

Take-home Quiz 8

Due Date: Monday April 6, 2020 23:59

Question 1

As we discussed in class, when using the (k -)nearest-neighbor algorithm for classification, an important choice we have to make is the distance, or *metric*, function we will use. Given the space of d -dimensional vectors \mathbb{R}^d , a function $D : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a valid metric for this space if and only if it satisfies all of the following properties for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^d$:

1. *Non-negativity*: $D(\mathbf{x}, \mathbf{y}) \geq 0$.
2. *Reflexivity*: $D(\mathbf{x}, \mathbf{y}) = 0$ if and only if $\mathbf{x} = \mathbf{y}$.
3. *Symmetry*: $D(\mathbf{x}, \mathbf{y}) = D(\mathbf{y}, \mathbf{x})$.
4. *Triangle inequality*: $D(\mathbf{x}, \mathbf{y}) + D(\mathbf{y}, \mathbf{z}) \geq D(\mathbf{x}, \mathbf{z})$.

A common class of metrics are defined using the so-called *p-norms* as:

$$D_p(\mathbf{x}, \mathbf{y}) \equiv \|\mathbf{x} - \mathbf{y}\|_p = \left(\sum_{k=1}^d |x_k - y_k|^p \right)^{1/p}, \quad (1)$$

for all values of $p \geq 1$, and where x_k indicates the k -th coordinate of vector \mathbf{x} . For $p = 2$, this is the regular Euclidean distance, whereas for $p = 1$ we obtain the so-called *Manhattan distance*. Prove that, for all values of $p \geq 1$, the function D_p is a valid metric.

Hint: For metrics defined through p -norms, the triangle inequality is equivalent to the following *convexity* property: For all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ such that $\|\mathbf{x}\|_p \leq 1$ and $\|\mathbf{y}\|_p \leq 1$, and for all $0 \leq \alpha \leq 1$, it holds that $\|\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}\|_p \leq 1$. It is easier to prove the convexity property instead of the triangle inequality, for which you can use the standard form of the triangle inequality from algebra: $|x + y| \leq |x| + |y|$. You are welcome to attempt to prove the triangle inequality directly if you want, for which you can look up the *Hölder's inequality*.

Question 2

Let us assume that we use the nearest-neighbor (NN) algorithm for classification of d -dimensional vectors into one of labels $\{1, \dots, L\}$, together with the *Euclidean distance* metric and a set of labeled training data points $\mathbf{x}_n \in \mathbb{R}^d, n \in \{1, \dots, N\}$. The training data points create a segmentation of the \mathbb{R}^d space into N so-called *Voronoi cells*: For each $n \in \{1, \dots, N\}$, the

Voronoi cell $\mathbf{V}_n \subset \mathbb{R}^d$ is the subset of \mathbb{R}^d such that any point $\mathbf{y} \in \mathbf{V}_n$ has training point \mathbf{x}_n as its nearest neighbor, and thus \mathbf{y} is classified as having the same label as \mathbf{x}_n . We can “color” each such cell by the label of the training point \mathbf{x}_n it is associated with. Figure 1 shows two Voronoi segmentation examples, for the cases $d = 2, 3$ and $L = 2$. Prove that all Voronoi cells are convex: this means that, if \mathbf{y} and \mathbf{z} belong to the same Voronoi cell, $\mathbf{y}, \mathbf{z} \in \mathbf{V}_n$, then so does any point on the linear segment connecting \mathbf{y} and \mathbf{z} .

Hint: It will be helpful to consider the case of two dimensions ($d = 2$) and two training data points ($N = 2$). Under these assumptions, the problem should reduce to a simple Euclidean geometry exercise. Once you have worked out this case, you can generalize to the case of two dimensions ($d = 2$) and multiple data points ($N \geq 2$). Finally, you can extend this to the full case of any number of dimensions ($d \geq 2$) and multiple data points ($N \geq 2$).

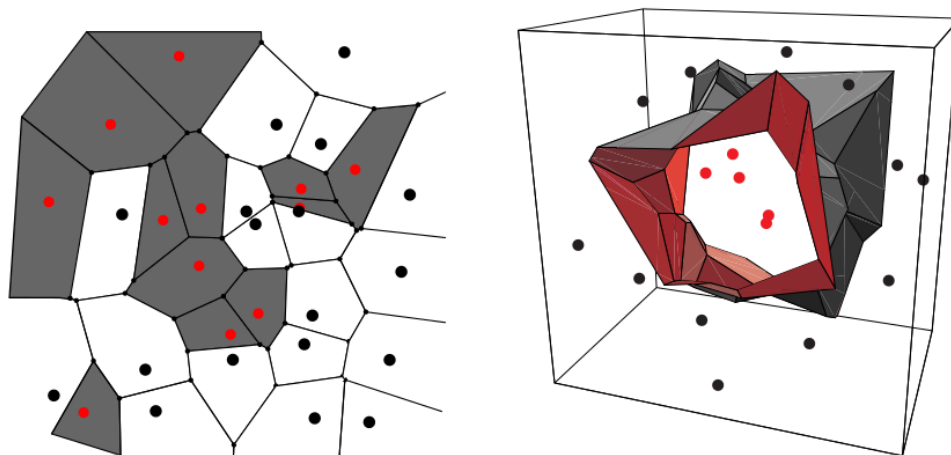


Figure 1: Example Voronoi segmentations induced by the nearest-neighbor algorithm, in two (left) and three (right) dimensions.

Instructions

1. **Integrity and collaboration:** Students are encouraged to work in groups but each student must submit their own work. If you work as a group, include the names of your collaborators in your write up. Plagiarism is strongly prohibited and may lead to failure of this course.
2. **Questions:** If you have any questions, please look at Piazza first. Other students may have encountered the same problem, and it may be solved already. If not, post your question on the discussion board. Teaching staff will respond as soon as possible.
3. **Write-up:** Your write-up should be typeset in \LaTeX and should consist of your answers

to the theory questions. Please note that we **DO NOT** accept handwritten scans for your write-up in quizzes.

4. **Submission:** Your submission for this assignment should be a PDF file, <andrew-id.pdf>, composed of your write-up. **Please do not submit ZIP files.**