

### friangular inequality holds for metrics



"as the crow flies" distance Example: Euclidean or

 $d(x,z) \le d(x,y) + d(y,z)$ 

Thursday 31 January 13 Horizontal 5



### <u> Itrametric</u>

- Euclidean distance makes a lot of sense when the population is homogeneous
- Ultrametric distance makes a lot of sense when the observables are heterogeneous, discontinuous
- Latter is especially useful for determining: anomalous, atypical, innovative cases

Thursday 31 January 13





## Analysis of semantics: 2. Hierarchy tracks anomaly and change

- Euclidean distance makes a lot of sense when the population is homogeneous
- Ultrametric distance makes a lot of sense when the observables are heterogeneous, discontinuous
- Latter is especially useful for determining: anomalous, atypical, innovative cases

### Analysis of semantics: |. Context - the collection of all interrelationships

- Euclidean distance makes a lot of sense when the population is homogeneous
- All interrelationships together provide context, relativities - and meaning

Thursday 31 January 13

## What is so special about hierarchy?

- Ultrametric spaces have interesting properties.
- Not just in data analysis and pattern recognition, but in physics at small scales, and in optimization.
- Ultrametric topology and p-adic number systems are closely associated.
- Next I will look at:
- (i) Quantifying inherent ultrametricity
- (ii) Computational implications.



ntonetiy ≙ Springe 12

Thursday 31 January 13	Thursday 31 January 13
10	15
<ul> <li>Relationship between subdominant ultrametric, and given dissimilarities.</li> <li>Rammal, Toulouse and Virasoro, Ultrametricity for physicists, Rev. Mod. Phys., 58, 765-788, 1986.</li> <li>Whether interval between median and max rank dissimilarity of every set of triplets is nearly empty. (Taking ranks provides scale invariance.)</li> <li>Lerman, Classification et Analyse Ordinale des Données, Dunod, 1981.</li> </ul>	<ul> <li>So: we take all possible triplets, i, j, k</li> <li>We look at their angles, and judge whether or not the ultrametric triangle properties are verified</li> <li>If so: #UM-triangles++</li> <li>Having examined all possible triangles, our α measure is: #UM-triangles / #triangles</li> <li>All triangles respect these ultrametric properties implies α = 1; no triangle does, then = 0</li> <li>For n objects, this is computationally prohibitive, so we sample i,j,k in practice (uniformly)</li> </ul>
Other Ways of Quantifying Ultrametricity – III	Quantifying ultrametricity – II
<ul> <li>Assume Euclidean space. Consider a triplet of points, that defines a triangle.</li> <li>Take smallest internal angle, a, in triangle ≤ 60 deg.</li> <li> and, for the two other internal angles, b and c, i   b - c   &lt; 2 deg. (arbitrary small angle),</li> <li>Then this triangle is ultrametric.</li> <li>We look for the overall proportion of such triangles in our data.</li> </ul>	<ul> <li>The distance between two objects or two terminals in the tree is the lowest rank which dominates them. Lowest or closest common ancestor distance.</li> <li>The ultrametric inequality holds for any 3 points (or terminals): <ul> <li>d(i, k) ≤ max {d(i,j), d(j,k)}</li> <li>Recall: the triangular inequality is: d(i,k) ≤ {d(i,j) + d(j,k)}</li> </ul> </li> <li>An ultrametric space is quite special: (i) all triangles are isosceles with small base, or equilateral; (ii) every point in a ball is its center; (iii) the radius of a ball equals the diameter; (iv) a ball is clopen; (v) an ultrametric space is always topologically 0-dimensional.</li> </ul>
Quantifying ultrametricity – I	Some Properties of Ultrametrics

