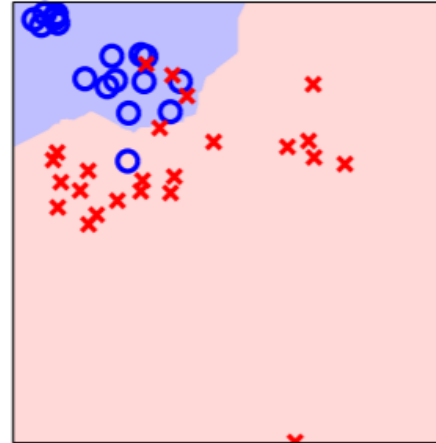
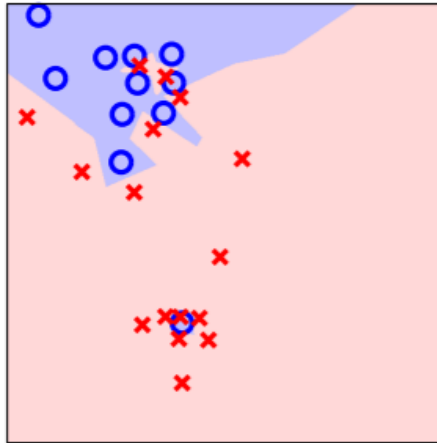
1-NN rule



21-NN rule



CNN →

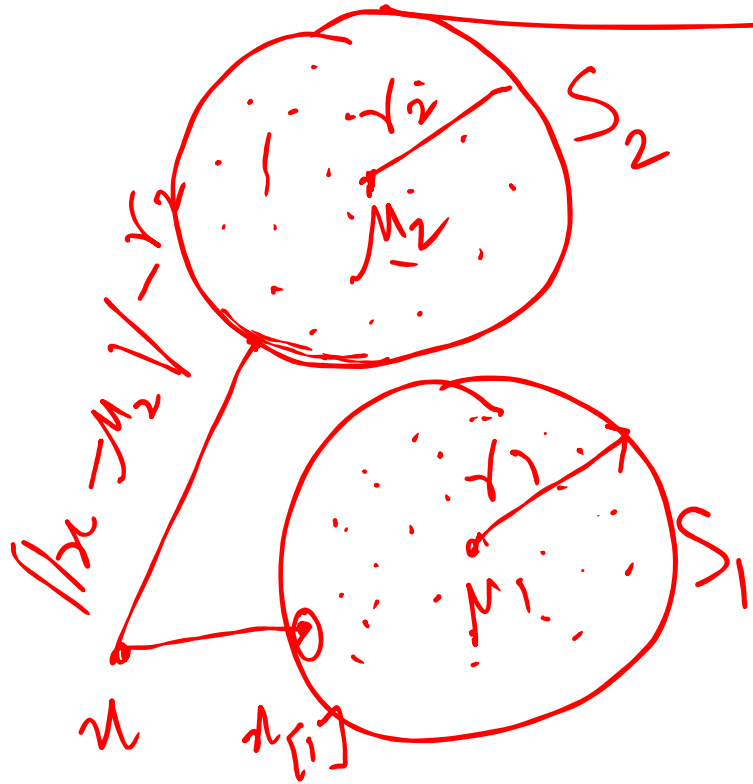


TS
↓
smallest subset?

NP-hard

k, N

Algorithmic efficiency \rightarrow finding NN
(Quickly)
Branch & Bound



1) Choose $S_1 \rightarrow$ Branch Step

2) BOUND CRITERIA

$$\text{if } \underbrace{\|x - \mu_1\| - r_1}_{\text{distance to S1 boundary}} \geq \|x - \hat{x}_{[1]}\|$$

