# Neighborhood Graphs
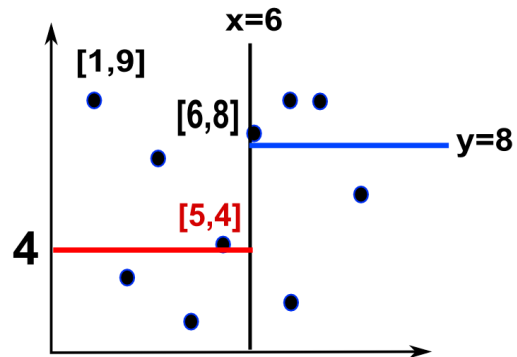
P

ε

x=6
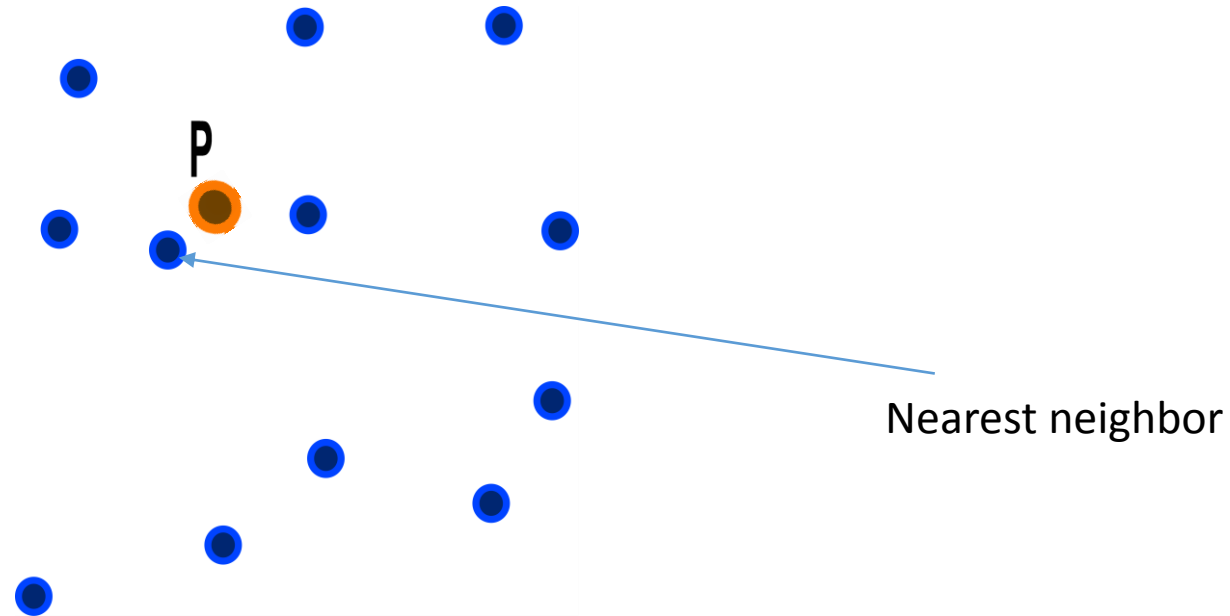
[1,9]

[6,8]

y=8

[5,4]

4

Mustafa Hajij

# Finding Nearest neighbors

Given: a set $X$ of n points in $R^d$ and distance function $d$. We are looking for efficient algorithms to for the following :

- Nearest Neighbor: for any query $p$, returns a point $x \in X$ minimizing $d(p, x)$.

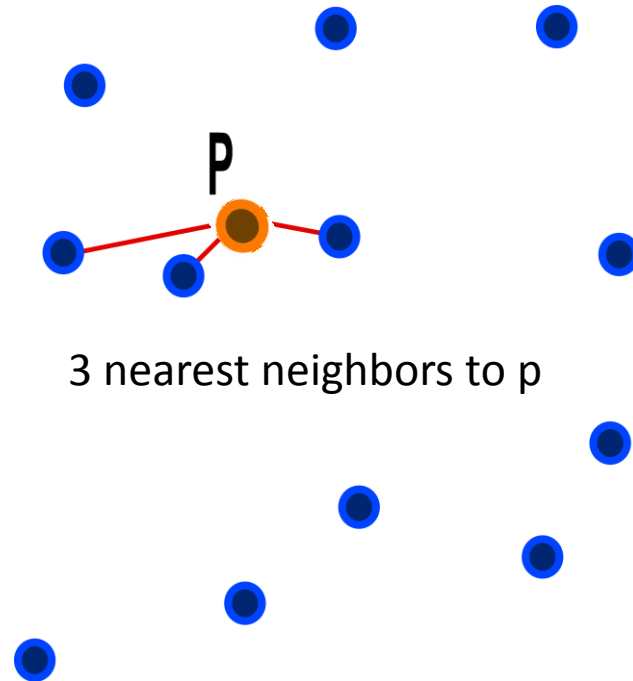- $\varepsilon$ -Near Neighbor: for any query $p$, returns all points $x \in X$ s.t. $d(p, x) \leq \varepsilon$ (if they exist)

# Nearest neighbor

Nearest Neighbor: for any query $p$, returns a point $x \in X$ minimizing $d(p, x)$



P

Nearest neighbor

Write a brute-force code for this algorithm

# K-Nearest neighbors
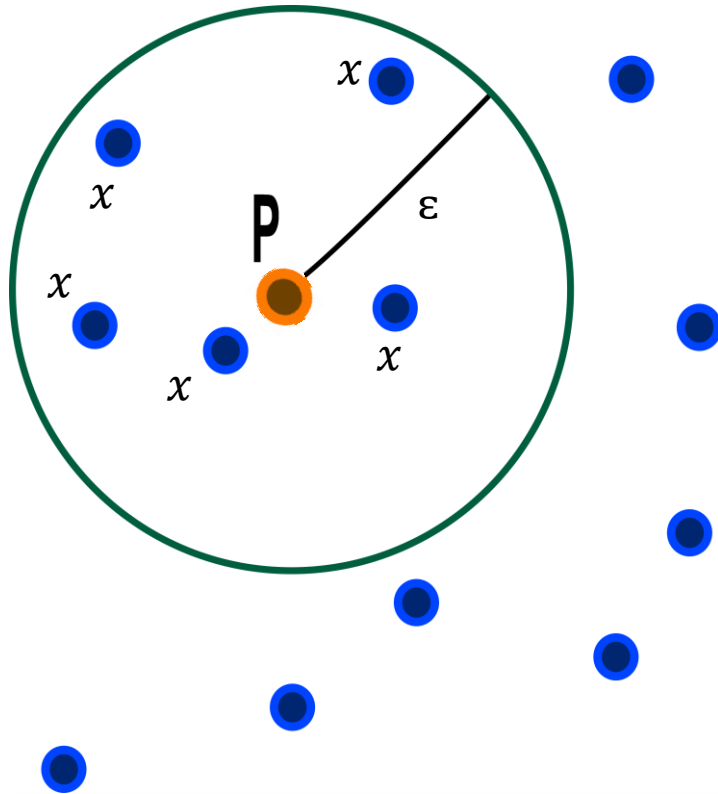
P

3 nearest neighbors to p

Application : if $p$ is an document, then the
The set of k nearest neighbors of $p$ represent
The k most similar articles to $p$

Write a brute-force code for this algorithm

Sklearn example

# Finding ε −Nearest neighbors

Finding all points x that that satisfy $d(p, x) \leq \varepsilon$



Write a brute-force code for this algorithm
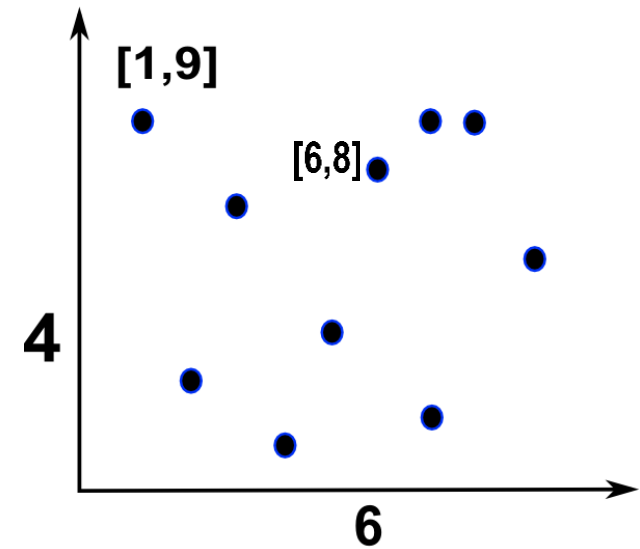
Sklearn example

# Algorithms?

# 2d tree

Suppose that we are given a data such as the data X=[[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]].