Thursday 31 January 13	– Then I will move to narrative analysis and synthesis.	– References.	-Hierarchical clustering via Baire distance using chemical compounds.	– Hierarchical clustering via Baire distance using SDSS (Sloan Digital Sky Survey) spectroscopic data.	<ul> <li>Clustering of large data sets.</li> </ul>	<ul> <li>Ultrametric topology, Baire distance.</li> </ul>	<ul> <li>First, agglomerative hierarchical clustering; then: "hierarchical encoding" of data.</li> </ul>	Applications in Search and Discovery	Thursday 31 January 13	17	<ul> <li>See: F Murtagh, "On ultrametricity, data coding, and computation", Journal of Classification, 21, 167-184, 2004</li> <li>Hall, P., Marron, J.S., and Neeman, A., "Geometric representation of high dimension low sample size data", JRSS B, 67, 427-444, 2005</li> <li>F. Delon, Espaces ultramétriques, J. Symbolic Logic, 49, 405-502, 1984</li> </ul>	<ul> <li>As dimensionality increases, so does ultrametricity.</li> <li>In very high dimensional spaces, the ultrametricity approaches being 100%.</li> <li>Relative density is important: high dimensional and spatially sparse mean the same in this context.</li> </ul>	Pervasive Ultrametricity
Thursday 31 January 13	20							Next: the Baire (ultra)metric	Thursday 31 January 13	18	<ul> <li>Adding a new terminal to a dendrogram is carried out in O(t) time.</li> <li>Cost of finding the ultrametric distance between two terminal nodes is twice the length of a traversal from root to terminals in the dendrogram. Therefore distance is computed in O(t) time.</li> <li>Nearest neighbor search in ultrametric space can be carried out in O(1) or constant time.</li> </ul>	<ul> <li>Consider a dendrogram: a rooted, labeled, ranked, binary tree. So: <i>n</i> terminals, <i>n</i>-1 levels.</li> <li>A dendrogram's root-to-terminal path length is <i>log<sub>2</sub>n</i> for a balanced tree, and <i>n</i>-1 for an imbalanced tree. Call the computational cost of such a traversal <i>O(t)</i> where <i>t</i> is this path length. It holds: 1 ≥ O(t) ≥ n-1.</li> </ul>	<b>Computational Implications</b>



0.13045037	0.120695	0.4634/806	145./4324
0 06343859	61604920 0	0 12690687	145 7ን943
0.19580211	0.16425499	0.53404553	145.73267
0.18679948	0.15610801	0.50961215	145.64568
0.41157582	0.46691701	0.63385916	145.6607
0.17476539	0.145909	0.53370196	145.42139
0.15175095	0.14611299	0.56416792	145.4339
phot. redshift	spec. redshift	DEC	RA
	- example	Data -	
			Thursday 31 January 13
	÷	ous onier datasets	(טרטט), מווע אמו ו
loan Digital Sky Survey	edshifts from the SI	nd photometric re	spectrometric a
hemical compounds,	distance to: c	ed the Baire	– We appli
shing scheme.)	ce: hierarchical has	the data. (Henc	linear scan of
ctly read from a	archy can be direc	ermore the hiera	with it. Furthe
ships associated	esent the relations	be used to repre	hierarchy can
It follows that a	ametric distance	stance is an ultr	– The Baire di
air of sequences.	efix, the closer a p	the common pre	– The longer 1
	O	reprinted 2002)]	Theory, Dover, 1979 (
n prefix IA. Lew. Basic Set	e longest commo	d in terms of the	metric defined
sequences with a	countable infinite	e consists of c	– Baire spac
	(עונו מ/וו		
· · · · · · · · · · · · · · · · · · ·	/ltra/m	N Duiro	) 5 +

Thursday 31 January 13	Thursday 31 January 13
	• I.e. numbers of observations sharing 6, 5, 4, 3 decimal digits.
dimensions?	<ul> <li>21.7% of z_spec and z_phot have at least 3 common prefix digits.</li> </ul>
<ul> <li>Since we are using digits of precision in our data (re)coding how do we handle high</li> </ul>	<ul> <li>We can find very efficiently where these</li> <li>82.8% of the astronomical objects are.</li> </ul>
chemoinformatics - which is high dimensional.	• I.e. numbers of observations sharing 6, 5, 4, 3, 2 decimal digits.
<ul> <li>Next - another case study, using</li> </ul>	<ul> <li>82.8% of z_spec and z_phot have at least 2 common prefix digits.</li> </ul>
	Framework for Fast Clusterwise Regression
Thursday 31 January 13	Thursday 31 January 13
On the left we have z_spec where three data peaks can be observed. On the right we have z_phot where only one data peak can be seen.	<ul> <li>Note: cluster-wise not spatially (RA, Dec) but rather within the data itself</li> </ul>
star. 70.000. 355.00. rolet 1.0000. 1.0000	<ul> <li>I.e., cluster-wise nearest neighbour regression</li> </ul>
	<ul> <li>Furthermore: determine good quality mappings of z_spect onto z_phot, and less qood quality mappings</li> </ul>
0.00 0.00 0.05 0.05	<ul> <li>Motivation - regress z_spect on z_phot</li> </ul>
Perspective Plots of Digit Distributions	

0.05 0.045 0.045 0.025 0.025 0.025 0.025 0.025 0.025 0.025



## Matching of Chemical Structures

- Clustering of compounds based on chemical descriptors or chemical representations, in the pharmaceutical industry.
- Used for screening large corporate databases
- Chemical warehouses are expanding due to mergers, about by combinatorial chemistry. acquisitions, and the synthetic explosion brought

