# 1    Space filling curves for NN-queries                          (1 P.)

1. Implement space-filling z-curves in a language of your choice. Your program has to take two files which contain points as input. The first file contains all the points of the base data set, while the second file contains query points. The program has to calculate and display:

   - The k-NN, of each query point, in the base data set, based on the actual distance.
   - The k-NN, of each query point, in the base data set, based on the z-curve distance.

   You can use the template in OLAT, which already parses the files and provides utility classes. Submit the code and the output of your program when executed with $k = 3$ and the two data files provided in OLAT. If you do not use the template, also submit instructions on how to compile and execute your program.

   *Note:* If you use code from external sources, provide the source as a comment.

**Solution:**

```
ZCurve exercise:
=========================================
Reading file: points.txt
Reading file: queries.txt
=========================================
Points:
(391,1) {81943}
(275,629) {600871}
(653,246) {322169}
(168,146) {50760}
(72,855) {668266}
(893,265) {464339}
(659,74) {287117}
(343,250) {113565}
(845,631) {866939}
(195,884) {686629}
(977,784) {1004289}
(677,350) {419513}
(9,27) {715}
(681,671) {837355}
(556,154) {296664}
=========================================
Queries:
(358,681) {629910}
(463,667) {643807}
(157,270) {147961}
=========================================
KNN (k = 3):
(358,681) {629910}: [(275,629) {600871}, (195,884) {686629}, (681,671) {837355}]
(463,667) {643807}: [(275,629) {600871}, (681,671) {837355}, (195,884) {686629}]
(157,270) {147961}: [(168,146) {50760}, (343,250) {113565}, (9,27) {715}]
=========================================
KNN ZCurve (k = 3):
(358,681) {629910}: [(275,629) {600871}, (72,855) {668266}, (195,884) {686629}]
(463,667) {643807}: [(72,855) {668266}, (195,884) {686629}, (275,629) {600871}]
(157,270) {147961}: [(343,250) {113565}, (391,1) {81943}, (168,146) {50760}]
=========================================

Process finished with exit code 0
```

2. Which differences can you see between the two results of your implementation? Explain why are or why are not the results the same.
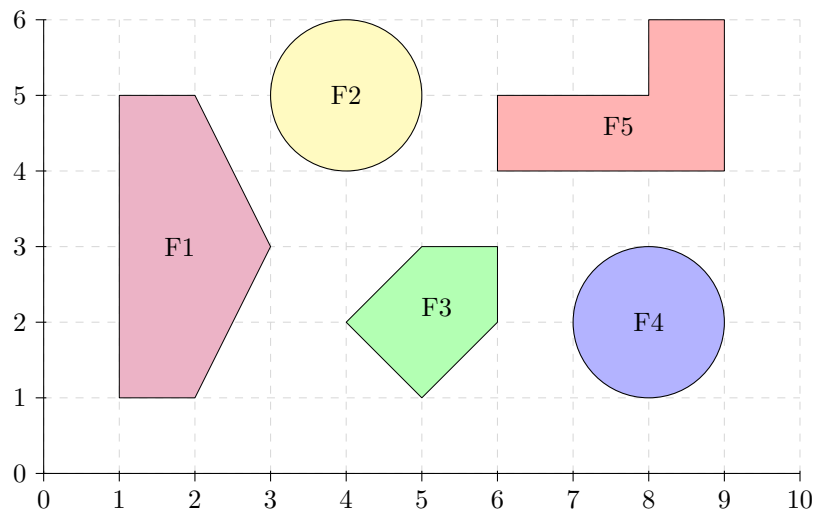
   **Solution:**
   Some of the actual k-NN are matching the k-NN according to the z-curve value, but not all of them. The z-curve mapping is only an approximation of the actual k-NN.

   Let's have a look at the points $A(0,0)$, $B(1,0)$, and $C(0,1)$. The real distance between $A$ and $B$ is 1. Also, the z-value is 1 in this case. The real distance between $A$ and $C$ is also 1, but in this case, the z-value is 2. Because of the shape of the z-curve, the z-value is only an estimation of the real distance.

# 2   R Tree                                                            (1 P.)
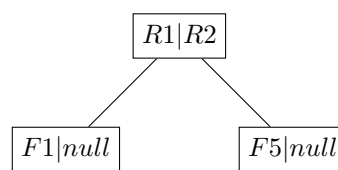
Given the following two-dimensional objects:



For the following R tree operations, explain exactly which steps are performed:

1. *Store the objects F1, F2, F3, F4, F5 in an initially empty R tree.* One node fits 1–2 entries.

