# Quadtree-kNN: Improving Flink's kNN with a quadtree

Daniel Blazevski

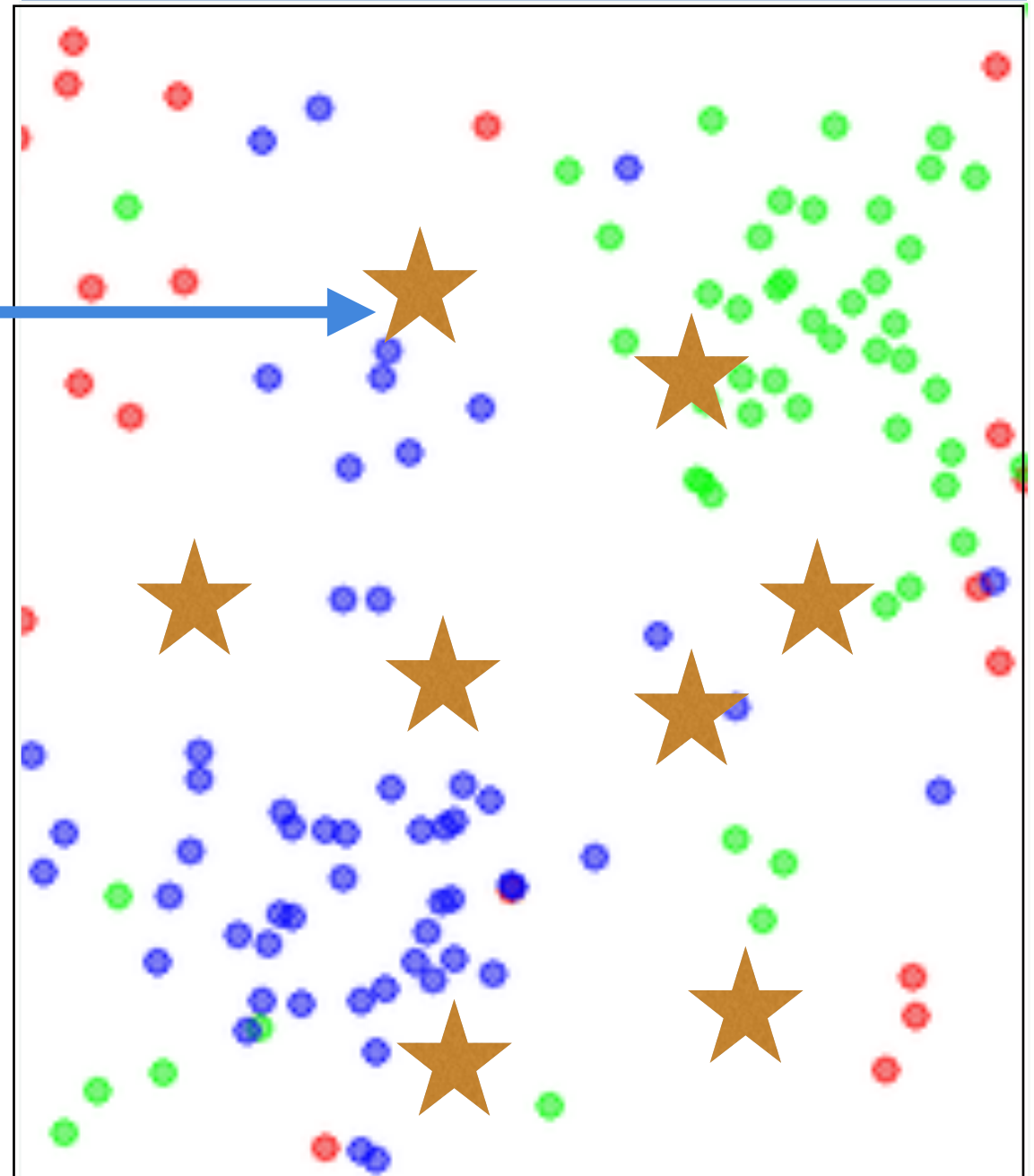Flink INSIGHT

# kNN

kNN: Which class do
the gold stars belong to?



