LITHOLOGICAL UNCERTAINTY EXPRESSED **BY NORMALIZED COMPRESSION DISTANCE** <u>Jānis JĀTNIEKS, Tomas SAKS, Aija DĒLIŅA, Konrāds POPOVS</u> Faculty of Geography and Earth Sciences, University of Latvia, janis.jatnieks@lu.lv, puma.lu.lv

Introduction

Lithological composition and structure of the highly Quaternary deposits is heterogeneous in nature, especially as described in borehole log data.

This work aims to develop a universal solution for quantifying uncertainty based on mutual information shared between the borehole logs. Such an approach presents tangible information directly useful in generalization of the geometry and lithology of the Quaternary sediments. Such generalization can be of use in regional groundwater flow models as a qualitative estimate of lithological uncertainty involving thousands pre-processing step. of borehole logs would be humanly impossible due to the compressors, such as prediction by amount of raw data involved. Our aim is to improve partial matching (PPM), used for parametrization of recharge in the Quaternary strata. This research however holds appeal for other areas of reservoir modelling, as demonstrated in the 2011 paper by Wellmann & Regenauer-Lieb.

class lithology sandstone 2 sandstone-low conductivity

3	sandstone-high conductivity	
4	limestone	
5	limestone-low conductivity	
6	dolomite-high conductivity	
7	dolomite	
8	domerite (clayey dolomite)	
9	clay	
10	loam	
11	sandy loam	
12	silt	
13	sand	
14	gravel	
15	gypsum	
16	peat	
17	gyttja and other biogenic sediments	
18	soil	
19	other	

 Table 1.
 Main lithology
classifier table. Quaternary deposit logs consist mainly of classes 9-14,16-19.

Methods

Our calculation of borehole log similarity relies on the concept of information distance, proposed by Bennet et al. in 1998.

This was developed into a practical data mining application by Cilibrasi in the 2007 dissertation. The resulting implementation called CompLearn utilities provide a calculation of the Normalized Compression Distance (NCD) metric.

 $NCD(x,y) = \frac{C(xy) - \min\{C(x), C(y)\}}{\max\{C(x), C(y)\}}$

The **C** in this formula denotes a length of output by a universal data compression program such as zlib for generating the popular zip archives, bzip2, 7zip, PPMd or, in principle, any other data compression library. The NCD is dissimilarity metric, expressing dissimilarity in range between 0 and 1 (Figure 3).

Data

For our experiments we used extracts of the Quaternary strata from general-purpose geological borehole log database, maintained by the Latvian Environment, Geology and Meteorology Centre, spanning the territory of Latvia (Takčidi 1998). Lithological codes were manually generalised into 20 rock types using a lithology classifier for all lithology codes in the database (Table 1).

For generating the results in this poster presentation, 3210 borehole logs were used. These logs had to have at least 4 distinct lithology layers, according to classifier in Table 1 and pass through Quaternary deposit cover.

The universal data compression algorithms, used in this way, estimate the mutual information content in the complex and data. This approach has proven to be universally successful for parameter free data mining in disciplines ranging from molecular biology, handwriting recognition, creation of language trees and a multitude of other surprisingly different applications (Cilibrasi, 2007).

To improve this approach for use in geology, it is beneficial to apply a transformation of borehole log data as a stream Text computing the NCD metric in our experiments, is highly dependent on context. We assign unique symbols for aggregate lithology types and serialize the borehole logs into text strings (such as in Fig. 2), where the string length represents normalized borehole depth (Fig. 1). This encoding ensures that lithology types as well as thickness structure and sequence of strata is comparable in a form, most native to the universal data compression software, that calculates the pairwise NCD dissimilarity matrix (Fig. 3).



Bedrock Fig. 1. An example of borehole log serialization transformation for a synthetic borehole with 3 layer transitions and 4 layers with different aggregate lithology types.

Borehole ID	Depth, m	Dept to t Quat
11390	8	
11390	12	
11390	15	
11390	24	
11390	28.8	
11390	56	
11390	58.8	
11390	71	
11390	78.7	
11390	79.7	

Table 2. Sample borehole log with lithology classes and depth, normalized to Quaternary thickness.