Thursday 31 January 13

З

# Chemoinformatics clustering

- 1.2 million chemical compounds, each characterized by 1052 boolean presence/absence values.
- Firstly we note that precision of measurement leads to greater ultrametricity (i.e. the data are more hierarchical).
- From this we develop an algorithm for finding equivalence "condensation". classes of specified precision chemicals. We call this: data
- Secondly, we use random projections of the 1052to k-means clustering outcomes. dimensional space in order to find the Baire hierarchy. We find that clusters derived from this hierarchy are quite similar

Thursday 31 January 13	Thursday 31 January 13
<ul> <li>Text synthesis - supporting collaborative narrative construction - book writing</li> </ul>	
<ul> <li>Characterizing structure and properties of CSI television episodes</li> </ul>	<ul> <li>A story is an expression of causality or connection</li> <li>Narrative connects facts or views or other units of information</li> </ul>
<ul> <li>lext attributes and their significance (and feasibility of mapping these onto desired outputs)</li> </ul>	<ul> <li>Narrative suggests a causal or emotional relationship between events.</li> </ul>
<ul> <li>Associate these, resp., with Euclidean metric and ultrametric</li> </ul>	<ul> <li>Semantics include time evolution of structures and patterns, including both: threads and commonality; and change, the exceptional, the anomalous.</li> </ul>
<ul> <li>Casablanca - analysis of emotion in scene 43</li> <li>Two senses of semantics - all interrelationships, and change over time</li> </ul>	<ul> <li>Semantics include web of relationships - thematic structures and patterns. Structures and interrelationships evolve in time.</li> </ul>
	<ul> <li>We must consider complex web of relationships.</li> </ul>
Topics	Analysis of Narrative Technical Issues Addressed
Thursday 31 January 13	Thursday 31 January 13
	S
- At issue throughout this work: embedding of our data in an ultrametric topole	<ul> <li>Second approach: use random projections of the high dimensional data, and then use the Baire distance.</li> </ul>
- The Baire method - we find - offers a fast alternative to k-means and a fortion traditional agglomerative hierarchical clustering	of these crude clusters. We call this "data condensation". For 20000 compounds, 1052 attributes, a few mins. needed in R.
- We are targeting clustering in massive data sets	read off equivalence classes of 0-distance compounds, with
- We are hashing, in a hierarchical or multiscale way, our data	errect of more compound values becoming the same for a given attribute. Through a heuristic (e.g. interval of row sum values),
- We obtain a hierarchy that can be visualized as a tree	• Limit precision of compound / attribute values. This has the
- Alternative viewpoint: we can cluster information based on the longest commprefix	<ul> <li>sum (hence "profile" in Correspondence Analysis terms).</li> <li>Two clustering approaches studied:</li> </ul>
- First viewpoint: encode the data hierarchically and essentially read off the clust	<ul> <li>Normalize chemical compounds by dividing each row by row</li> </ul>
- We have a new way of inducing a hierarchy on data	
Summary Remarks on Search and Discovery	
· · ·	



#### subdivid Ana ided lysis of Casablanca's "Mii into Climax", Scene 43 "beats" (subscenes d-Act

- McKee divides this scene, relating to Ilsa and Rick seeking black market exit visas, into 11 "beats"
- Beat I is Rick finding Ilsa in the market

•

- Beats 2, 3, 4 are rejections of him by Ilsa
- Beats 5, 6 express rapprochement by both
- Beat 7 is guilt-tripping by each in turn
- Beat 8 is a jump in content: Ilsa says she will leave Casablanca soon
- In beat 9, Rick calls her a coward, and Ilsa calls him a fool
- In beat 10, Rick propositions her
- In beat 11, the climax, all goes to rack and ruin: Ilsa says she was married to Laszlo all along. Rick is stunned

#### Thursday 31 January 13



#### between May and shot by Warner August 1942

You will not find a treasure like this in all Morocco, Mademoiselle Only seven hundred francs. At the linen stall, IIsa examines a tablecloth which an Arab vendor is endeavoring to sell. He holds a sign which reads "700 francs."

ARAB

Rick walks up behind Ilsa.

EXT. BLACK MARKET - DAY

Casablanca

Movie

Brothers









- Scene 43 in Casablanca (out of 77 scenes).
- analyze 11 subscenes ("beats"). Crucial mid-point scene. Following McKee, I will
- Right, first three subscenes (in blue, brown, red).

Thursday 31 January 13



ILSA It doesn't matter, thank you.

She looks briefly at Rick, then politely formal. RICK You're being cheated.

away. Her manner



"virtually perfect".