KD-Tree organizes a data set such as the data X as a tree for the propose of fast data inquiry.
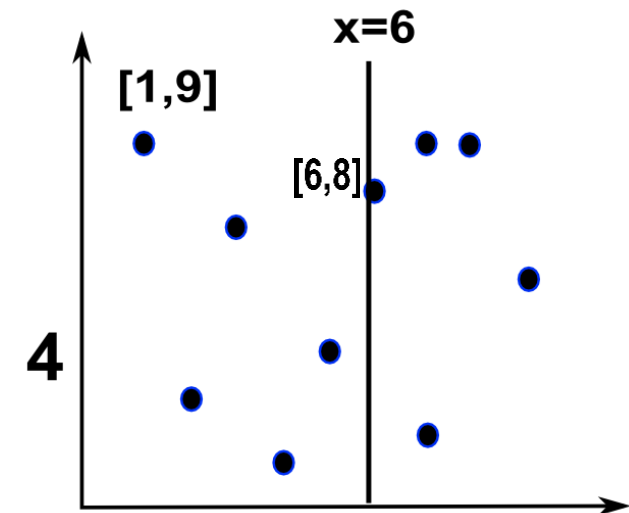
# 2d-tree: Example

2D-Tree organizes a data set such as the data X as a tree for the propose of fast data inquiry.

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

# 2d-tree: Example

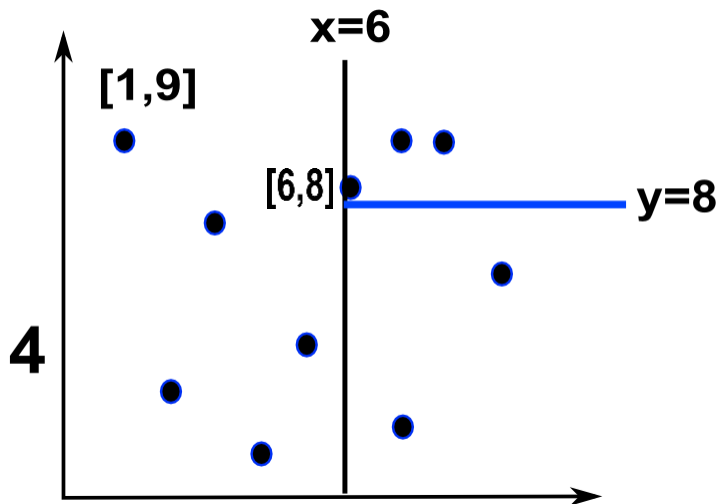[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

↓

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

x=6

[1,9]

[6,8]

4

# 2d-tree: Example

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=8
[6,8]

[7,2],[9,6]

[6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

4

# 2d-tree: Example

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]
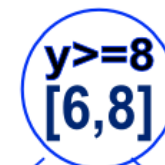
y>=8
[6,8]

[2,3],[4,1]      [1,9],[3,7],[5,4]

[7,2],[9,6]      [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

[5,4]

4

# 2d-tree: Example

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]



x>=6
[6,8]

[1,9],[2,3],[4,1],
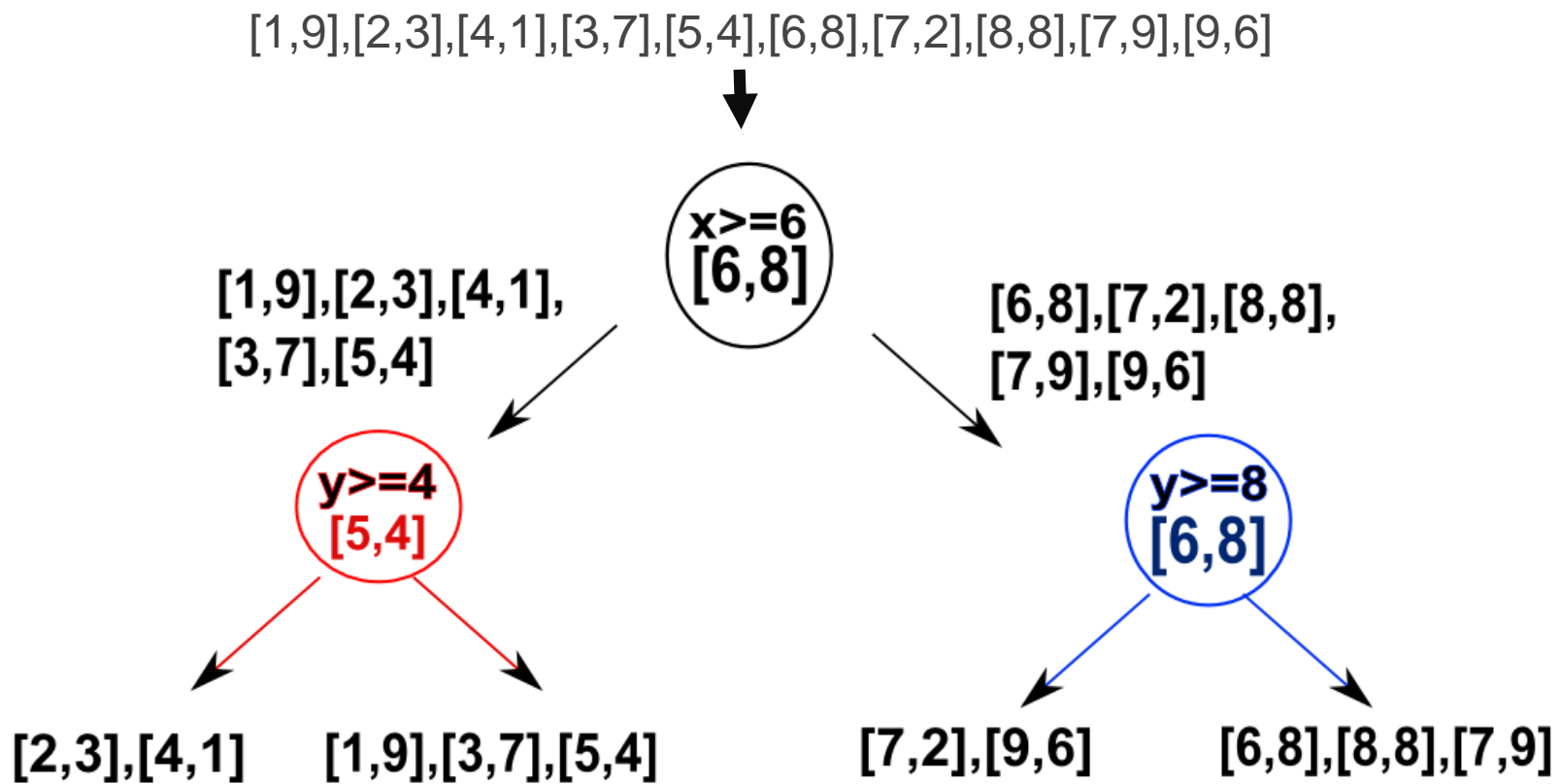[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]     [1,9],[3,7],[5,4]

[7,2],[9,6]     [6,8],[8,8],[7,9]

Points here satisfy all the conditions as we
travel the tree from top to the leaf

x=6
[1,9]
[6,8]
y=8
[5,4]
4

sklearn

# 2d-tree: Example

What information do we store at the each node in the tree ?

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]     [1,9],[3,7],[5,4]

[7,2],[9,6]     [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

[5,4]

4

sklearn

# 2d-tree: Example

What information do we store at the each node in the tree ?
At each node we store the following data :
1- The dimension we split on. In the 2d example below we store if it is the first coordinate or the second one.

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]    [1,9],[3,7],[5,4]

[7,2],[9,6]    [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

[5,4]

4

# 2d-tree: Example

What information do we store at the each node in the tree ?

At each node we store the following data :

1- The dimension we split on. In the 2d example below we store if it is the first coordinate or the second one.
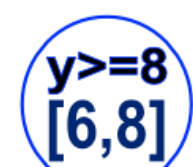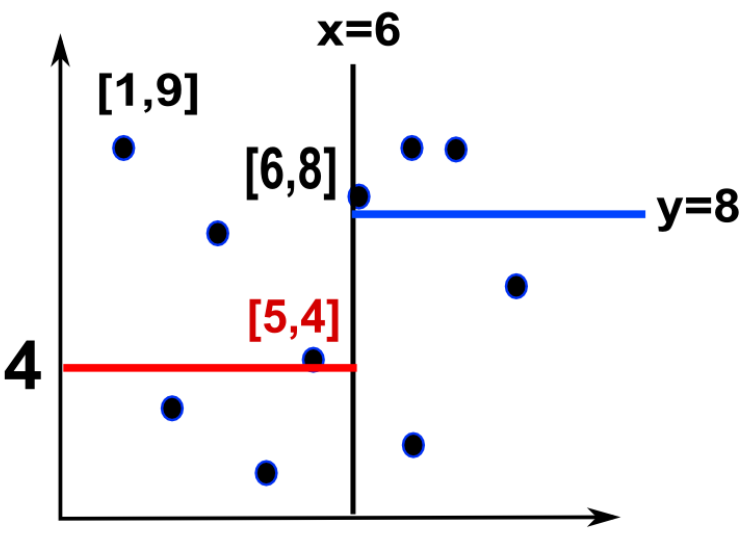
2- the split value (the value of the coordinate)

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

# 2d-tree: Example

What information do we store at the each node in the tree ?
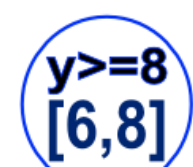
At each node we store the following data :

1- The dimension we split on. In the 2d example below we store if it is the first coordinate or the second one.