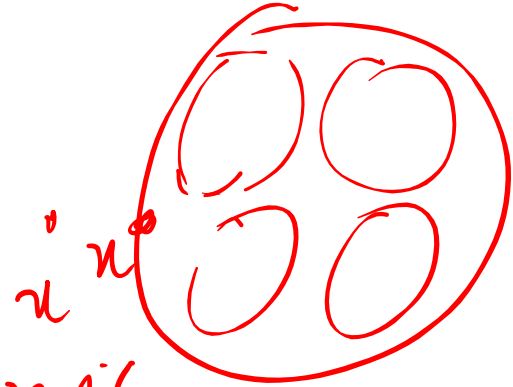
No need to search S_2

Observations

- i) Produces the NN (guaranteed).
- ii) Recursively
- iii) Worst Case $O(N)$

Relan \rightarrow W.h.p \rightarrow logarithmic

E_{out}



Finding the Nearest Neighbor

1. S_1, S_2 are 'clusters' with centers μ_1, μ_2 and radii r_1, r_2 .

2. **[Branch]** Search S_1 first $\rightarrow \hat{\mathbf{x}}_{[1]}$.

3. The distance from \mathbf{x} to any point in S_2 is at least

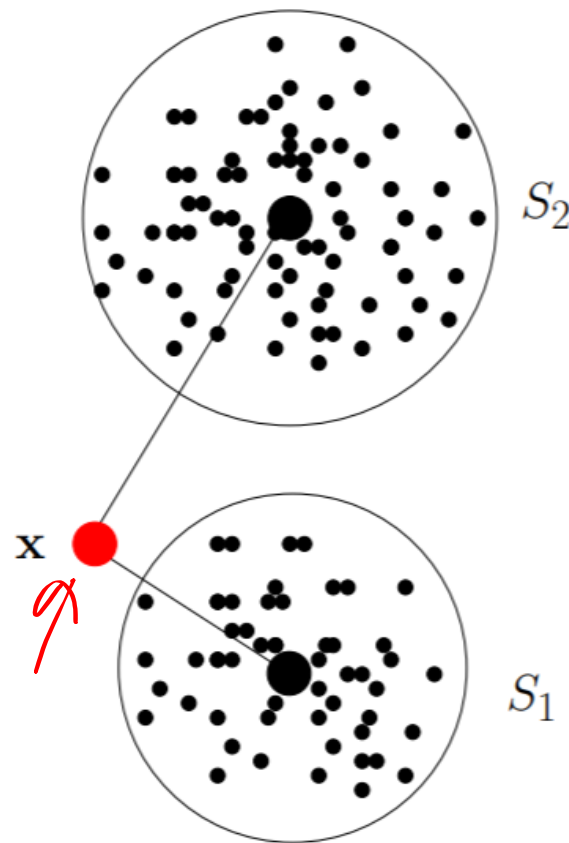
$$\|\mathbf{x} - \mu_2\| - r_2$$

4. **[Bound]** So we are done if

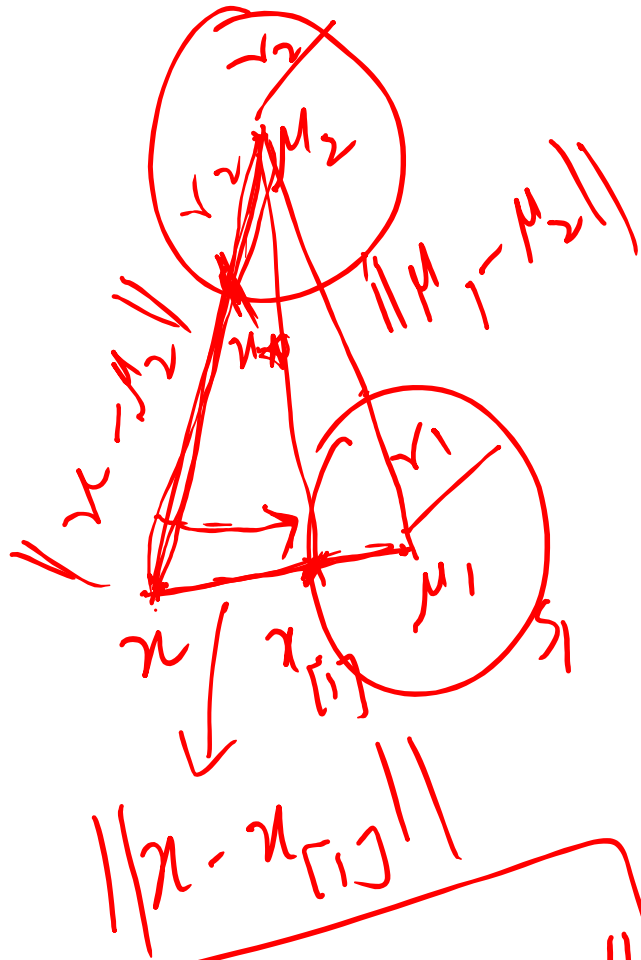
$$\|\mathbf{x} - \hat{\mathbf{x}}_{[1]}\| \leq \|\mathbf{x} - \mu_2\| - r_2$$