Daniel Blazevski

INSIGHT

# kNN

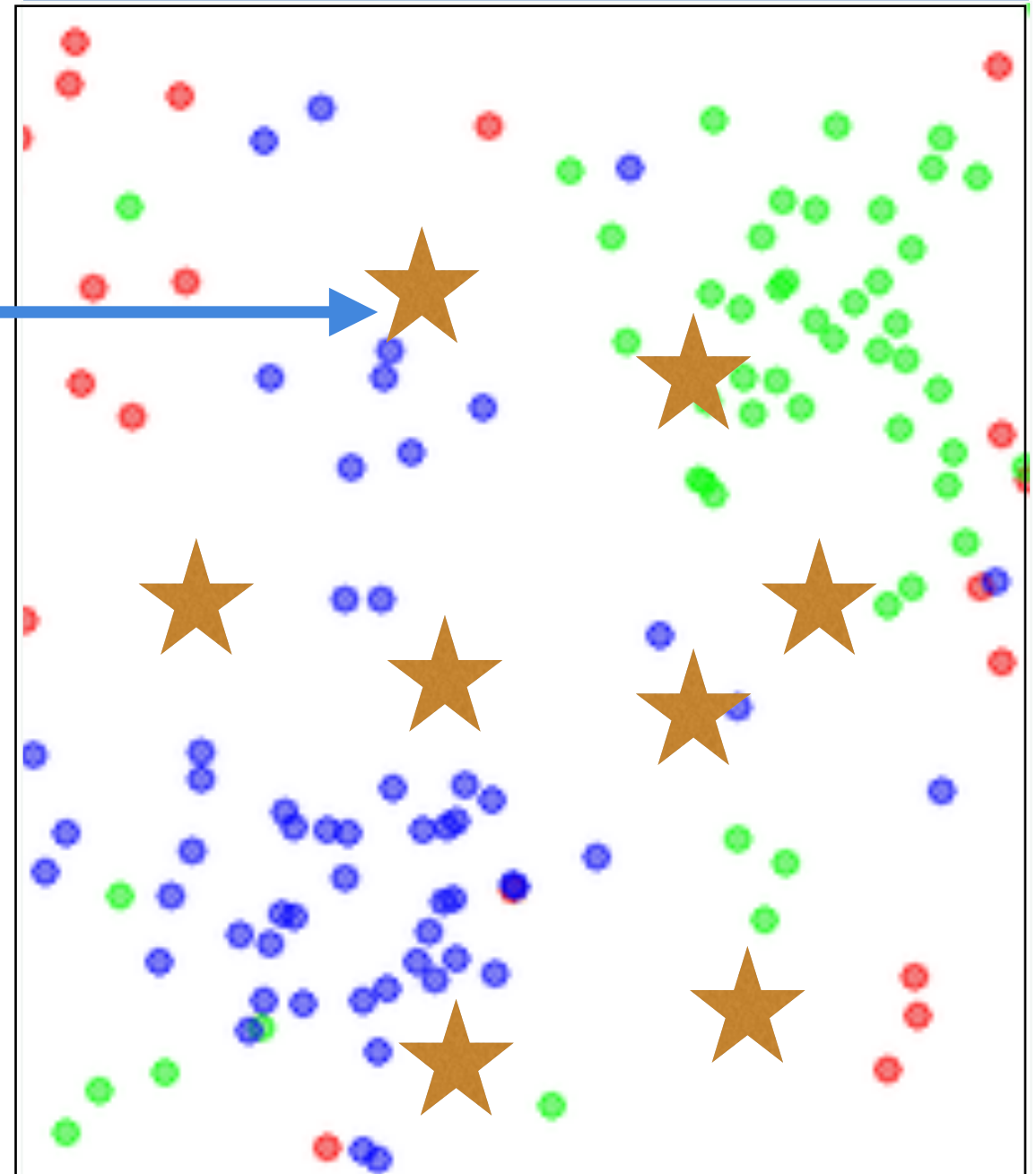kNN: Which class do
the gold stars belong to?
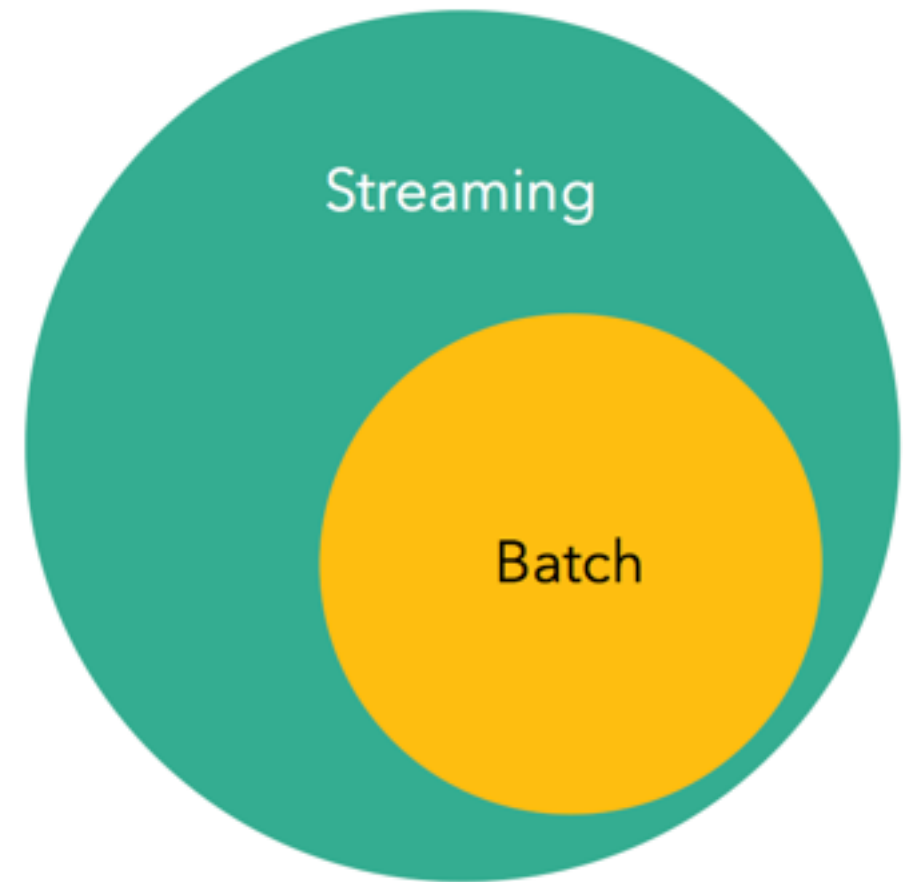
· Look at the k-nearest
  neighbors (kNN)

INSIGHT

# kNN

kNN: Which class do
the gold stars belong to?

