# The K-Nearest Neighbors Algorithm

Team 2

Computerpraktikum Teil 2

4th February 2026

# Presentation Overview

- **Problem Statement**
- **Code Structure**
- **Ball Tree**
  - Constructed Recursively
  - Constructed Iteratively
- **Speed Ups**
- **Error Reduction Strategies**
- **References**

# Problem Statement

# Problem Statement

- Implement a method for **binary classification** in Python
- Given labeled data

$$D = \{(y_i, x_i)\}_{i=1}^{n}, \ y_i \in \{-1, 1\}, \ x_i \in [-1, 1]^d$$

- Learn a classifier $f_D : [-1, 1]^d \to \{-1, 1\}$
- Goal: **minimize the misclassification rate** on unseen test data $D'$
- Approach: $k$-**nearest neighbors** with cross validation to select $k^*$

# Code Structure

# classify.py

- **Setup**
  - Import standard libraries
  - Custom `BallTree` class for nearest neighbor search

- **Data Loading**
  - `load_data`: reads CSV files and stores samples as $(y, x)$ pairs (label and feature vector)

- **Cross Validation**
  - `run_cross_validation`: performs $l$-fold cross validation
  - Evaluates $k = 1, \ldots, K_{\max}$ using cumulative majority voting
  - Computes fold-wise and averaged validation error rates

- **Training**
  - Selection of optimal $k^*$ based on minimal validation error
  - Construction of an ensemble of BallTrees

- **Testing & Output**
  - Ensemble voting for final predictions on the test set
  - Writes predictions, logs, and error curves to disk
  - Optional visualization for 2D datasets

## ball_tree.py

- **Tree Construction**
  - Data is recursively partitioned into hyperspherical regions (balls)
  - Each node stores a center and radius enclosing its data
  - Leaf nodes store the actual data points
  - Tree is built iteratively using a stack (no recursion)
- **Distance Optimization**
  - Precomputation of point norms $\|x\|^2$ for fast distance evaluation
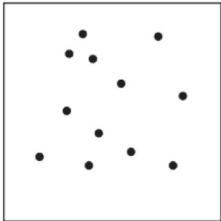  - Uses the identity $\|x - y\|^2 = \|x\|^2 - 2\langle x, y \rangle + \|y\|^2$
- **Query Algorithm**
  - Best-first search using a priority queue ordered by distance to node centers
  - Local and global pruning based on node radius and current worst neighbor
  - Early termination when no closer neighbors are possible

# Speed Ups

# Error Reduction Strategies

# Ball Tree

(A)

(A)  (B)
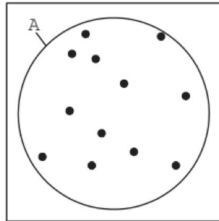
**At each node:**

(A)  (B)

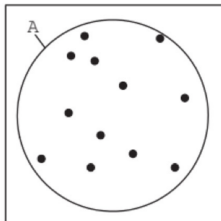**At each node:**

1. Compute center and radius

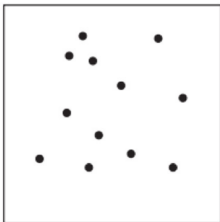(A)          (B)

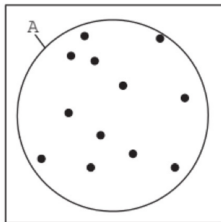**At each node:**
1. Compute center and radius
2. Choose two distant points

(A)          (B)          (C)

**At each node:**

1. Compute center and radius
2. Choose two distant points
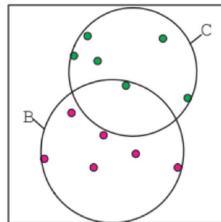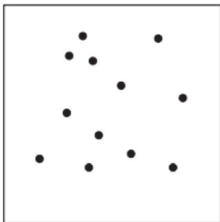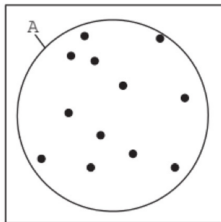3. Split points based on points
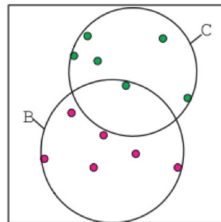
# Ball Tree: Conceptual



(A)          (B)          (C)

(D)

**At each node:**

1. Compute center and radius
2. Choose two distant points
3. Split points based on points
4. Repeat until leaf size reached

# Ball Tree: Conceptual
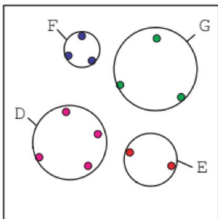


(A)     (B)     (C)

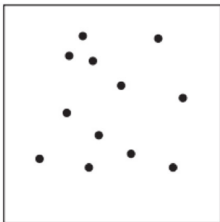(D)     (E)

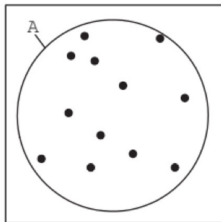**At each node:**
1. Compute center and radius
2. Choose two distant points
3. Split points based on points
4. Repeat until leaf size reached