   **Solution:**
   The objects get inserted one after the other in theory they can be inserted in any node. I decided for the following insertion order: F1, F5, F4, F2, F3. F1 and F5 get inserted into the tree.
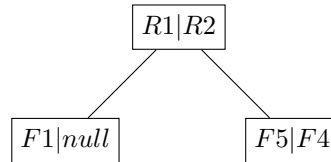


   R1 is the smalest rectangle that fits around F1 and R2 is the smalest rectangle that fits arround F5 Insertion of F4:
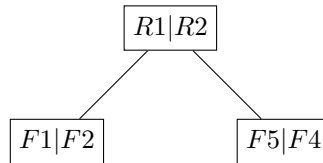   Calculate the arrea differences of R1 if F4 gets inserted into R1 and calculate the area difference of R2 if F4 gets inserted.

$$A_{R1} = 8 \cdot 4 - 2 \cdot 4 = 24$$
$$A_{R2} = 5 \cdot 3 - 6 \cdot 4 = 9$$

   Insert F4 int R2 because the area increase is smaler then when inserted into R2

$$\boxed{R1|R2}$$

$$\boxed{F1|null} \qquad \boxed{F5|F4}$$

Insertion of F2:
Calculate the arrea differences of R1 if F2 gets inserted into R1 and calculate the area difference of R2 if F2 gets inserted.

$$A_{R1} = 4 \cdot 5 - 2 \cdot 4 = 12$$
$$A_{R2} = 5 \cdot 3 - 3 \cdot 5 = 15$$

Insert F2 int R1 because the area increase is smaler then when inserted into R1

$$\boxed{R1|R2}$$

$$\boxed{F1|F2} \qquad \boxed{F5|F4}$$

Insertion of F3:
Calculate the area differences of R1 if F2 gets inserted into R1 and calculate the area difference of R2 if F2 gets inserted.

$$A_{R1} = 5 \cdot 5 - 4 \cdot 5 = 5$$
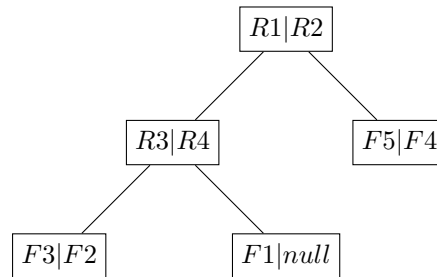$$A_{R2} = 5 \cdot 5 - 3 \cdot 5 = 10$$

Insert F2 int R1 because the area increase is smaler then when inserted into R1. The maximum of objects in R1 is already reached. In order to insert F2 R1 needs to be split into two rectangles. The size of the two newly created rectangels R3 and R4 should be minimal in the optimal case. Calculation of the area of all posible split combinations:
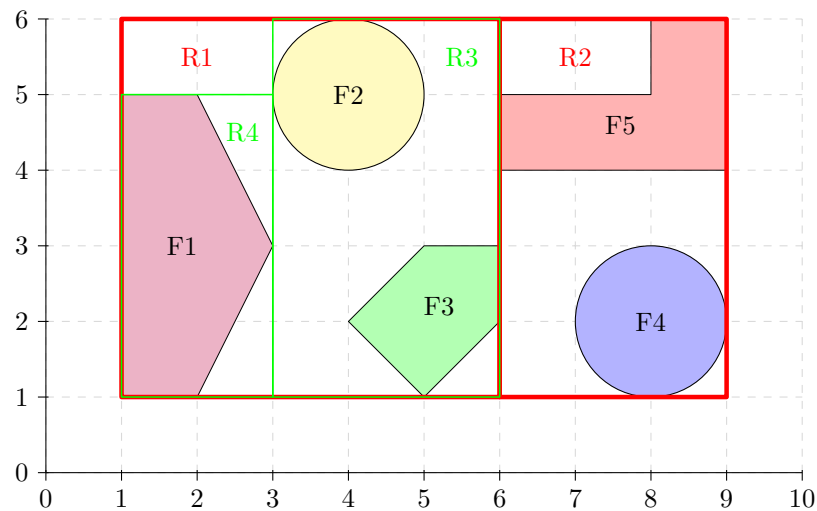
$$R3[F1, F2], R4[F3] \; A = 4 \cdot 5 + 2 \cdot 2 = 24$$
$$R3[F1, F3], R4[F2] \; A = 5 \cdot 4 + 2 \cdot 2 = 24$$
$$R3[F3, F2], R4[F1] \; A = 3 \cdot 5 + 2 \cdot 4 = 23$$