RAB Ah, the lady is a friend of Rick & For friends of Rick we have a small discourt. Did I say seven hundred france? You can have it for

Reaching under the counter, he takes out a sign reading "200 francs", and replaces the other sign with it.

condition to receive you when you called

ILSA It doesn't matter RICK I'm sorry I was ir on me last night.

ARAB Ah, for special friends of Rick's we have a special discount. One hundred francs.

He replaces the second sign with a third which reads "100 francs."







Thursday 31 January 13



- Back to a deeper look at Casablanca
- We have taken comprehensive but qualitative discussion by McKee and sought qualitative and algorithmic implementation

### McKee, Methuen, 1999

Casablance is based on a range of miniplots.

"virtually perfect" composition is McKee: its

"sensory surface" of the underlying Text is the semantics



Thursday 31 January 13

## Our way of analyzing semantics

- We discern story semantics arising out of the orientation of narrative
- This is based on the web of interrelationships
- We examined caesuras and breakpoints in the flow of narrative
- With CSI scripts: characterization

Style analysis of scene 43 based on McKee Monte Carlo tested against 999 uniformly randomized sets of the beats

- In the great majority of cases (against 83% we find the style in scene 43 to be and more of the randomized alternatives) characterized by:
- small variability of movement from one beat to the next
- greater tempo of beats
- high mean rhythm

Thursday 31 January 13

# CSI: Crime Scene Investigation

ranscripts of 3 episodes first aired by CBS Oct. 2000



# Support environment for collaborative, distributed creating of narrative

- Pinpointing anomalous sections
- Assessing homogeneity of style over successive iterations of the work
- Scenario experimentation and planning
- This includes condensing parts, or elaborating
- Similarity of structure relative to best practice in chosen genre

## **Text Synthesis**

- Aristotle's Poetics (c. 350 BC)
- "Outlines and episodization" "Stories ... should first be set out in universal terms ... on that basis, one should then turn the story into episodes and elaborate it."
- "... reasoning is the speech which the agents use to argue a case or put forwards an opinion"





The Delivery

A Y I N G

Hierarchy, as well as geometry (Euclidean factor space as in Correspondence Analysis) for both understanding and working in complex systems. In this presentation: applications to search and discovery, information retrieval

### References

- F. Murtagh, Correspondence Analysis and Data Coding with R and Java, Chapman and Hall/CRC Press, 2005. See chapter 5 on text analysis. Software at www.correspondances.info
- F. Murtagh, A. Ganz, S. McKie, J. Mothe and K. Englmeier, "Tag Clouds for Displaying Semantics: The Case of Filmscripts", Information Visualization Journal, forthcoming. 9, 253-262, 2010.
- F. Murtagh, "The Correspondence Analysis platform for uncovering deep structure in data and information", Sixth Boole Lecture, Computer Journal, 53, 304-315, 2010.
- F. Murtagh, A. Ganz and S. McKie, "The structure of narrative: the case of film scripts", Pattern Recognition, 42, 302-312, 2009. (See discussion in Z. Merali, "Here's looking at you, kid. Software promises to identify blockbuster scripts.", Nature, 453, p. 708, 4 June 2008.)
- F. Murtagh, A. Ganz and J. Reddington, "New methods of analysis of narrative and semantics in support of interactivity", Entertainment Computing, 2, 115-121 2011.

#### Thursday 31 January 13

### References

F. Murtagh, "Thinking ultrametrically", in D. Banks, L. House, F.R. McMorris, P. Arabie and W. Gaul, Eds., Classification, Clustering, and Data Mining Applications, Springer, 3-14, 2004.

 F. Murtagh, "On ultrametricity, data coding, and computation", Journal of Classification, 21, 167-184, 2004.

 F. Murtagh, "Identifying the ultrametricity of time series", European Physical Journal B, 43, 573-579, 2005.

– F. Murtagh, G. Downs and P. Contreras, "Hierarchical clustering of massive, high dimensional data sets by exploiting ultrametric embedding". SIAM Jnl. on Scientific Computing, Vol. 30, No. 2, pp. 707–730. February 2008.

 P. Contreras and F. Murtagh. "Fast, linear time hierarchical clustering using the Baire metric". Journal of Classification, 29, 118-143, 2012.

– F. Murtagh and P. Contreras, "Fast, linear time, m-adic hierarchical clustering for search and retrieval using the Baire metric, with linkages to generalized ultrametrics, hashing, formal concept analysis, and precision of data measurement", p-Adic Numbers, Ultrametric Analysis and Applications, 4, 45-56, 2012.

 – P. Contreras and F. Murtagh, "Linear time Baire hierarchical clustering for enterprise information retrieval", International Journal of Software and Informatics, 6, 363-380, 2012.