- Look at the k-nearest neighbors (kNN)
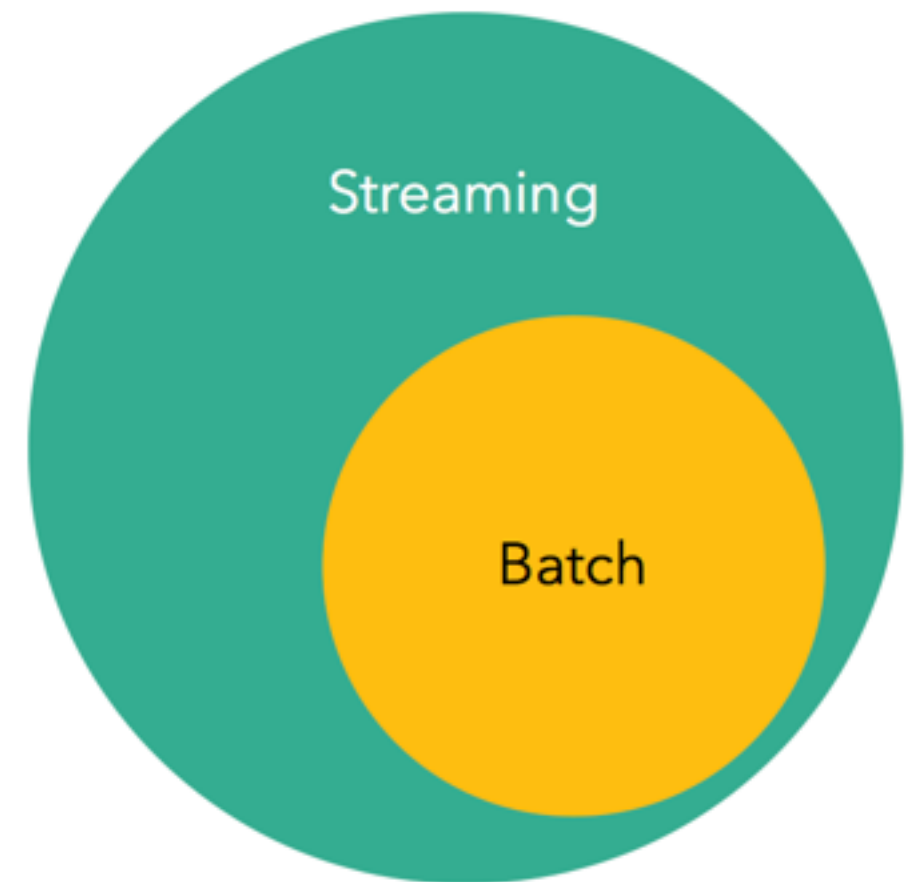
- Wide array of applications in data science

INSIGHT

# Flink

Philosophy: Batch is a subset of Streaming

INSIGHT

# Flink

## Philosophy: Batch is a subset of Streaming

- Distributed data processing tool

- December 2014: Became a top-level Apache project!

**Streaming**

**Batch**

Flink / **FLINK-1745**

## Add exact k-nearest-neighbours algorithm to machine learning library

Agile Board

**Details**

| | |
|---|---|
| Type: | New Feature |
| Status: | IN PROGRESS |
| Priority: | Major |
| Resolution: | Unresolved |
| Affects Version/s: | None |

**People**

Assignee:

Daniel Blazevski

Reporter:

Till Rohrmann
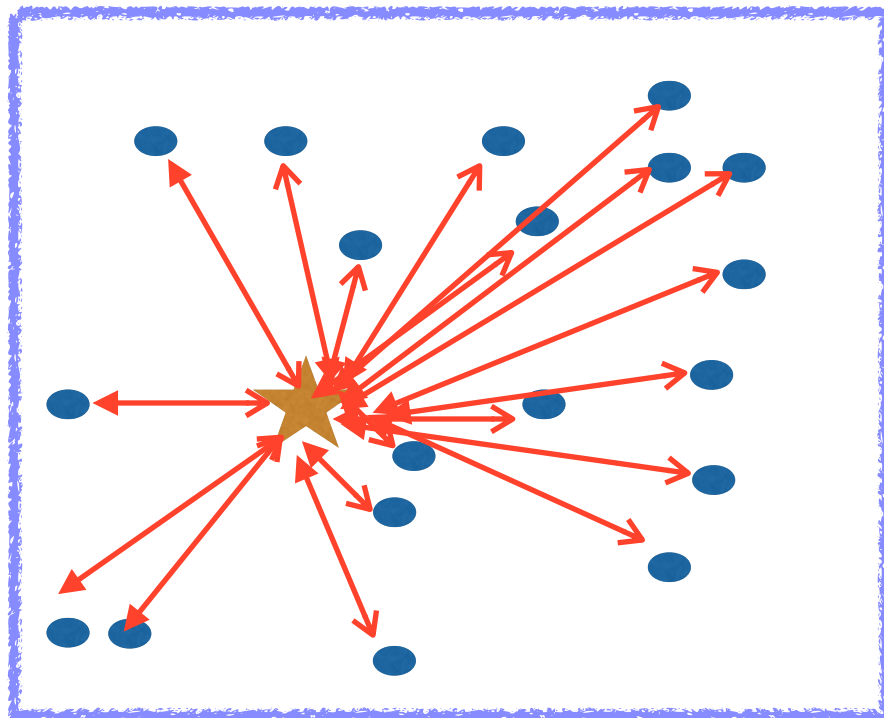
Daniel Blazevski

INSIGHT

# Strategy

- Currently in Flink:
  Compute all pairwise
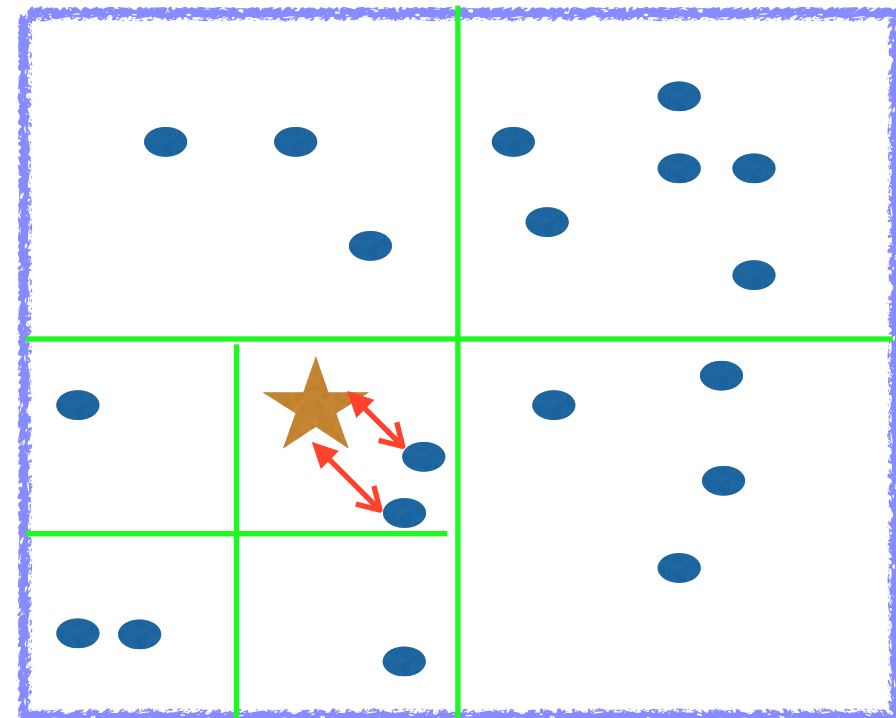  distances — expensive!