2- the split value (the value of the coordinate)

3-the smallest bounding box that contains all points within that node

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]



Example of the bounding box associated to a tree node

# 2d-tree: Example

What information do we store at the each node in the tree ?

At each node we store the following data :

1- The dimension we split on. In the 2d example below we store if it is the first coordinate or the second one.

2- the split value (the value of the coordinate)

3-the smallest bounding box that contains all points within that node

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

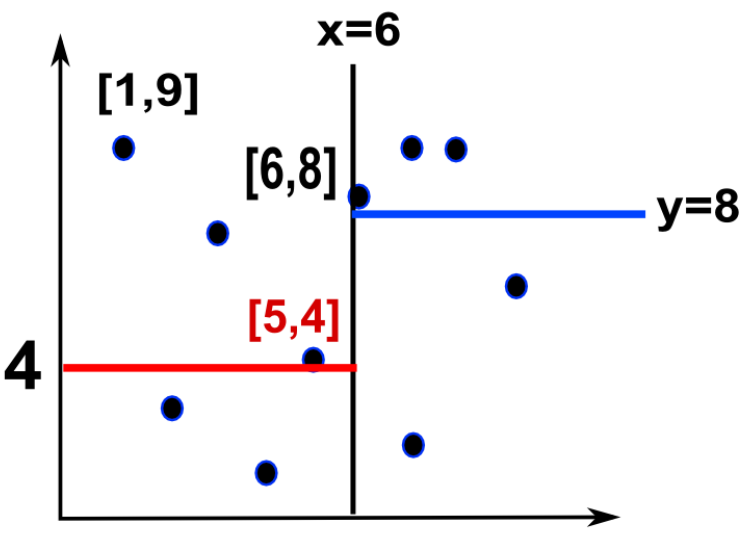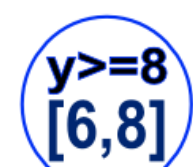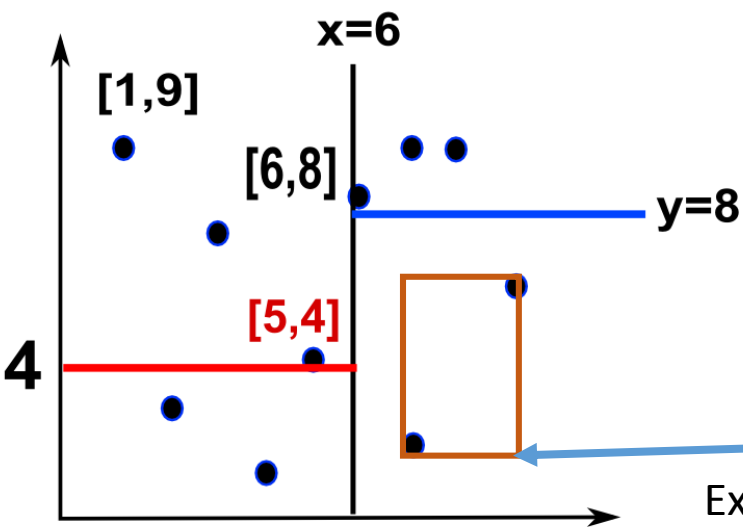y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]    [1,9],[3,7],[5,4]

[7,2],[9,6]    [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

[5,4]

4

Example of the bounding box associated to a tree node

sklearn

# 2d-tree: Example

What information do we store at the each node in the tree ?
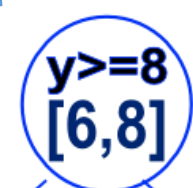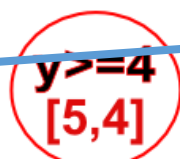
At each node we store the following data :

1- The dimension we split on. In the 2d example below we store if it is the first coordinate or the second one.

2- the split value (the value of the coordinate)

3-the smallest bounding box that contains
   all points within that node

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

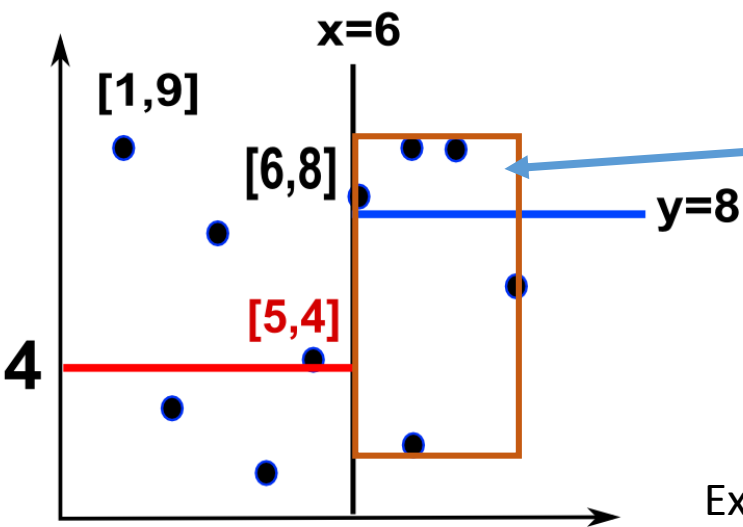[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

x=6

[1,9]

[6,8]

y=8

[5,4]

4

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]    [1,9],[3,7],[5,4]

[7,2],[9,6]    [6,8],[8,8],[7,9]

Example of the bounding box associated to a tree node

sklearn

# 2d-tree: Example

Now suppose that we want to find the nearest neighbor from the data X to a given point .
In that case we navigate in the KD-tree to the region that "hopefully" contains the nearest neighbor

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]     [1,9],[3,7],[5,4]

[7,2],[9,6]     [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

[5,4]

4

# 2d-tree: Example

Consider the point p=[7.7.5]. To find the closest point from the set X to p. We navigate down the tree following the conditions on the nodes

# 2d-tree: Example

Consider the point p=[7.7.5]. To find the closest point from the set X to p. We navigate down the tree following the conditions on the nodes

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]
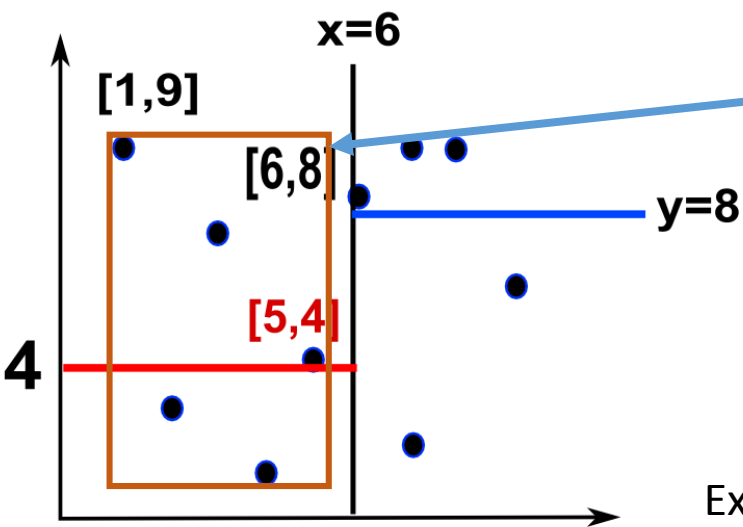
[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]    [1,9],[3,7],[5,4]

[7,2],[9,6]    [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

[5,4]

y=8

[7,7.5]

4

sklearn

# 2d-tree: Example

Consider the point p=[7.7.5]. To find the closest point from the set X to p. We navigate down the tree following the conditions on the nodes

When we arrive to a leaf we check the distance from the point p to every other point in the leaf.

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]      [1,9],[3,7],[5,4]

[7,2],[9,6]      [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

[7,7.5]

[5,4]

4

sklearn

# 2d-tree: Example

Consider the point p=[7.7.5]. To find the closest point from the set X to p. We navigate down the tree following the conditions on the nodes

When we arrive to a leaf we check the distance from the point p to every other point in the leaf. Note that [1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]
This is the leaf that contains the point.