A *branch and bound* algorithm

Can be applied *recursively*



When do we gain? $r_1 + r_2 \rightarrow$ Inter, Intra



$$\|x - \hat{x}_{[1]}\| \leq \|x - \mu_1\| + r_1$$

any point in cluster

$$\|x - \mu_1\| + r_1 \leq \|x - \mu_2\| - r_2$$

→ Gain! → Sufficient condition.

$$\|x - \mu_1\| + r_1 \leq \|x - \mu_2\| - r_2$$

$$r_1 + r_2 \leq \|x - \mu_2\| - \|x - \mu_1\|$$

small

$$\approx \|\mu_1 - \mu_2\|$$

$$r_1 + r_2 < \|\mu_1 - \mu_2\|$$

When Does the Bound Hold?

Bound condition: $\|\mathbf{x} - \hat{\mathbf{x}}_{[1]}\| \leq \|\mathbf{x} - \boldsymbol{\mu}_2\| - r_2$.

$$\|\mathbf{x} - \hat{\mathbf{x}}_{[1]}\| \leq \|\mathbf{x} - \boldsymbol{\mu}_1\| + r_1$$

So, it suffices that

$$r_1 + r_2 \leq \|\mathbf{x} - \boldsymbol{\mu}_2\| - \|\mathbf{x} - \boldsymbol{\mu}_1\|.$$

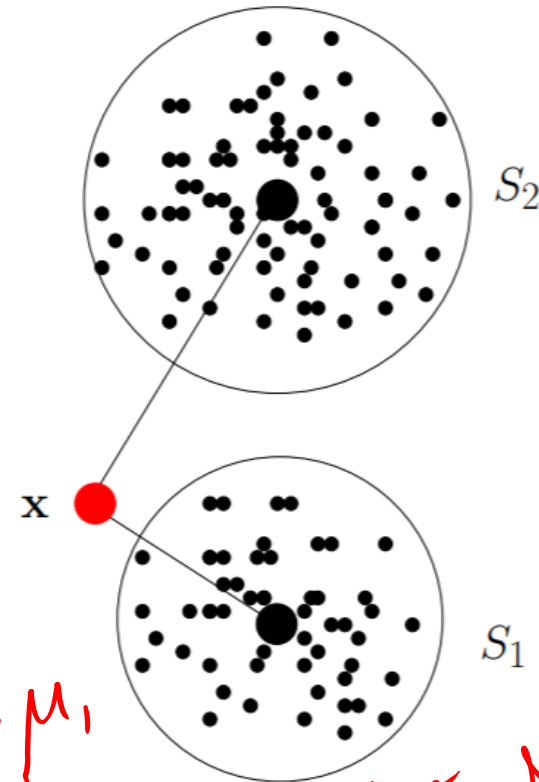
$\|\mathbf{x} - \boldsymbol{\mu}_1\| \approx 0$ means $\|\mathbf{x} - \boldsymbol{\mu}_2\| \approx \|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|$.

