# Lecture notes - Clustering and Persistence SF2704

Teacher: Wojciech Chachólski

Jim Holmström
Ariel Ekgren
Gabriel Isheden

{jimho, ekgren, ...}@kth.se

October 12, 2013

## 1  introduction

The big goal for this lecture series is to understand metrics by searching for a metric between metric spaces. This has been proven to be impossible (ref), but dispite the lack of a grand metric of metrics we can still try to search within (common word for ultrametrics, barcodes etc) to at least get a grasp on certain parts.

## 2  notes

### note 1

Skeleton for note 1

### note 2

*Metric space* $(X, d)$ where $X$ is a finite set and $d$ is a distance between the points in this finite set.

$$|X| < \infty \tag{1}$$

$$d(x,x) = 0, d(x,y) = d(y,x), d(x,y) + d(y,z) \leq d(x,z) \forall x, y, z \in X \tag{2}$$

An intresting submetric is the *ultrametric* which has a strong triangle inequality contraint

$$\max\{d(x,y), d(y,z)\} \leq d(x,z) \forall x, y, z \in X \tag{3}$$

also called ultrametric inequality.

Stated in another way, for an ultrametric any 3 points, i.e. a triangle, will have the following property

```
<image of triangle with the sides (a, a, b <= a)>
```

P $(X)$ is the set of partitions of $X$. A partition of $X$

$$\sigma = \{u_i\} = \bigsqcup_i u_i \in \mathrm{P}\ (X) \tag{4}$$

where $u_i \in \sigma$ is called a block.

To be able to order partitionings we define

$$\sigma \leq \tau \Leftrightarrow \forall u \in \sigma \exists v \in \tau : u \subset v \tag{5}$$

A *clustering* is just a function $\Psi$ in the form of a algorithm or procedure which maps a metric space $(X, d)$ to a partition of $X$

$$\Psi : (X, d) \to \mathrm{P}\ (X) \tag{6}$$

The function $\Phi(X, d) \in \mathrm{P}\ (X)$ can have 3 important properties

**Scale invariant** Changing the scale for the distance does not change the partitioning

$$\Psi(X, d) = \Psi(X, \alpha d) \forall \alpha \in \mathbb{R}_{++} \tag{7}$$

**Rich** $\Psi(X, d) \twoheadrightarrow \mathrm{P}\ (X)$ i.e. surjective or onto the set of partitions of $X$.

$$\forall \sigma \mathrm{P}\ (X) \exists d : \Psi(X, d) = \sigma \tag{8}$$

**Consistent** $d'$ is such that you decrease the intrablock distance and increase the extrablock distance for all blocks. This will correspond to making the clusters more distinct. Let $x \sim_{\Psi(X,d)} y$ denote that $x$ and $y$ belongs to the same block of the clustering $\Psi(X, d)$.

$$d'(x, y) : \begin{cases} \leq d(x, y) & x \sim_{\Psi(X,d)} y \\ \geq d(x, y) & x \nsim_{\Psi(X,d)} y \end{cases} \tag{9}$$

(not in notes; but it should be that $\Psi(X, d) = \Psi(X, d') \forall d'$ fullfilling the above property)

## note 9

Skeleton for note 9.

# 3   exercises

## exercise 1

# 4   code