INSIGHT

# Strategy

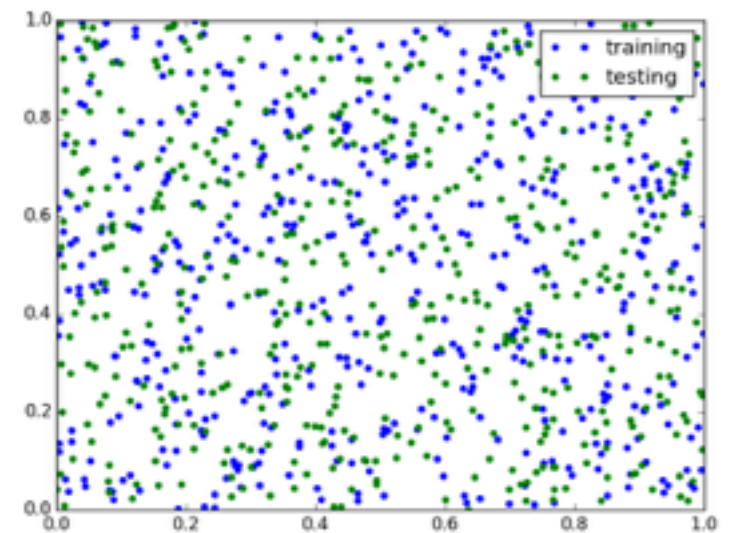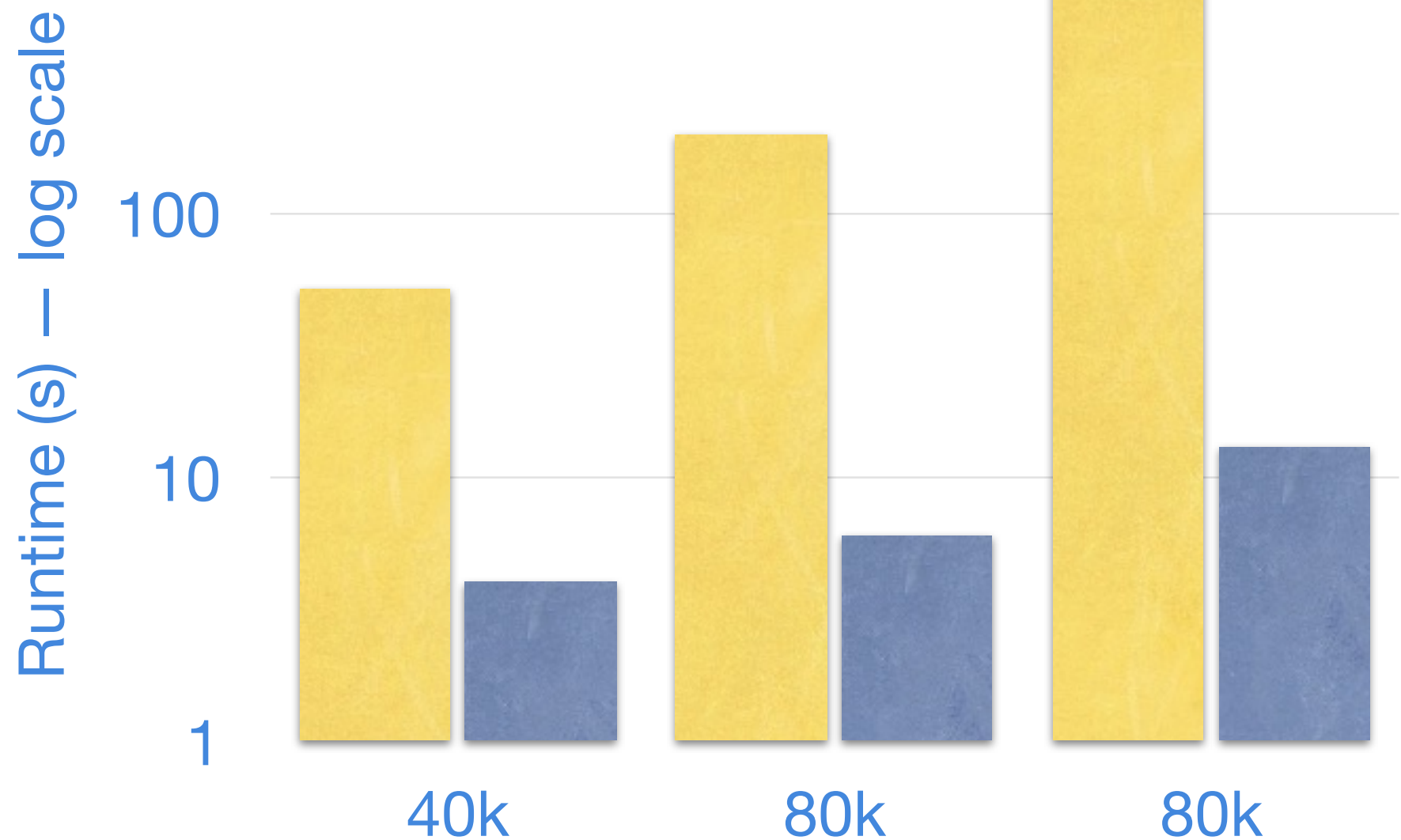- Currently in Flink: Compute all pairwise distances — expensive!

