

Exercise 1

(i)

To show that D_{Hamming} is a metric we have to show that D_{Hamming} satisfies the usual conditions for a metric:

- (a) $D_{\text{Hamming}}(x, y) \geq 0$ and $D_{\text{Hamming}}(x, y) = 0$ iff $x = y$
- (b) $D_{\text{Hamming}}(x, y) = D_{\text{Hamming}}(y, x)$
- (c) $D_{\text{Hamming}}(x, z) \leq D_{\text{Hamming}}(x, y) + D_{\text{Hamming}}(y, z)$ for any $x, y, z \in F^{(m)}$

Definition of D_{Hamming} :

The Hamming distance $D_{\text{Hamming}}(x, y)$ between two vectors $x, y \in F^{(m)}$ is the number of coefficients in which they differ.

$$D_{\text{Hamming}}(x, y) = |x - y| = \sum_{i=1}^m |x_i - y_i|$$

Proof:

- (a) $D_{\text{Hamming}}(x, y) = 0$ iff x, y agree in all coordinates and this happens iff $x = y$.
 $D_{\text{Hamming}}(x, y) \geq 0$ is given because we sum up the absolute value of $x - y$.
- (b) Because of the absolute value of $x - y$ in the calculation of $D_{\text{Hamming}}(x, y)$ $x - y = y - x$ gives us the same value. So $D_{\text{Hamming}}(x, y) = D_{\text{Hamming}}(y, x)$.
- (c) $D_{\text{Hamming}}(x, y)$ is equal to the minimal number of coordinate changes necessary to get from x to y . In its turn, $D_{\text{Hamming}}(y, z)$ is equal to the minimal number of coordinate changes necessary to get from y to z .
 So $D_{\text{Hamming}}(x, y) + D_{\text{Hamming}}(y, z)$ changes will get us from x to z .
 Hence $D_{\text{Hamming}}(x, z) \leq D_{\text{Hamming}}(x, y) + D_{\text{Hamming}}(y, z)$ which is the minimal number of coordinate changes necessary to get from x to z .

Since we can show that all conditions for a metric are fulfilled, D_{Hamming} is a metric.

(ii)

To show that D_{Δ} is a metric on the power set 2^U of U , we have to show that D_{Δ} satisfies the usual conditions for a metric:

- (a) $D_{\Delta}(S_1, S_2) \geq 0$ and $D_{\Delta}(S_1, S_2) = 0$ iff $S_1 = S_2$
- (b) $D_{\Delta}(S_1, S_2) = D_{\Delta}(S_2, S_1)$
- (c) $D_{\Delta}(S_1, S_3) \leq D_{\Delta}(S_1, S_2) + D_{\Delta}(S_2, S_3)$ for any $S_1, S_2, S_3 \subseteq U$

Definition of D_{Δ} :

$$D_{\Delta}(S_1, S_2) = |S_1 \Delta S_2| = |(S_1 \setminus S_2) \cup (S_2 \setminus S_1)| \text{ for any } S_1, S_2 \subseteq U$$

Proof:

- (a) For $D_{\Delta}(S_1, S_2) = 0$ we need $(S_1 \setminus S_2) = \emptyset$ and $(S_2 \setminus S_1) = \emptyset$ because only the union of two empty sets gives us an empty set again. $(S_1 \setminus S_2) = \emptyset$ iff $S_1 \subseteq S_2$ and $(S_2 \setminus S_1) = \emptyset$ iff $S_2 \subseteq S_1$.

Because of $S_1 \subseteq S_2$ and $S_2 \subseteq S_1 \Rightarrow S_1 = S_2$. $D_{Hamming}(x, y) \geq 0$ you can see easily.

$$(b) \quad D_{\Delta}(S_1, S_2) = |S_1 \Delta S_2| = |(S_1 \setminus S_2) \cup (S_2 \setminus S_1)| = |(S_2 \setminus S_1) \cup (S_1 \setminus S_2)| = |S_2 \Delta S_1| = D_{\Delta}(S_2, S_1)$$

(c)

$$\begin{aligned} D_{\Delta}(S_1, S_3) &\leq D_{\Delta}(S_1, S_2) + D_{\Delta}(S_2, S_3) \\ &\Leftrightarrow \\ (S_1 \Delta S_3) &\subseteq (S_1 \Delta S_2) \cup (S_2 \Delta S_3) \\ &\Leftrightarrow \\ [(S_1 \setminus S_3) \cup (S_3 \setminus S_1)] &\subseteq [(S_1 \setminus S_2) \cup (S_2 \setminus S_1)] \cup [(S_2 \setminus S_3) \cup (S_3 \setminus S_2)] \\ &\Leftrightarrow \\ [(S_1 \setminus S_3) \cup (S_3 \setminus S_1)] &\subseteq [(S_1 \cup S_2) \setminus (S_1 \cap S_2)] \cup [(S_2 \cup S_3) \setminus (S_2 \cap S_3)] \\ &\Leftrightarrow \\ [(S_1 \cup S_3) \setminus (S_1 \cap S_3)] &\subseteq [(S_1 \cup S_2 \cup S_3) \setminus (S_1 \cap S_2 \cap S_3)] \\ &\Leftrightarrow \\ [(S_1 \setminus S_3) \cup (S_3 \setminus S_1)] &\subseteq [(S_1 \setminus S_3) \cup (S_3 \setminus S_1)] \cup [(S_2 \setminus S_3) \cup (S_3 \setminus S_2)] \\ &\Leftrightarrow \\ (S_1 \Delta S_3) &\subseteq (S_1 \Delta S_3) \cup (S_2 \Delta S_3) \\ &\Leftrightarrow \\ D_{\Delta}(S_1, S_3) &\leq D_{\Delta}(S_1, S_3) + D_{\Delta}(S_2, S_3) \\ &\Leftrightarrow \\ 0 &\leq D_{\Delta}(S_2, S_3) \end{aligned}$$