Fig. 2. Normalized borhole log from Table 2 in serialized text form. This encodes the thickness, sequence and lithology type of layers in a form well suited for text compression programs such as PPMd (Shkarin 2002) that are used for calculating the NCD metric.



Heatmap of Normalized Compression Distance matrix calculated using Shkarin's **PPMd** compression algorithm on 3120 serialized borehole logs.

The NCD results can be used for studying the structure of Quaternary sediments from the perspective of similarity, according to universal entropy coding copression algorithms (Cilibrasi 2007).

lithological structure in hydrogeological modelling as it allows for used complete-link hierarchical We minimization of uncertainty in the Quaternary strata. clustering, as implemented in C Clustering Output for delineation of regions of similar lithological structure library by Hoon et al. on two experiment **minimizes uncertainty** by delineating regions with quantifiably more matrices – NCD and NCD + range similar Quaternary lithological structure using Voronoi tesselation and normalized Euclidean matrix of borehole NCD as a measure of dissimilarity. locations (NCD+E). The differene of spatial This work paves way for practical application of NCD metric to basin clustering results is shown in Figure 6 and modelling, where these results can be used as input in hydrogeological structure of NCD+E clustering shown in models for building a better representation of Quaternary lithological dendrogram in Figure 5. structure.



th normalized Lithology thickness of Class ternary cover 0.1 sand 0.151 sand 0.188 sil 0.301 c 0.361 silt 0.703 loam 0.738 gravel 0.891 0.987 gravel



Fig. 4 (left). Spatial context of modelling territory. Spatial cluster solutions in Figure 6 cover the territory of Latvia.



Results

The current implementation provides cluster membership information for all boreholes in clustering solution as well as numerical determination of the most representative borehole section for each cluster. This information can be used for generalization of 3D

Acknowledgments. This work was supported by the European Social Fund project "Establishment of interdisciplinary scientist group and modelling system for ground-water research",

2009/0212/1DP/1.1.1.2.0/09/APIA/VIAA/060



Fig. 5. Complete-link agglomerative hierarhical clustering dendrogram of Normalized Compression Distance matrix + range normalized borehole Euclidean distance matrix.

Fig. 6 (left & right) Examples of spatial clustering solutions in the territory of Latvia, created dissolving by Voronoi adjacent polygons with matching membership generated from borehole locations. Random colors.

Left - Spatial clusters, created from flattening the hierarhical clustering solution into 4,6 and 12 logical clusers, using the NCD matrix distances.

Right - Spatial clusters, created from 4,6 and 12 clusters, using the NCD matrix distance with normalized sum Euclidean distances borehole between locations. EPSG:25884 coordinate projected system. Dendrogram in Fig. 5. Note that there are more spatial clusters than logical clusters in the dendrogram.



References

- 1. Bennett, C. H., Gacs P., Li M., Vitanyi P., Zurek W. 1998, Information Distance, IEEE Transactions on
- Information Theory, 44(4), 1407-1423., IEEE. 2. Cilibrasi, R. 2007., Statistical Inference Through Data Compression, ILLC Dissertation Series DS 2007-01, Institute for Logic, Language and Computation, Universiteit van Amsterdam.
- 3. Hoon M., Imoto, S., Miyano, S., 2010, The C Clustering Library, The University of Tokyo, Institute of Medical Science, Human Genome Center.
- 4. Shkarin D., 2002, PPM: one step to practicality Proceedings of the Data Compression Conference 2002. IEEE
- 5. Takčidi, E. 1999. Datu bāzes "Urbumi" dokumentācija [Documentation of the database "Boreholes"]. Valsts ģeoloģijas dienests, Rīga. [In Latvian].
- Wellmann J.F., Regenauer-Lieb K., 2011 Uncertainties have a meaning: Information entropy as a quality measure for 3-D geological models, Tectonophysics, Elsevier (in press).