**Leaf nodes store data points**

# Ball Tree: Recursive

```python
def BallTree:

    def __init__(self, data, leaf_size=1):
        self.leaf_size = leaf_size
        self.points = data
        self.left = None
        self.right = None
        self.center = self._computer_center(data)
        self.radius = self._compute_radius(data, self.center)
```

# Ball Tree: Recursive

```python
def BallTree:

    def __init__(self, data, leaf_size=1):
        self.leaf_size = leaf_size
        self.points = data
        self.left = None
        self.right = None
        self.center = self._computer_center(data)
        self.radius = self._compute_radius(data, self.center)

        if len(data) > self.leaf_size:
            self._split()
```

# Ball Tree: Recursive

```python
def BallTree:

    def __init__(self, data, leaf_size=1):
        self.leaf_size = leaf_size
        self.points = data
        self.left = None
        self.right = None
        self.center = self._computer_center(data)
        self.radius = self._compute_radius(data, self.center)

        if len(data) > self.leaf_size:
            self._split()

    def _split(self):

        ... # Split points

        self.left = BallTree(left_points, self.leaf_size)
        self.right = BallTree(right_points, self.leaf_size)
        self.points = None
```
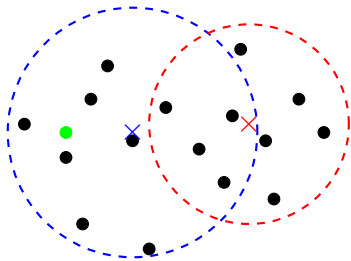
```
def BallTree:

    ...
```

```
def BallTree:

    ...
```

```
def BallTree:

    ...
```

```python
def BallTree:

    ...

    def query(self, x, k, knn=[]):
        lower_bound = max(0.0, distance(x, self.center) - self.radius)
```

```
def BallTree:

    ...

    def query(self, x, k, knn=[]):
        lower_bound = max(0.0, distance(x, self.center) - self.radius)
        if len(knn) == k and lower_bound > knn[0][0]:
            return
```

```python
def BallTree:

    ...

    def query(self, x, k, knn=[]):
        lower_bound = max(0.0, distance(x, self.center) - self.radius)
        if len(knn) == k and lower_bound > knn[0][0]:
            return

        if self.left is None and self.right is None:

            ... # Pointsearch in leaf node
```

```python
def BallTree:

    ...

    def query(self, x, k, knn=[]):
        lower_bound = max(0.0, distance(x, self.center) - self.radius)
        if len(knn) == k and lower_bound > knn[0][0]:
            return

        if self.left is None and self.right is None:

            ... # Pointsearch in leaf node

        if distance(x, self.left.center) < distance(x, self.right.center):
            self.left.query(x, k, knn)
            self.right.query(x, k, knn)
        else:
            ...
```
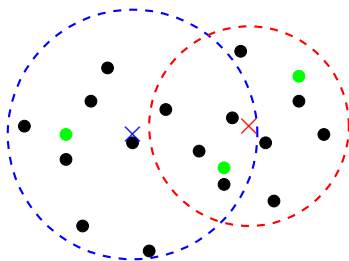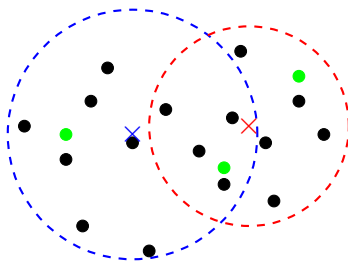
1. Use **built-in functions** (C-optimized).
   For example, instead of list comprehension:

```
def vector_sum(vectors):
    return [sum(components) for components in zip(*vectors)]
```

1. Use **built-in** **functions** (C-optimized).
   For example, instead of list comprehension:

   ```
   def vector_sum(vectors):
       return [sum(components) for components in zip(*vectors)]
   ```

2. **Avoid defining** often used functions.

# Ball Tree: Recursive Speed

1. Use **built-in** **functions** (C-optimized).
   For example, instead of list comprehension:

   ```python
   def vector_sum(vectors):
       return [sum(components) for components in zip(*vectors)]
   ```

2. **Avoid defining** often used functions.
3. Save often used variables **locally**.

1. Use **built-in functions** (C-optimized). ✓
   For example, instead of list comprehension:

   ```python
   def vector_sum(vectors):
       return [sum(components) for components in zip(*vectors)]
   ```

2. **Avoid defining** often used functions.
3. Save often used variables **locally**.

1. Use **built-in** **functions** (C-optimized). ✓
   For example, instead of list comprehension:

   ```
   def vector_sum(vectors):
       return [sum(components) for components in zip(*vectors)]
   ```

2. **Avoid defining** often used functions. ✗
3. Save often used variables **locally**.

1. Use **built-in functions** (C-optimized). ✓
   For example, instead of list comprehension:

   ```python
   def vector_sum(vectors):
       return [sum(components) for components in zip(*vectors)]
   ```

2. **Avoid defining** often used functions. ✗
3. Save often used variables **locally**. ✗

1. Use **built-in functions** (C-optimized). ✓
   For example, instead of list comprehension:

   ```
   def vector_sum(vectors):
       return [sum(components) for components in zip(*vectors)]
   ```

2. **Avoid defining** often used functions. ✗
3. Save often used variables **locally**. ✗

# Recursive Method not optimal !!!

# References

author (XXXX). "title". In.