We show in (a) that $0 \leq D_{\Delta}(S_2, S_3)$ is true for any $S_2, S_3 \subseteq U$, so if you go back the equation chain $D_{\Delta}(S_1, S_3) \leq D_{\Delta}(S_1, S_2) + D_{\Delta}(S_2, S_3)$ for any $S_1, S_2, S_3 \subseteq U$ is true.

Since we can show that all conditions for a metric are fulfilled, D_{Δ} is a metric on the power set 2^U of U .

Exercise 2

(a) We have to show that for each $n \in \mathbb{N}$ with $n \geq 3$ we can find a set of P with $|P| = n$ in such a way that all cells of the Voronoi diagram are unbounded.

Let P be a set of vertices of a convex n -polygon A . So $|P| = n$. Now we take a point M inside of A . Now we can construct our Voronoi diagram in such a way that we draw n half-lines with M as the vertex for all. To draw a half-line we need another point which lays on the half-line. For the first half-line we chose the middle point of the line segment P_1 and P_2 . So we go on till we chose our last middle point between P_n and P_1 . So we construct a Voronoi diagram that is consistent with:

- each cell contains exactly one point from P

- for all $q \in \mathbb{R}^2$ that lie in the cell containing $p \in P$ it holds that $d(p, q) < d(p', q)$ for all $p' \in P \setminus \{p\}$

And all cells are unbounded because we only use half-lines.

We found this construction without any limitations, so we found a way for all n .

(b) Show: A Voronoi diagram with $n \geq 3$ has at most a number of vertices equal to $2n-5$ and the number of its edges is at most $3n-6$.

Proof:

We want to use Euler's formula, but we need to handle the half-lines. So we add a "vertex at infinity" for all unbounded edges.

Now we can use Euler's formula:

Let v be the number of Voronoi vertices (including the vertex at infinity). Let e be the number of Voronoi edges. Let n be the number of Voronoi faces.

Now

$$(v+1) - e + n = 2$$

The sum of vertex degrees is $2e$. Since each vertex has degree at least 3, we have $2e \geq 3(v+1)$, which gives a) $v \leq \frac{2}{3}e - 1$ or b) $e \leq \frac{3}{2}(v+1)$.

Plugging a) into Euler's formula, we get $2 \leq \frac{2}{3}e - e + n \Leftrightarrow e \leq 3n - 6$

Plugging b) into Euler's formula, we get $2 \leq (v+1) - \frac{3}{2}(v+1) + n \Leftrightarrow v \leq 2n - 5$

Exercise 3

For our implementation of kNN, we used Python 3.5 with scipy and numpy libraries. In the following, we will present our results shortly.

Part (2): k-NN classifier with $k = 1, 3, 5$

The accuracy results of the kNN classifier on the normalized data were as follows:

$k = 1$: 0.9072

$k = 3$: 0.8927

$k = 5$: 0.896

It is interesting that $k=1$ yields the best result, although it obtains its classification from checking only one neighbor. This suggests that there are often two samples with the same classification positioned very closely together in the feature space. (For this case, two emails with very similar content and similar classification as spam or not spam.)

Part (3): Decision Tree with depth 5

We used our previously implemented TDIDT algorithm to build a Decision Tree for the spam data. The tree was cut off at depth 5. It achieved an accuracy of 0.9078 on the test data, which is slightly better than the k-NN result for $k=1$. You can find a visualization of the Decision Tree at the bottom of the page. The most important attribute, according to the tree, is number 53. (We numbered the attributes 1-57 for this purpose.)

The Decision Tree is more easily understandable to humans, since it gives a clear structure of ordered choices everyone can comprehend. The k-NN model is a more abstract one, as its multi-feature-distances are not as intuitively interpretable.

Part (4.4): k-NN classifier, using only the attributes from the Decision Tree

For the group-dependent task, we ran the k-NN algorithm on a modified training- and test-set with only the features that appear in the Decision Tree. The accuracies were as follows:

$k = 1$: 0.8999

$k = 3$: 0.9045

$k = 5$: 0.9026

Interestingly, here the accuracy is worst for $k=1$ and best for $k=3$, which is rather what we would have expected already in (2).