[1,9],[2,3],[4,1],[3,7],[5,4]

x>=6
[6,8]

[6,8],[7,2],[8,8],[7,9],[9,6]

x=6

[1,9]

[6,8]

y=8
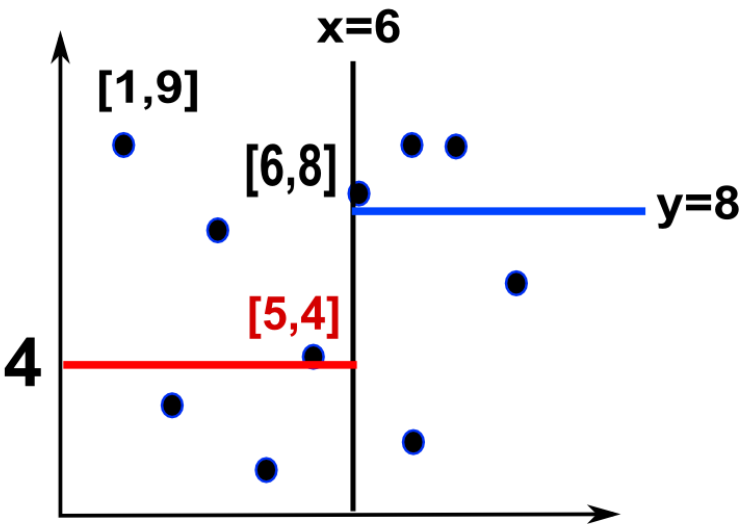
y>=4
[5,4]

y>=8
[6,8]

[7,7.5]
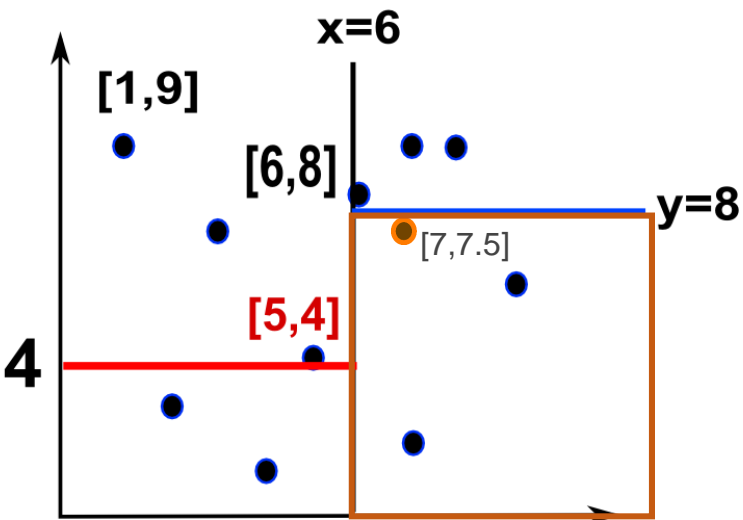
[5,4]

4

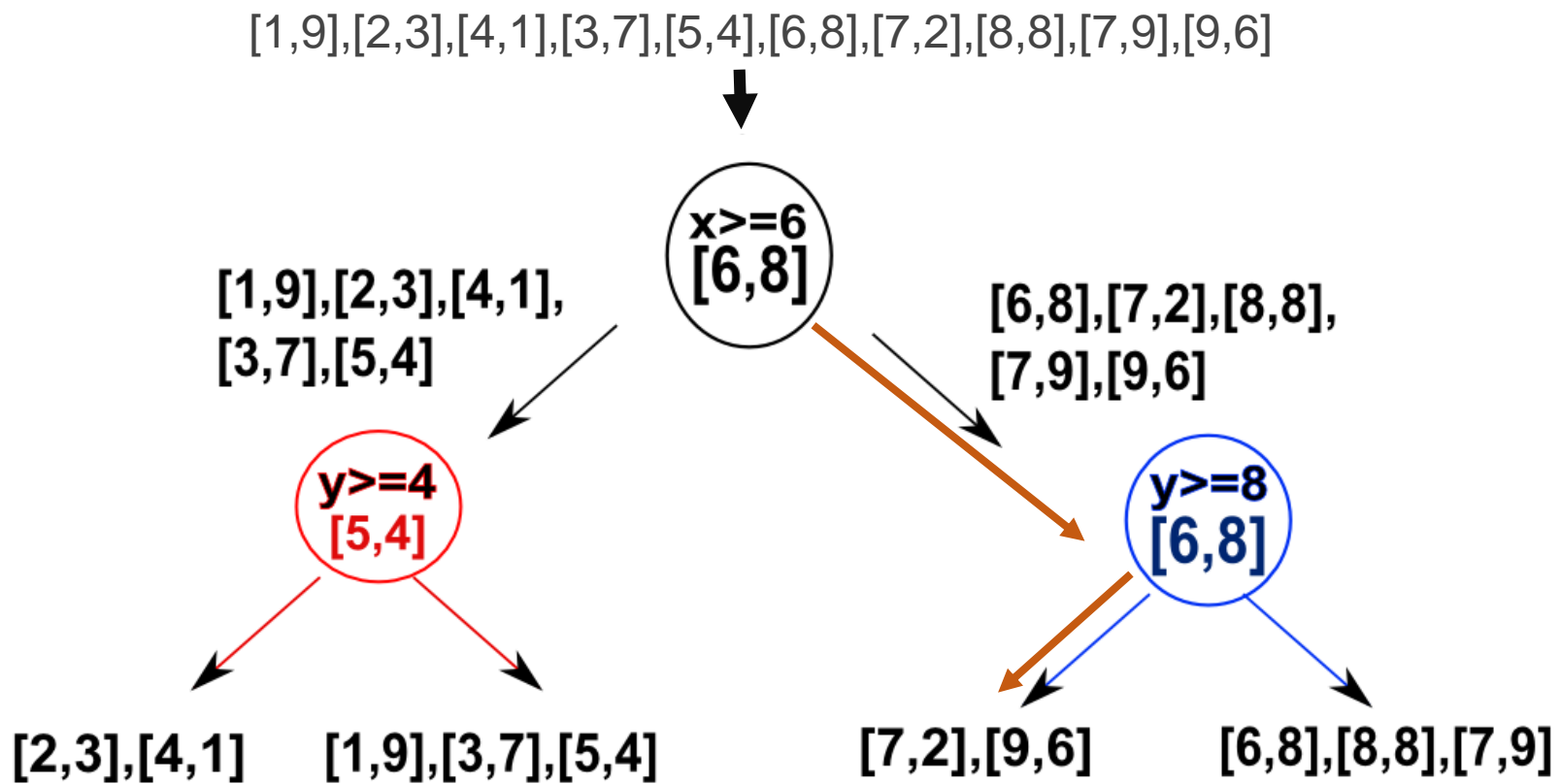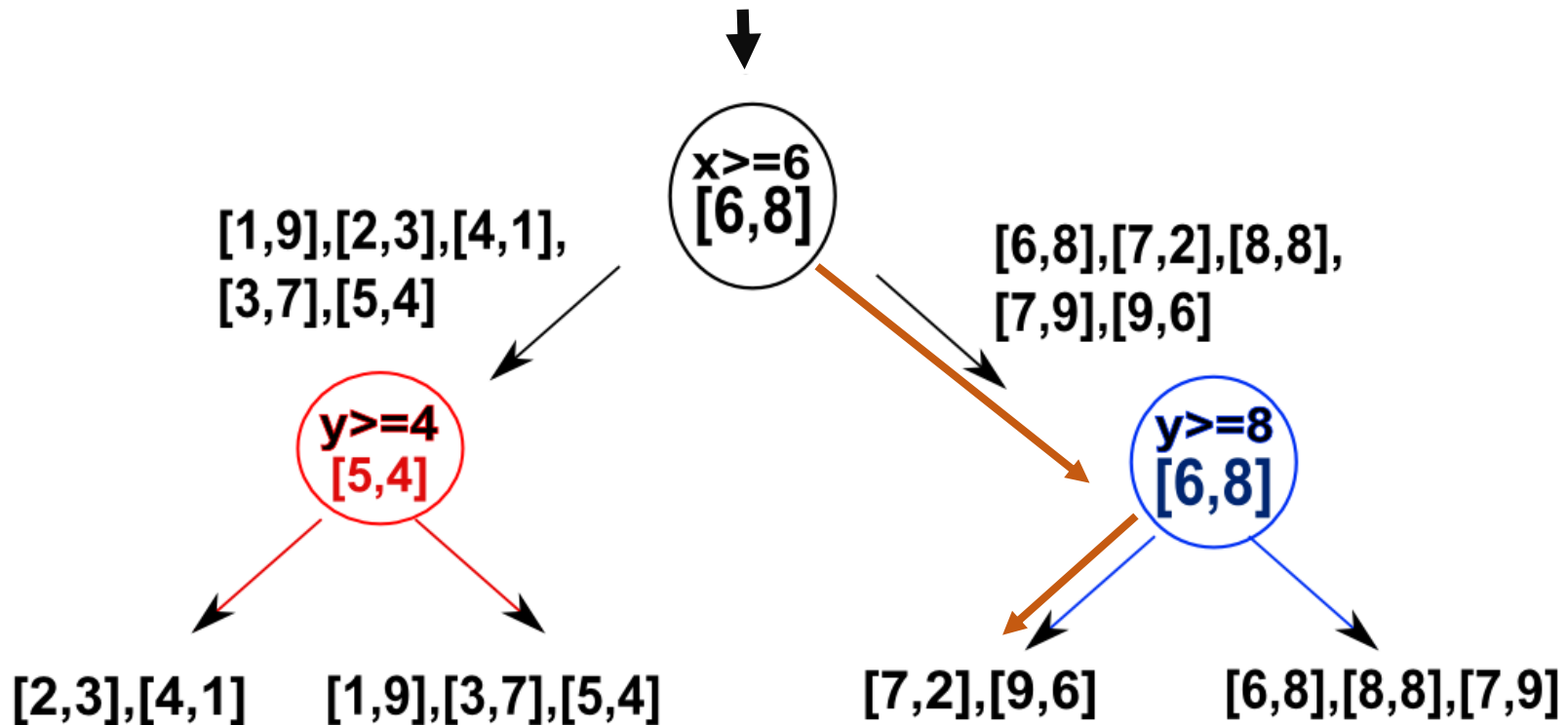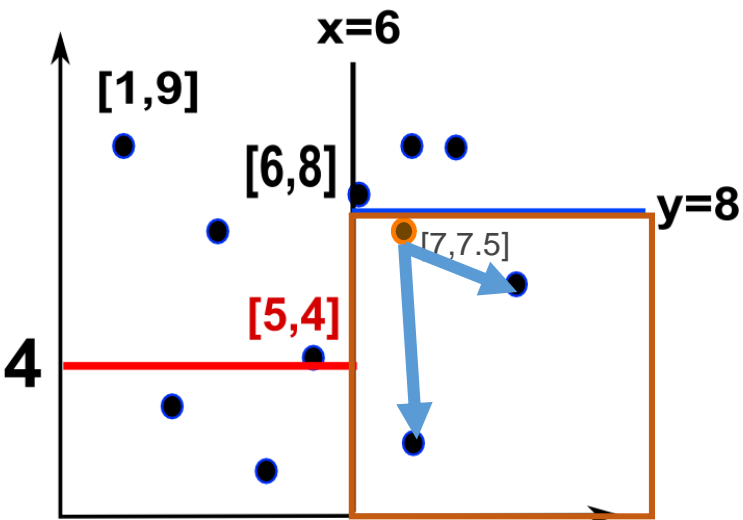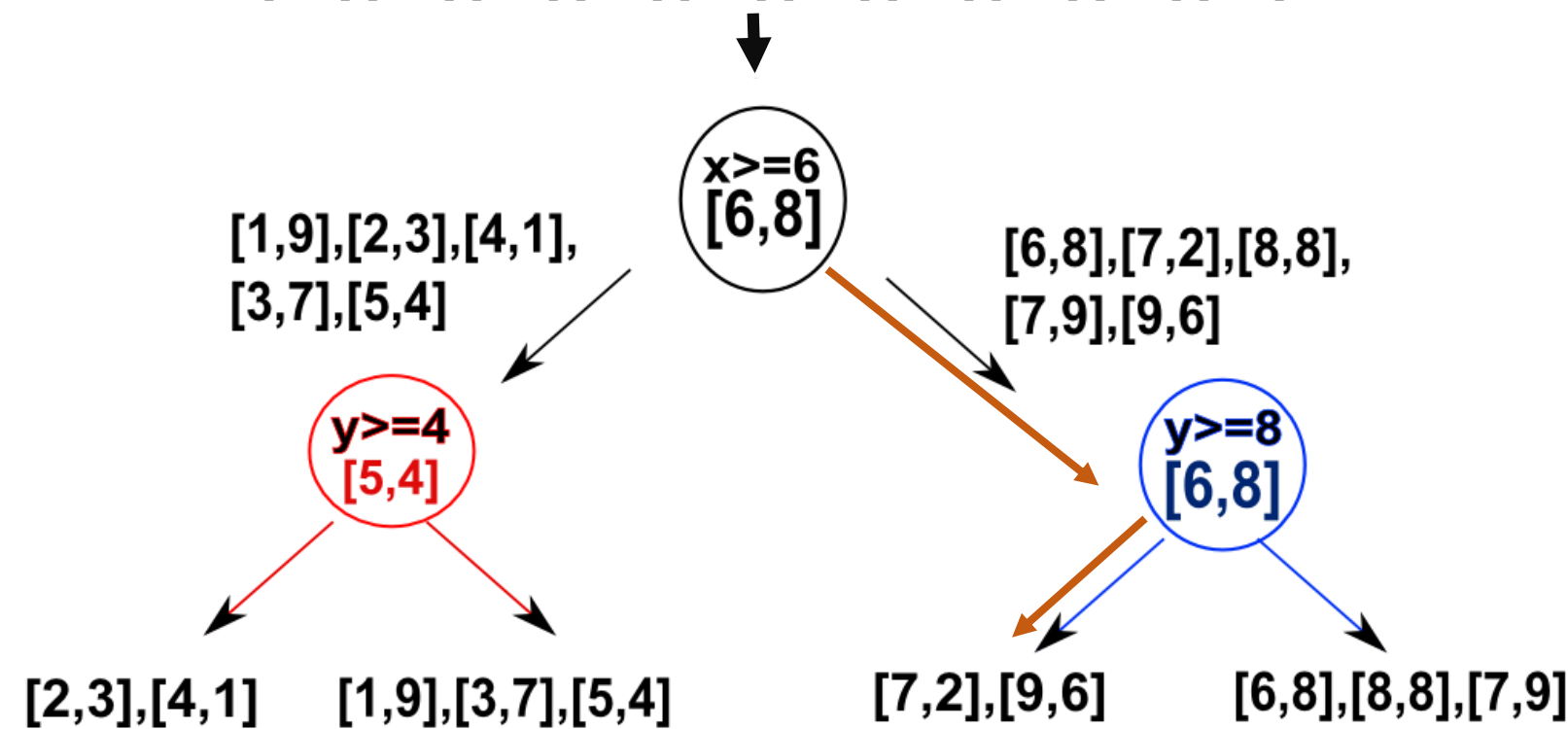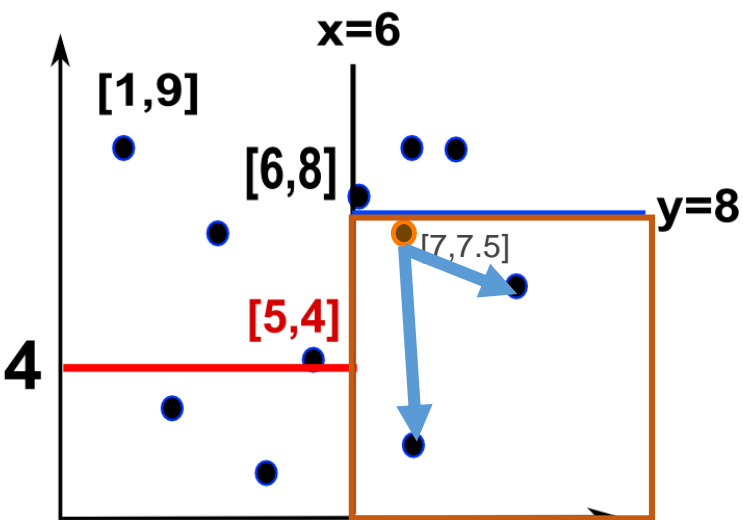[2,3],[4,1]    [1,9],[3,7],[5,4]