According to the Calculation the last option is to optimal one because the area of R3 and R4 will be minimal. The R-Tree looks now like:

$$R1|R2$$

$$R3|R4 \qquad F5|F4$$

$$F3|F2 \qquad F1|null$$

In the two dimensional plane the R-Tree looks like the following:



2. *Find all objects, containing the point (6,2).*

   **Solution:**

   Step by step execution:

   > Scan the root of the tree (R1,R2)
   > Check if (6,2) is in R1 ⇒ true ⇒ put R3 and R4 in the que
   > Check if (6,2) is in R2 ⇒ true ⇒ put F5 and F4 in the que
   > Check if (6,2) is in R3 ⇒ true ⇒ put F2 and F3 in the que
   > Check if (6,2) is in R4 ⇒ false
   > Check if (6,2) is in F5 ⇒ false
   > Check if (6,2) is in F4 ⇒ false
   > Check if (6,2) is in F2 ⇒ false
   > Check if (6,2) is in F3 ⇒ true ⇒ leafnode ⇒ add F3 to the query result
   > (6,2) is only in object F3

3. *Find all objects that are positioned completely in the rectangle Q, which is defied by the points* $(2,1)(9,3)$.

4. *Find all objects, intersecting with Q′:* $(2,2)(4,4)$.

   We assume that sharing exactly one point, also counts as intersecting.

## 3   Index Structures in Metric Space                                    (1 P.)

Show, for each of the following distance functions, that the properties for being a metric are fulfilled, or provide a counterexample.

1. For two vectors (or points) $a = (x_a, y_a)$ and $b = (x_b, y_b)$:

   $d(a, b) \mapsto |x_a - x_b| + |y_a - y_b|$.

2. For two vectors (or points) $a = (x_a, y_a)$ and $b = (x_b, y_b)$: $d(a, b) \mapsto (x_a - x_b)^2 + (y_a - y_b)^2$.

3. For two strings $s_1$ and $s_2$, where $S$ is the set of characters of the string (e.g.: $S(\text{``Codd''}) = \{\text{`C'}, \text{`o'}, \text{`d'}\}$): $d(s_1, s_2) \mapsto 2 * |S(s_1) \cap S(s_2)| \, / \, (|S(s_1)| + |S(s_2)|)$.
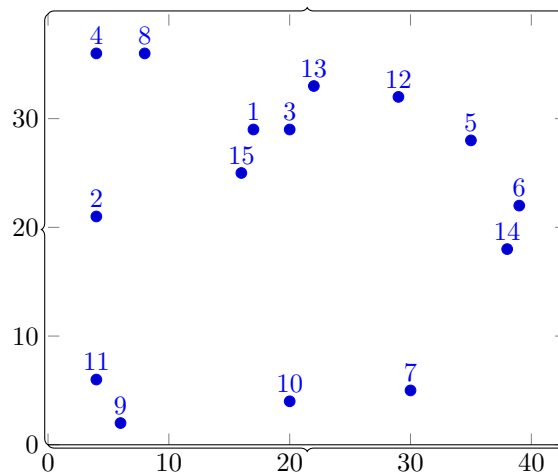
# 4 Misc Metric Indexing (1 P.)

Below are the points used for the following questions:

$$1 : (17, 29) \quad 2 : (4, 21) \quad 3 : (20, 29) \quad 4 : (4, 36) \quad 5 : (35, 28)$$

$$6 : (39, 22) \quad 7 : (30, 5) \quad 8 : (8, 36) \quad 9 : (6, 2) \quad 10 : (20, 4)$$

$$11 : (4, 6) \quad 12 : (29, 32) \quad 13 : (22, 33) \quad 14 : (38, 18) \quad 15 : (16, 25)$$

1. **GH Tree partitioning:** Create a GH partitioning, such that the leaf nodes of the tree have at most 2 elements. Draw the tree and draw the partitioning into the plot below.



2. **VP Tree:** Given the following VP tree, search for query point (25,4) with $\varepsilon = 11$. Describe which parts of the tree you pruned and why.