It suffices that

$$r_1 + r_2 \leq \|\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1\|.$$

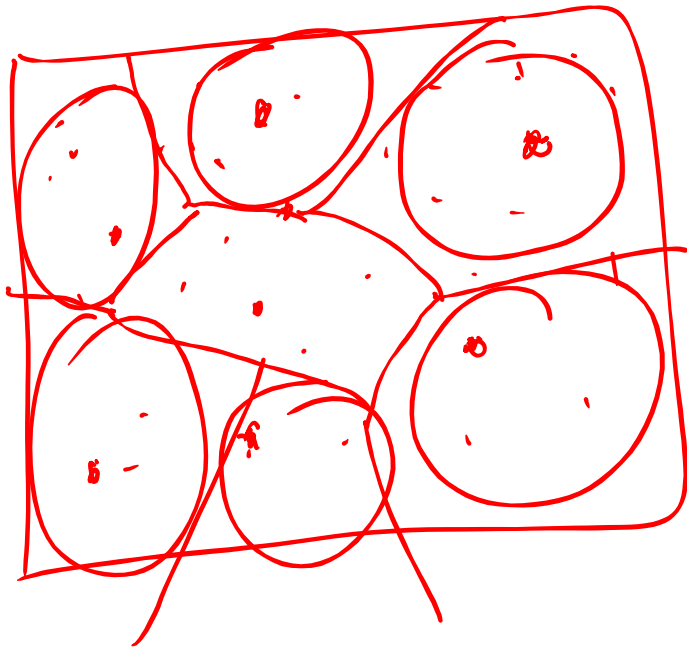
$$\hookrightarrow r_1 + r_2 \ll \mu_2 - \mu_1$$

within cluster spread should be less than between cluster spread



tight & well separated.

Lloyd's Algo ✓



- 1) Pick center randomly.
- 2) Pick next as far as possible, Voronoi.

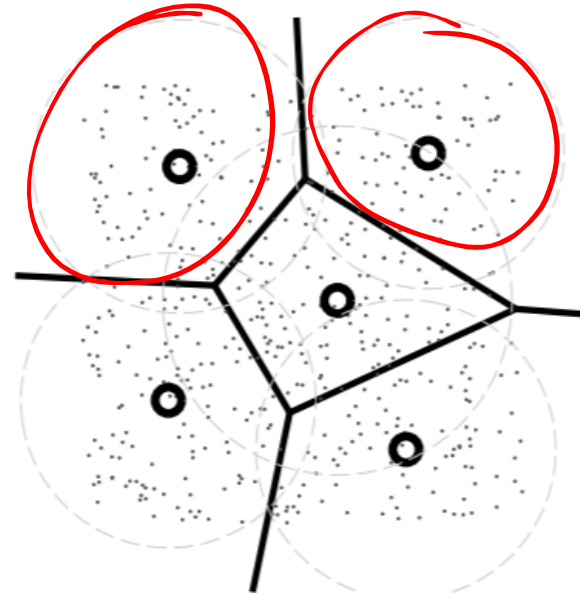
- 3) Update the center to the actual center

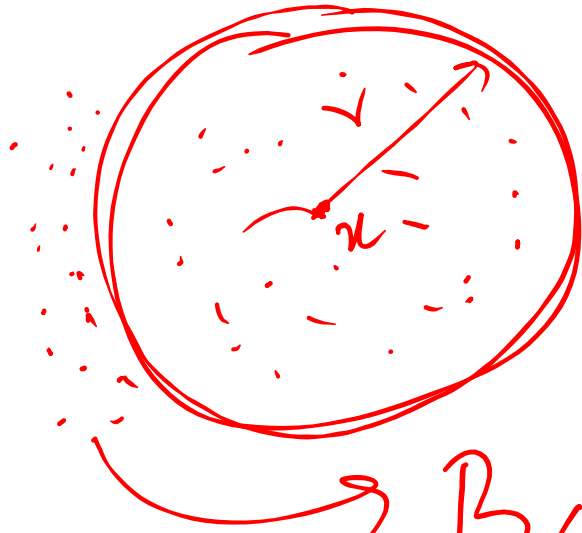
$$r_1 + r_2 < < \|\mu_1 - \mu_2\|$$

Finding Clusters – Lloyd's Algorithm

1. Pick well separated centers for each cluster.
2. Compute Voronoi regions as the clusters.
3. Update the Centers.
4. Update the Voronoi regions.
5. Compute centers and radii:

$$\mu_j = \frac{1}{|S_j|} \sum_{\mathbf{x}_n \in S_j} \mathbf{x}_n; \quad r_j = \max_{\mathbf{x}_n \in S_j} \|\mathbf{x}_n - \mu_j\|.$$





Bad \rightarrow Weights

Radial Basis functions.

Thanks!