[7,2],[9,6]    [6,8],[8,8],[7,9]

# 2d-tree: Example

Consider the point *p*=[7.7.5]. To find the closest point from the set X to p. We navigate down the tree following the conditions on the nodes

When we arrive to a leaf we check the distance from the point *p* to every other point in the leaf. Note that   [1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]
This is the leaf that contains the point.

We then choose the point that is closest within this bucket and record its distance to the point *p*.

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]     [1,9],[3,7],[5,4]

[7,2],[9,6]     [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

[7,7.5]

[5,4]

4

# 2d-tree: Example

What we do from this point is backtracking and try other branches of the tree.

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]



x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]     [1,9],[3,7],[5,4]

[7,2],[9,6]     [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

[7,7.5]

y=8

[5,4]

4

sklearn

# 2d-tree: Example

What we do from this point is backtracking and try other branches of the tree.

At each leaf we visit, we again measure the distance to every node within that leaf.

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]      [1,9],[3,7],[5,4]      [7,2],[9,6]      [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

[5,4]

4

# 2d-tree: Example

What we do from this point is backtracking and try other branches of the tree.

At each leaf we visit, we again measure the distance to every node within that leaf.
If we find a shortest distance we record this distance along with the point that gave us that.

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]        [1,9],[3,7],[5,4]

[7,2],[9,6]        [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

[5,4]

4

# 2d-tree: Example

What we do from this point is backtracking and try other branches of the tree.

Now when we go to check the other branch we notice that the distance to the bounding box is greater than the distance to the nearest neighbor found so far, so we ignore that branch from the search. Since there are no more branches to search in this case we return the point that we have found.

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]

x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]      [1,9],[3,7],[5,4]

[7,2],[9,6]      [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

[5,4]

4

# 2d-tree: Example

What we do from this point is backtracking and try other branches of the tree.