- My work: Partition training set using a Quadtree

# Performance
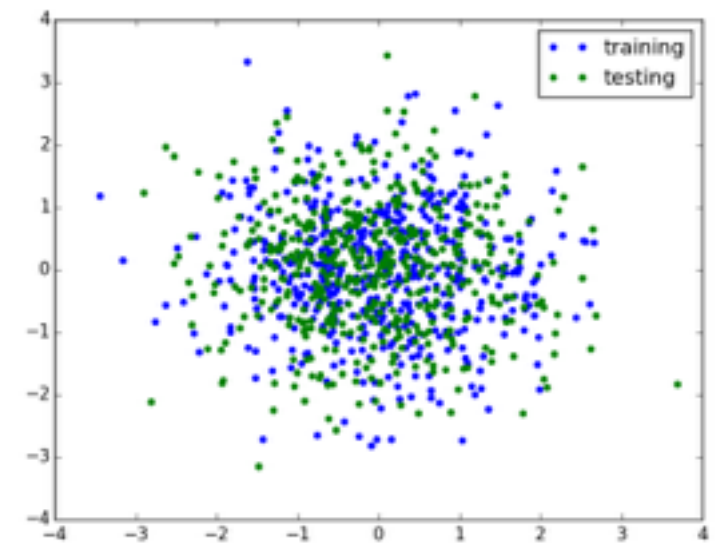


Brute-force   Quadtree

2D uniform data

Runtime (s) — log scale

100

10

1

40k          80k          80k

Total number of points: N_test = N_train

# Performance



Brute-force   Quadtree

Normal distribution

N_test = N_train = 20k

Spatial Dimension

Runtime (s) — log scale

2D   4D   6D

100
10
1

INSIGHT

# Performance

Brute-Force    Quadtree

Normal distribution

Runtime (s) — log scale

10000

100

1

40k    80k    160k    320k

Total number of points: N_test = N_train

Spatial dimension = 6

Daniel Blazevski

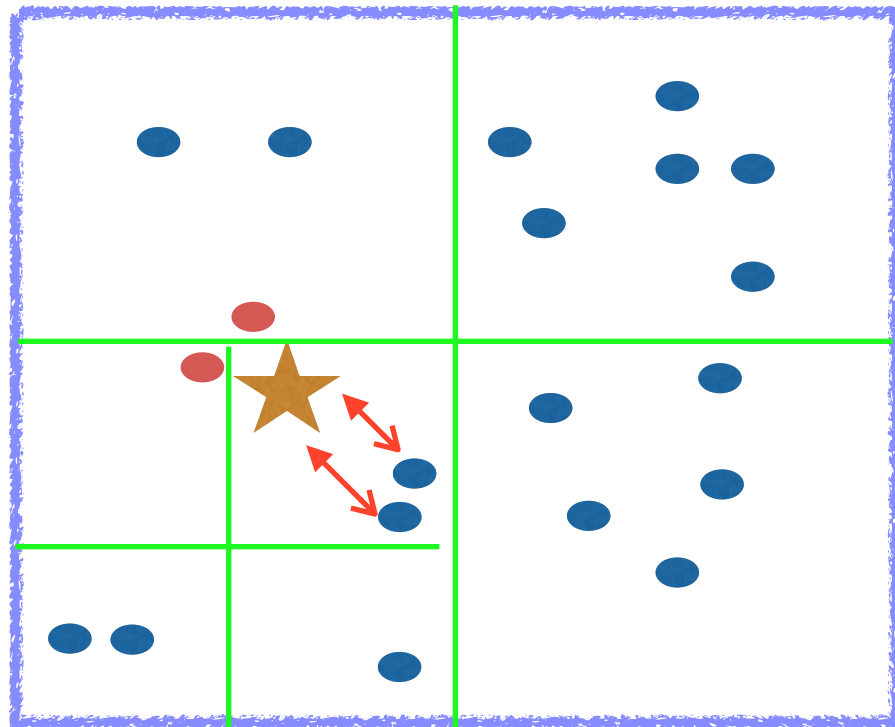INSIGHT

# Challenge

Nearest neighbors may not
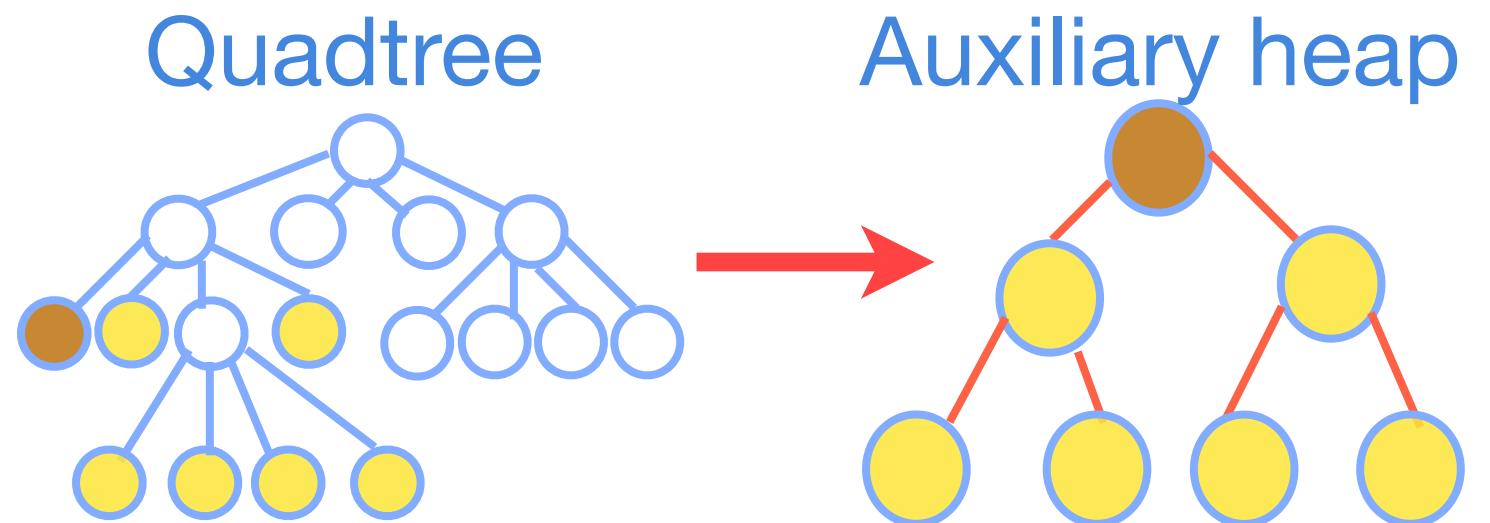be in minimal bounding box
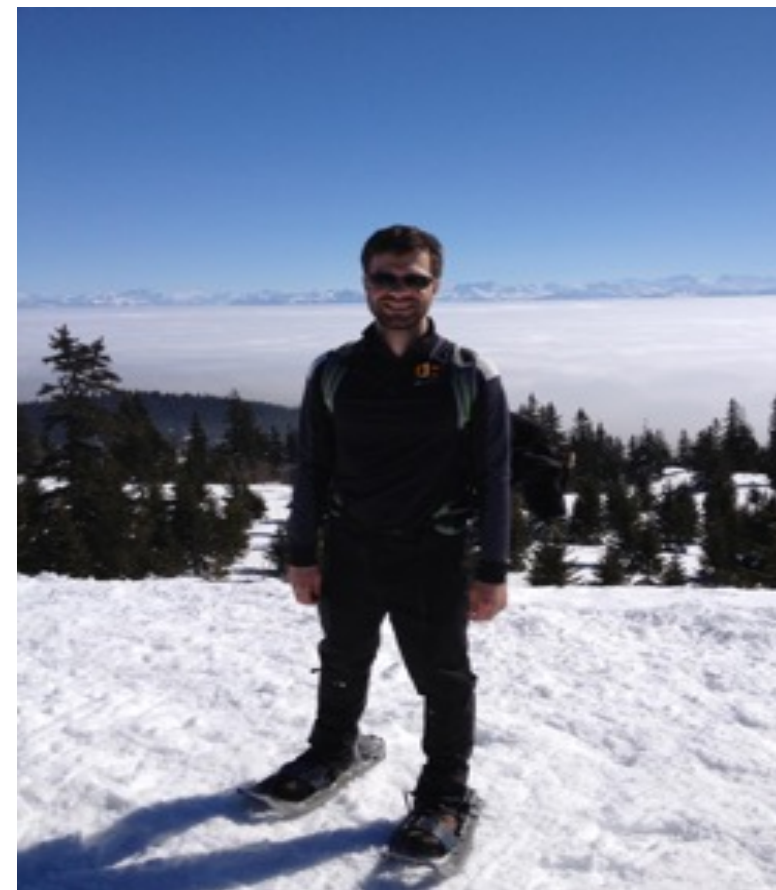
# Challenge