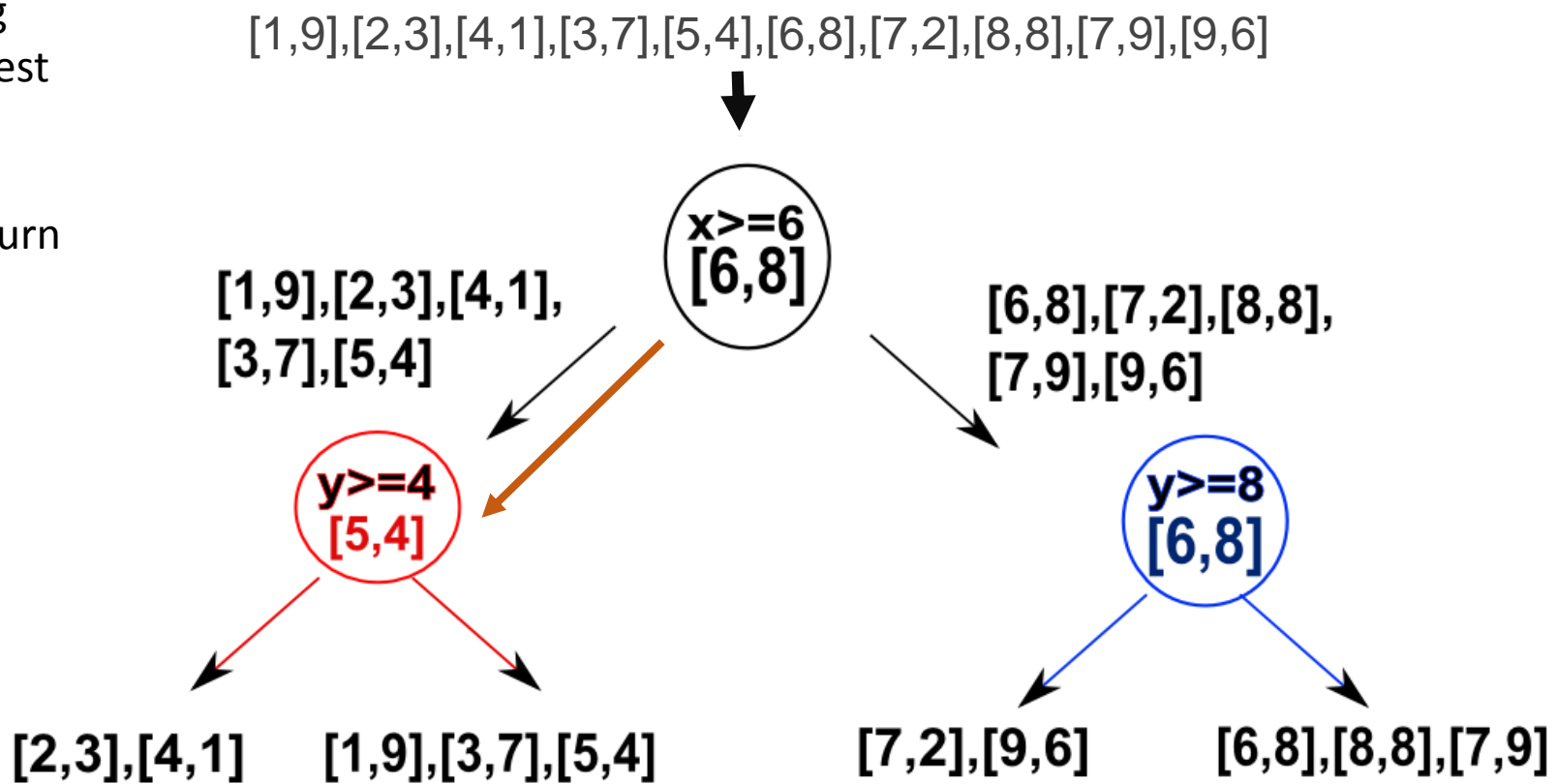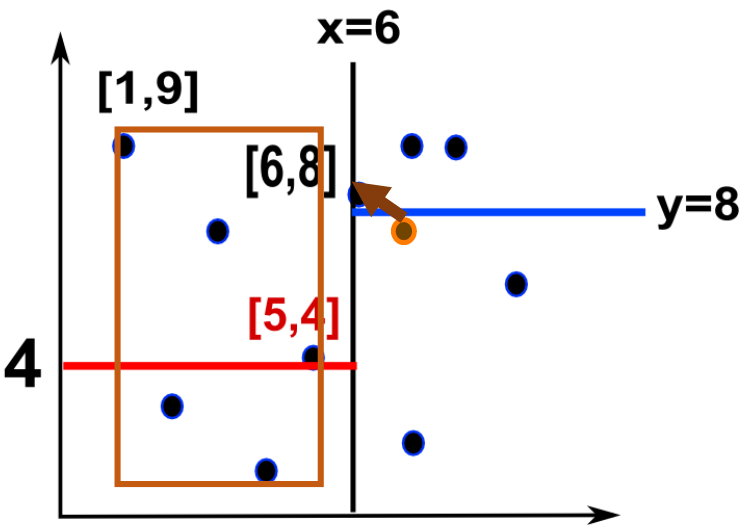Now when we go to check the other branch we notice that the distance to the bounding box is greater than the distance to the nearest neighbor found so far, so we ignore that branch from the search. Since there are no more branches to search in this case we return the point that we have found.

[1,9],[2,3],[4,1],[3,7],[5,4],[6,8],[7,2],[8,8],[7,9],[9,6]
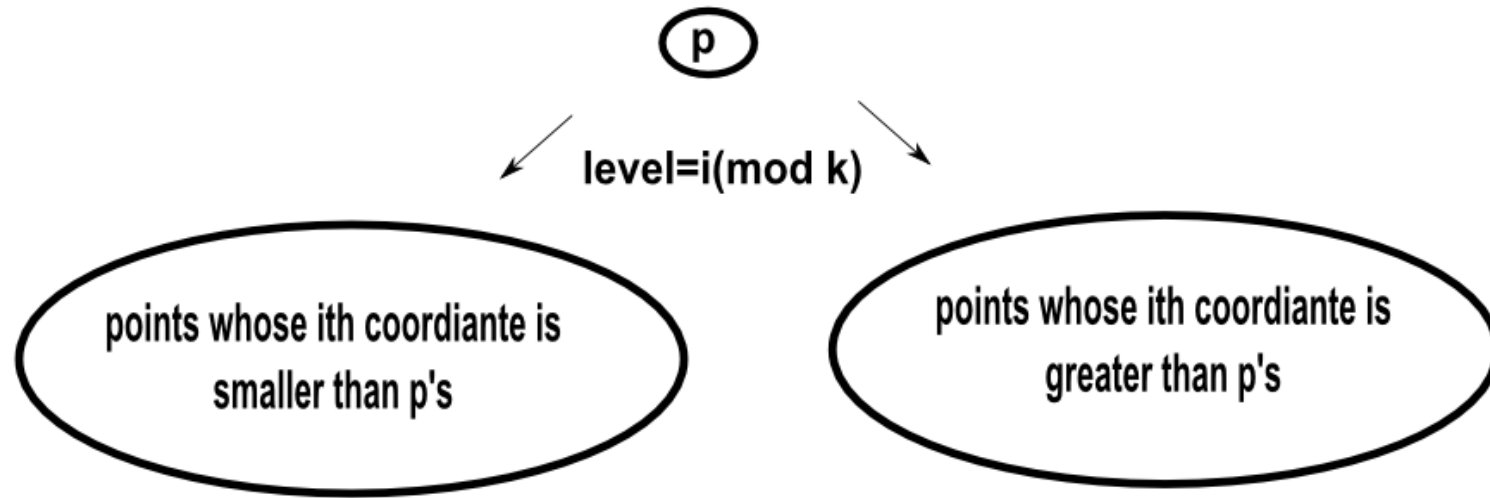
x>=6
[6,8]

[1,9],[2,3],[4,1],
[3,7],[5,4]

[6,8],[7,2],[8,8],
[7,9],[9,6]

y>=4
[5,4]

y>=8
[6,8]

[2,3],[4,1]      [1,9],[3,7],[5,4]

[7,2],[9,6]      [6,8],[8,8],[7,9]

x=6

[1,9]

[6,8]

y=8

[5,4]

4

See another

# Kd-tree: Example



$p$

level=i(mod k)

points whose ith coordiante is smaller than p's

points whose ith coordiante is greater than p's

we keep recursively divising the k-space in this manner

KD-tree is only efficient for small d

# Sklearn Implementations

1-Brute force : Query here costs $O(n*d)$ where $n$ is the size of the data and $d$ is the dimension of the space (# of features)
2-Ball tree : This is approximately $O(d \log(n))$
3-Kd-tree : Depends largely on d. Efficient for $d < 20$ (around $O(d*log(n))$ )

Link for a complete description

# Neighborhood graphs and their relatives

# K-Nearest Neighbor Graph (KNN Graph)

Suppose that we are given a set of points $X = \{p_1, p_2, \ldots, p_n\}$ in $R^d$ with a distance function $d$ defined one them.

Sklearn implementation

# K-Nearest Neighbor Graph (KNN Graph)

Suppose that we are given a set of points $X = \{p_1, p_2, \dots, p_n\}$ in $R^d$ with a distance function $d$ defined one them.

For a fixed integer k, connect the points x, y in $X$ if either $d(x,y) \leq d(x, x_k)$ or $d(x,y) \leq d(y, y_k)$ where $x_k, y_k$ are the $k^{th}$ nearest neighbors of $x, y$ respectively. Doing do for all points we obtain a graph $G(X, E)$ where $E$ is the set of edges connect the points in $X$ as described above.

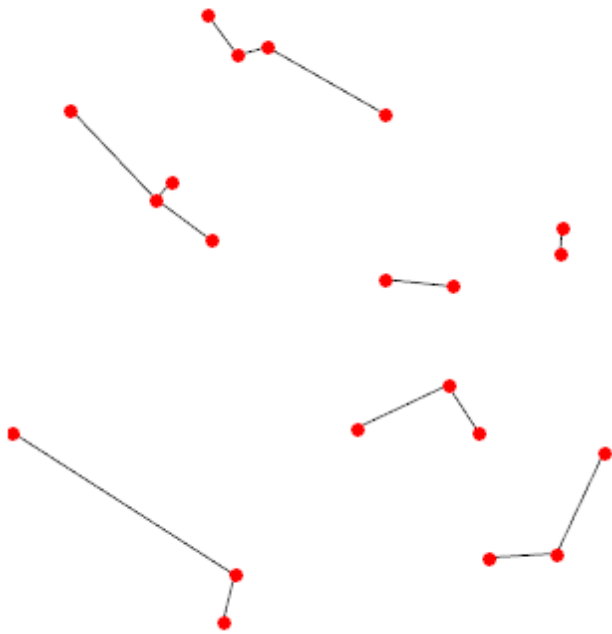Sklearn implementation

# K-Nearest Neighbor Graph (KNN Graph)

Suppose that we are given a set of points $X = \{p_1, p_2, \ldots, p_n\}$ in $R^d$ with a distance function $d$ defined one them.

For a fixed integer k, connect the points x, y in $X$ if either $d(x,y) \leq d(x, x_k)$ or $d(x,y) \leq d(y, y_k)$ where $x_k, y_k$ are the $k^{th}$ nearest neighbors of $x, y$ respectively. Doing do for all points we obtain a graph $G(X, E)$ where $E$ is the set of edges connect the points in $X$ as described above.
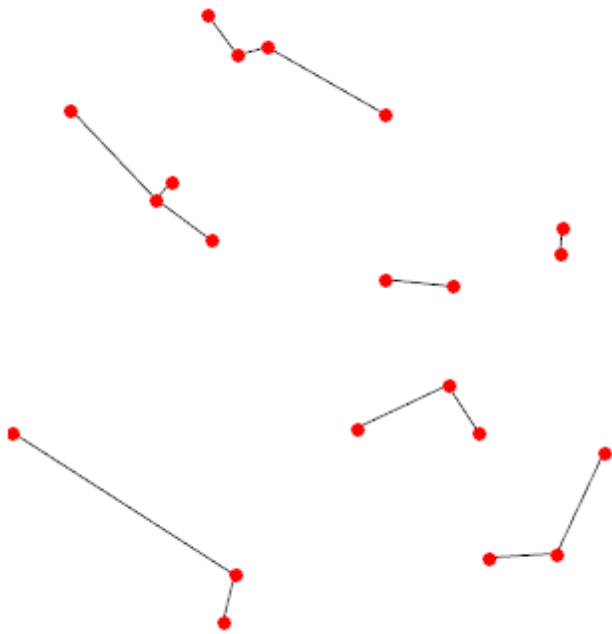
Example of 1-NN graph

Sklearn implementation     The graph that we obtain is called the K-nearest neighbor graph or K-NN graph.

# K-Nearest Neighbor Graph (KNN Graph)

Suppose that we are given a set of points $X = \{p_1, p_2, \ldots, p_n\}$ in $R^d$ with a distance function $d$ defined one them.

For a fixed integer k, connect the points x, y in $X$ if either $d(x, y) \leq d(x, x_k)$ or $d(x, y) \leq d(y, y_k)$ where $x_k, y_k$ are the $k^{th}$ nearest neighbors of $x, y$ respectively. Doing do for all points we obtain a graph $G(X, E)$ where $E$ is the set of edges connect the points in $X$ as described above.
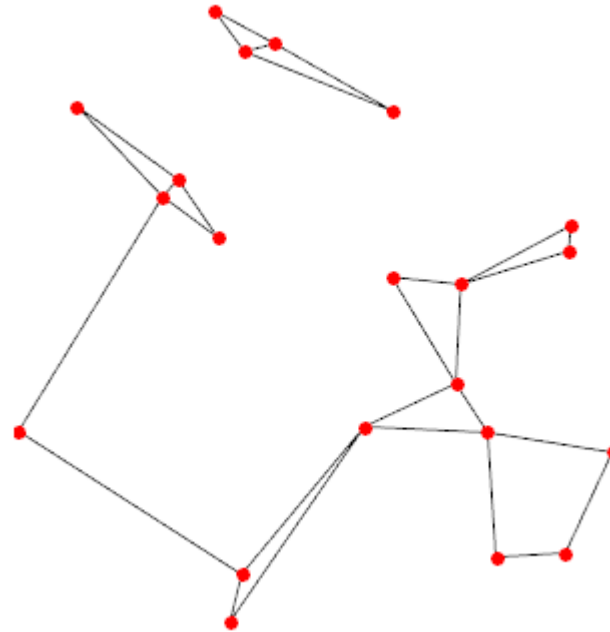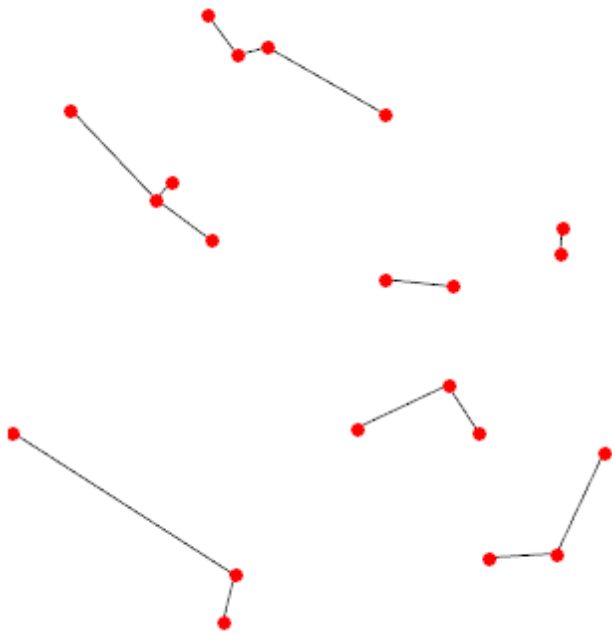


Example of 1-NN graph

Example of 2-NN graph

Sklearn implementation

The graph that we obtain is called the K-nearest neighbor graph or K-NN graph.

# K-Nearest Neighbor Graph (KNN Graph)

Suppose that we are given a set of points $X = \{p_1, p_2, \ldots, p_n\}$ in $R^d$ with a distance function $d$ defined one them.

For a fixed integer k, connect the points x, y in $X$ if either $d(x, y) \leq d(x, x_k)$ or $d(x, y) \leq d(y, y_k)$ where $x_k, y_k$ are the $k^{th}$ nearest neighbors of $x, y$ respectively. Doing do for all points we obtain a graph $G(X, E)$ where $E$ is the set of edges connect the points in $X$ as described above.



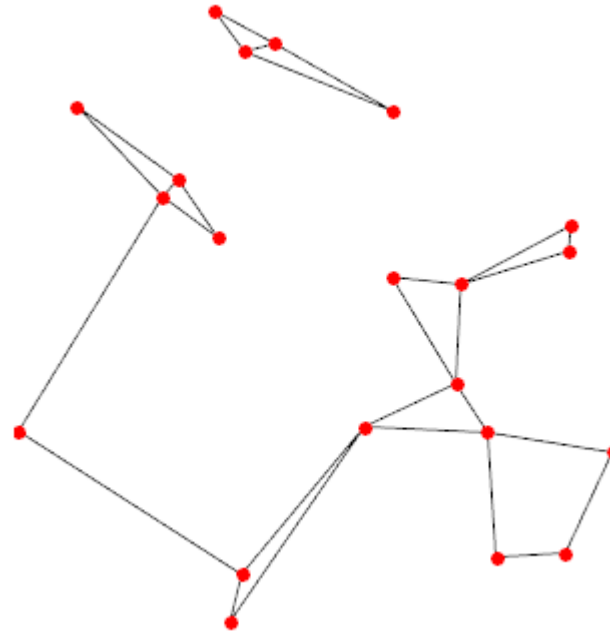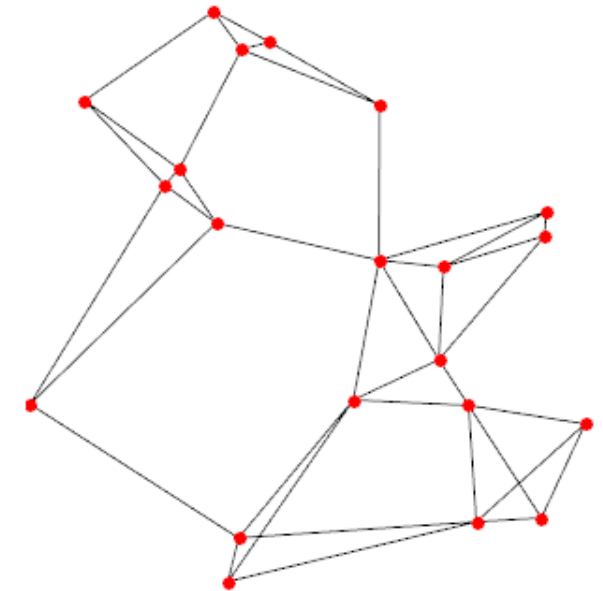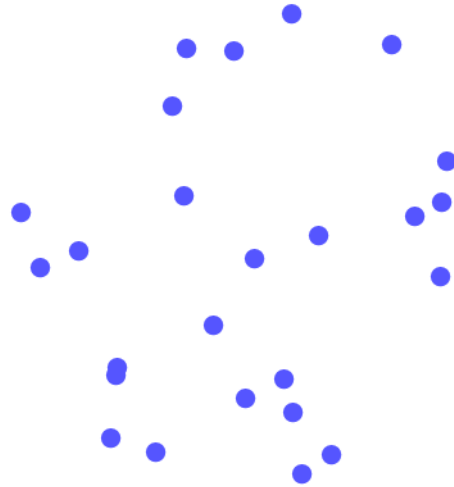Example of 1-NN graph          Example of 2-NN graph          Example of 3-NN graph