Nearest neighbors may not be in minimal bounding box

- min-heap on siblings' leaf nodes:
  `PriorityQueue[(Double, Node)]`

- Priority = "Dist(star,box)"
  `node.minDist(obj:DenseVector)`

Quadtree

Auxiliary heap

Scala  Flink

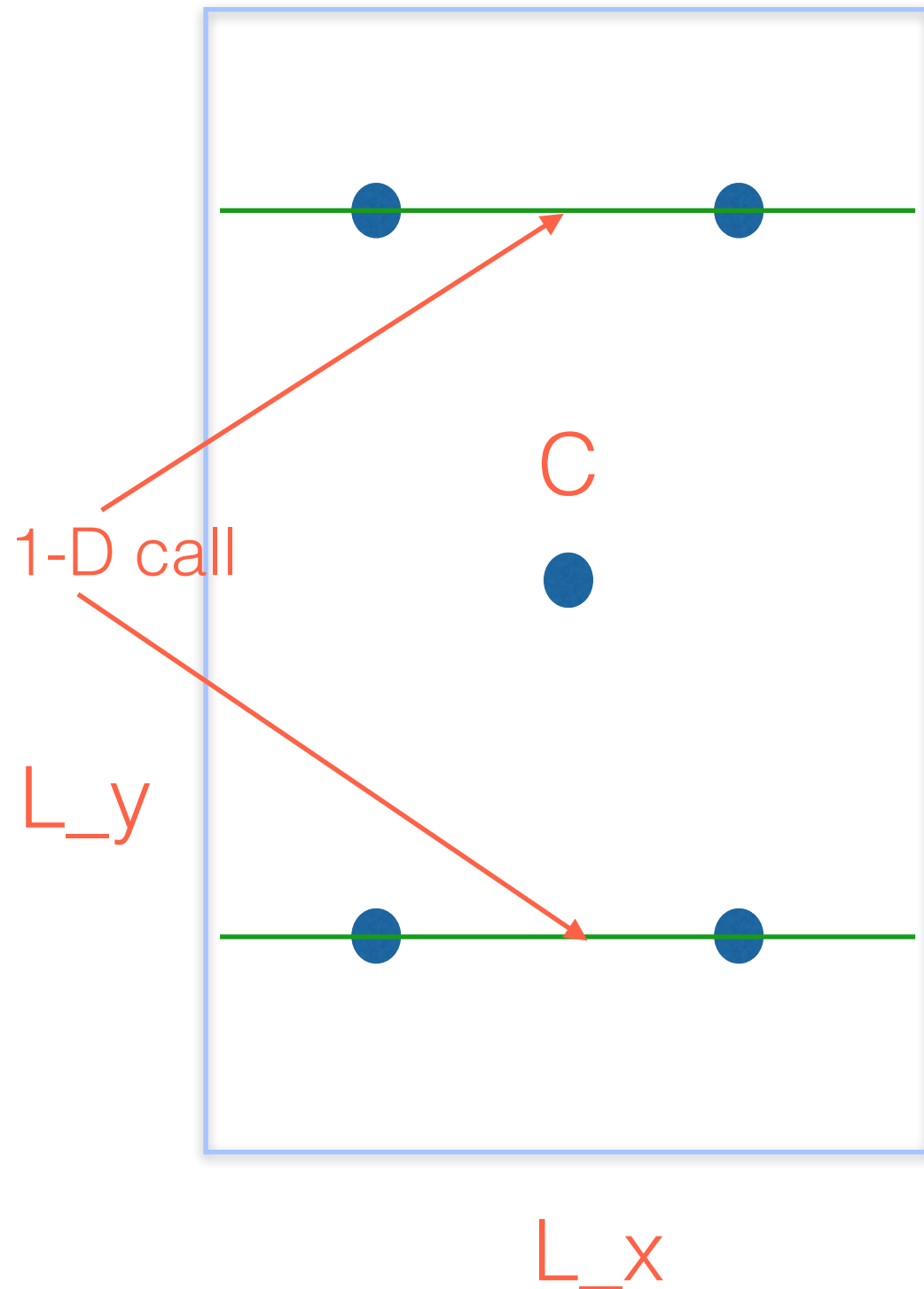Daniel Blazevski

INSIGHT

# About me

- PhD in Math from UT Austin

  - aerospace and fusion energy

- Most recently: Oak Ridge National Laboratory

- Enjoy the outdoors

- daniel.blazevski@gmail.com

- github.com/danielblazevski

- project website: bit.ly/quadtree-flink

INSIGHT

# Back-up slides

INSIGHT

# Partitioning n-dim Box



- Box defined by center C, and Length vector L

- When partitioning, new L is easy: L <− L/2

- Have 2^(d) new centers!