Sklearn implementation          The graph that we obtain is called the K-nearest neighbor graph or K-NN graph.

# ε- neighborhood graph

Suppose that we are given a set of points $X = \{p_1, p_2, \ldots, p_n\}$ in $R^d$ with a distance function $d$ defined one them.

For a fixed ε, connect the points $x, y$ if $d(x, y) \leq$ ε.

Doing do for all points we obtain a graph $G(X, E)$ where $E$ is the set of edges connect the points in $X$ as described above.

Sklearn implementation

# ε- neighborhood graph

Suppose that we are given a set of points $X = \{p_1, p_2, \ldots, p_n\}$ in $R^d$ with a distance function $d$ defined one them.

For a fixed $\varepsilon$, connect the points $x, y$ if $d(x, y) \leq \varepsilon$.

Doing do for all points we obtain a graph $G(X, E)$ where $E$ is the set of edges connect the points in $X$ as described above.
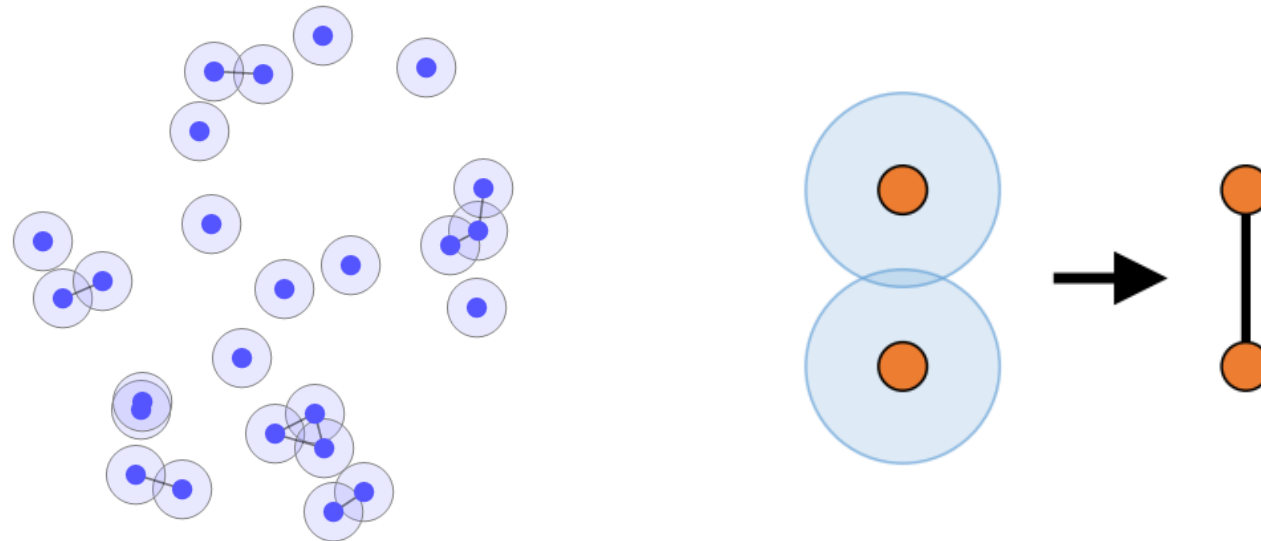
This is *equivalent* to :
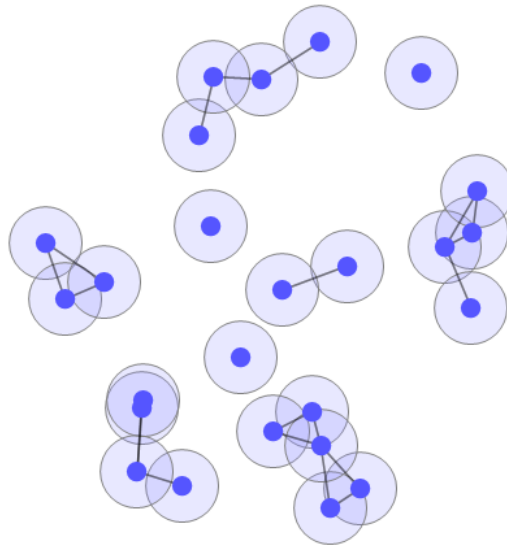Insert an edge if the disks around the points intersect each other.

Note that $d(x, y) \leq \frac{\varepsilon}{2}$ if and only if the two balls surrounding x and y intersect

# ε- neighborhood graph

Suppose that we are given a set of points $X = \{p_1, p_2, \ldots, p_n\}$ in $R^d$ with a distance function $d$ defined one them.

For a fixed $\varepsilon$, connect the points $x, y$ if $d(x, y) \leq \varepsilon$.

Doing do for all points we obtain a graph $G(X, E)$ where $E$ is the set of edges connect the points in $X$ as described above.
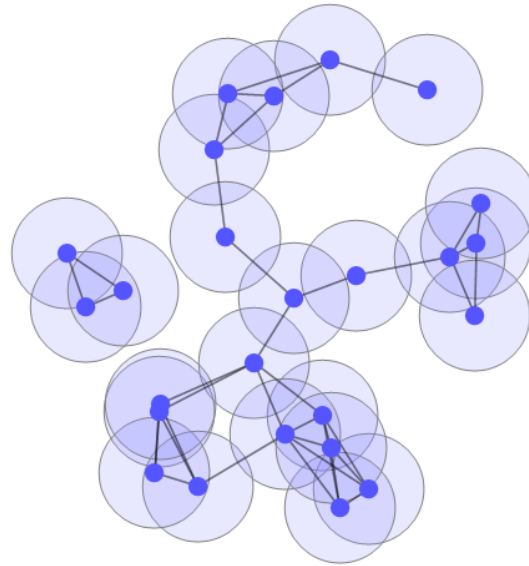


Trying different ε

Sklearn implementation

# ε- neighborhood graph

Suppose that we are given a set of points $X = \{p_1, p_2, \dots, p_n\}$ in $R^d$ with a distance function $d$ defined one them.

For a fixed ε, connect the points $x, y$ if $d(x, y) \leq \varepsilon$.

Doing do for all points we obtain a graph $G(X, E)$ where $E$ is the set of edges connect the points in $X$ as described above.



Trying different ε

# Euclidian Minimal Spanning Tree (EMST)

Suppose that we are given a set of points $X = \{p_1, p_2, \ldots, p_n\}$ in $R^d$ with a distance function $d$ defined one them.

The EMST of X is a minimal spanning tree where the weight of the edge between each pair of points is the distance between those two points.
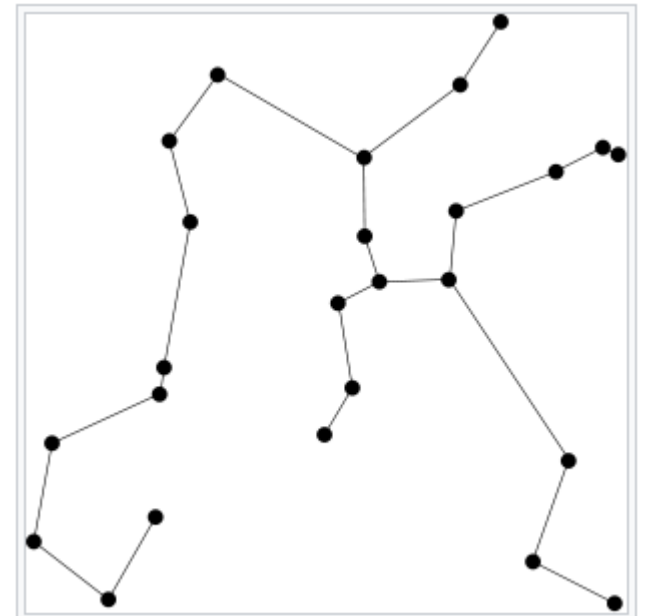


Image source-Wikipedia article