- Use recursion by shifting up and down in the last coordinate

```
cPart ++=
partitionBox(cPartDown,L.take(dim-1),dim-1)

cPart ++=
partitionBox(cPartUp,L.take(dim-1),dim-1)
```

Daniel Blazevski
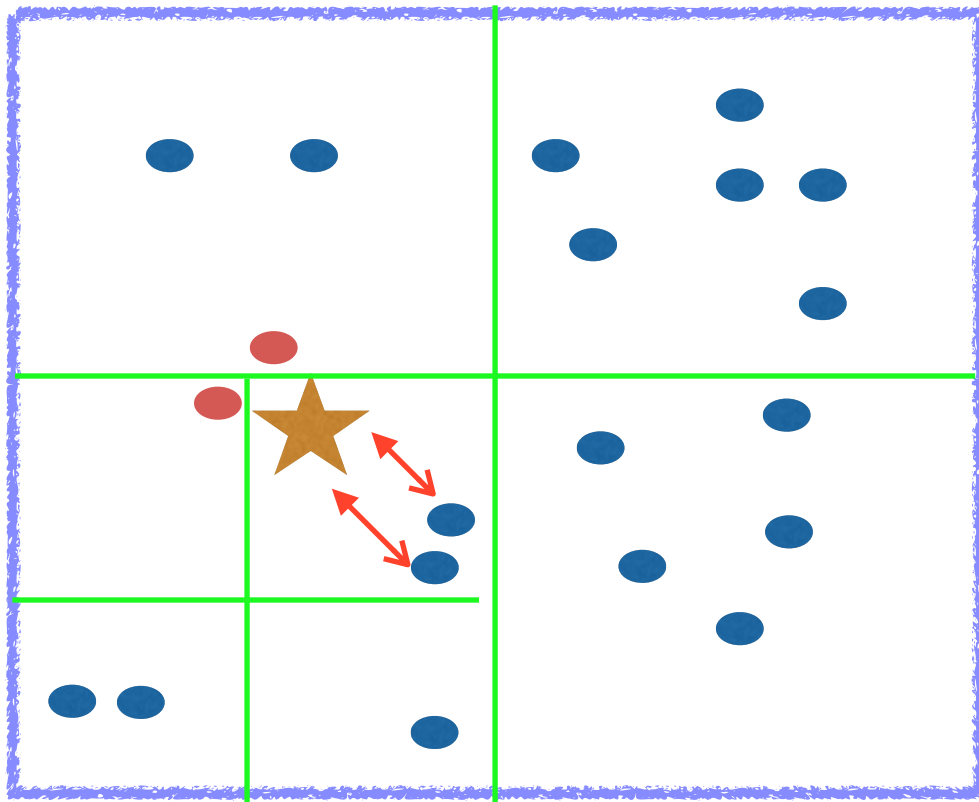
# Challenges

Nearest neighbors may not
be in minimal bounding box



- Use a min-heap on leaf nodes:
`PriorityQueue[(Double, Node)]`

- Priority = "Dist(star,box)"
`node.minDist(obj:DenseVector)`

- Pop leaf nodes + objects until at least k "near" points; look at MAX distance, R

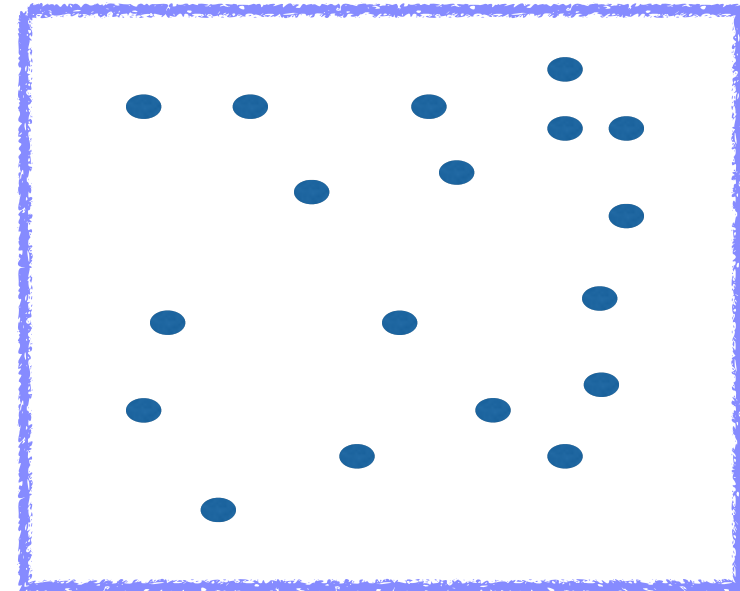- Then search all boxes with Dist(star,box)<R

# Implementation

- Scala code; utilized Flink's existing structures

```scala
class QuadTree
(minVec:ListBuffer[Double],
maxVec:ListBuffer[Double])
```

```scala
def insertRecur(obj:DenseVector,
n:Node)
```

```scala
def searchRecurSiblingQueue
(obj:DenseVector,n:Node,
nodeBuff:ListBuffer[Node])
```

Partition training set on nodes, then form a quadtree on each node