

Andreas Holzinger
VO 709.049 Medical Informatics
09.12.2015 11:15-12:45

Lecture 09

Interactive Information Visualization and Visual Analytics

a.holzinger@tugraz.at
Tutor: markus.plass@student.tugraz.at
<http://hci-kdd.org/biomedical-informatics-big-data>



A. Holzinger 709.049 1/98 Med Informatics L09

Schedule

- 1. Intro: Computer Science meets Life Sciences, challenges, future directions
- 2. Back to the future: Fundamentals of Data, Information and Knowledge
- 3. Structured Data: Coding, Classification (ICD, SNOMED, MeSH, UMLS)
- 4. Biomedical Databases: Acquisition, Storage, Information Retrieval and Use
- 5. Semi structured and weakly structured data (structural homologies)
- 6. Multimedia Data Mining and Knowledge Discovery
- 7. Knowledge and Decision: Cognitive Science & Human-Computer Interaction
- 8. Biomedical Decision Making: Reasoning and Decision Support
- 9. Intelligent Information Visualization and Visual Analytics
- 10. Biomedical Information Systems and Medical Knowledge Management
- 11. Biomedical Data: Privacy, Safety and Security
- 12. Methodology for Info Systems: System Design, Usability & Evaluation

A. Holzinger 709.049

2/98

Med Informatics L09

Keywords of the 9th Lecture 

- Data visualization
- Flow cytometry
- Human-Computer Interaction (HCI)
- Information visualization
- Interactive information visualization
- k-Anonymization
- Longitudinal data
- Multivariate data
- Parallel coordinates
- RadViz
- Semiotics
- Star plots
- Temporal data analysis
- Visual analytics
- Visual information

A. Holzinger 709.049 3/98 Med Informatics L09

Advance Organizer (1/2)

- **Biological data visualization** = as branch of bioinformatics concerned with visualization of sequences, genomes, alignments, phylogenies, macromolecular structures, systems biology, etc.
- **Clustering** = Mapping objects into disjoint subsets to let appear similar objects in the same subset;
- **Data visualization** = visual representation of complex data, to communicate information clearly and effectively, making data useful and usable;
- **Information visualization** = the interdisciplinary study of the visual representation of large-scale collections of non-numerical data, such as files and software, databases, networks etc., to allow users to see, explore, and understand information at once;
- **Multidimensional scaling** = Mapping objects into a low-dimensional space (plane, cube etc.) in order to let appear similar objects close to each other;
- **Multi-Dimensionality** = containing more than three dimensions and data are multivariate;
- **multivariate** = encompassing the simultaneous observation and analysis of more than one statistical variable; (Antonym: univariate = one-dimensional);

A. Holzinger 709.049

4/98

Med Informatics L09

Advance Organizer (2/2) 

- **Parallel Coordinates** = for visualizing high-dimensional and multivariate data in the form of N parallel lines, where a data point in the n-dimensional space is transferred to a polyline with vertices on the parallel axes;
- **RadViz** = radial visualization method, which maps a set of m-dimensional points in the 2-D space, similar to Hooke's law in mechanics;
- **Semiotics** = deals with the relationship between symbology and language, pragmatics and linguistics. Information and Communication Technology deals not only in words and pictures but also in ideas and symbology;
- **Semiotic engineering** = a process of creating a semiotic system, i.e. a model of human intelligence and knowledge and the logic for communication and cognition;
- **Star Plot** = aka radar chart, spider web diagram, star chart, polygon plot, polar chart, or Kiviat diagram, for displaying multivariate data in the form of a two-dimensional chart of three or more quantitative variables represented on axes starting from the same point;
- **Visual Analytics** = focuses on analytical reasoning of complex data facilitated by interactive visual interfaces;
- **Visualization** = a method of computer science to transform the symbolic into the geometric, to form a mental model and foster unexpected insights;

A. Holzinger 709.049 5/98 Med Informatics L09

Learning Goals: At the end of this 9th lecture you ...

- ... have some background on visualization, visual analytics and content analytics;
- ... got an overview about various possible visualization methods for multivariate data;
- ... got an introduction into the work of and possibilities with parallel coordinates;
- ... have seen the principles of RadViz mappings and algorithms;
- ... are aware of the possibilities of Star Plots;
- ... have seen that visual analytics is intelligent Human-Computer Interaction at its finest;

A. Holzinger 709.049

6/98

Med Informatics L09

Slide 9-1 Key Challenges

TU Graz

- How to understand high-dimensional spaces?
- The transformation of results from high-dimensional space \mathbb{R}^N into \mathbb{R}^2
- From the complex to the simple
- Low integration of visual analytics techniques into the clinical workplace
- Sampling, modelling, rendering, perception, cognition, decision making
- Trade-off between time and accuracy
- How to model uncertainty

A. Holzinger 709.049 7/98 Med Informatics L09

Visualization is an essential part of Data Science

TU Graz

Interactive **Data Mining** **Knowledge Discovery**

HCI **KDD**

6 Data Visualization 2 Machine Learning 1 Data Preprocessing Data Fusion

3 Graph-based Data Mining 4 Topological Data Mining 5 Entropy-based Data Mining

Privacy, Data Protection, Safety and Security 7

A. Holzinger 709.049 8/98 Med Informatics L09

TU Graz

Verbal Information versus Visual Information

A. Holzinger 709.049 9/98 Med Informatics L09

TU Graz

Problem: Context!

A. Holzinger 709.049 10/98 Med Informatics L09

Semantic Ambiguity – Missing Context

TU Graz

Cell Energizer battery Car frame Man in jail Blood sample

Radio Mast Transceiver Radio Mast Transceiver Radio Mast Transceiver Radio Mast Transceiver

11/98 Med Informatics L09

TU Graz

A picture is worth a thousand words?

A. Holzinger 709.049 12/98 Med Informatics L09

Slide 9-7: Example: Ribbon Diagram of a Protein Structure

Magnani, R., et al. 2010. Calmodulin methyltransferase is an evolutionarily conserved enzyme that trimethylates Lys-115 in calmodulin. *Nature Communications*, 1, 43.

A. Holzinger 709.049 13/98 Med Informatics L09

Slide 9-8 "Is a picture really worth a thousand words?"

A. Holzinger 709.049 14/98 Med Informatics L09

Informatics as Semiotics Engineering

A. Holzinger 709.049 15/98 Med Informatics L09

Slide 9-9 Three examples for Visual Languages

Ware, C. (2004) *Information Visualization: Perception for Design (Interactive Technologies)* 2nd Edition. San Francisco, Morgan Kaufmann.

Holzinger, A., Searle, G., Auinger, A. & Ziefle, M. (2011) Informatics as Semiotics Engineering: Lessons learned from Design, Development and Evaluation of Ambient Assisted Living Applications for Elderly People. *Universal Access in Human-Computer Interaction. Context Diversity. Lecture Notes in Computer Science (LNCS 6767)*. Berlin, Heidelberg, New York, Springer, 183-192.

A. Holzinger 709.049 16/98 Med Informatics L09

Slide 9-10 Informatics as Semiotics Engineering

- 1. Physical: is it present?
 - Signals, traces, components, points, ...
- 2. Empirical: can it be seen?
 - Patterns, entropy, codes, ...
- 3. Syntactic: can it be read?
 - Formal structure, logic, deduction, ...
- 4. Semantic: can it be understood?
 - Meaning, proposition, truth, ...
- 5. Pragmatic: is it useful?
 - Intentions, negotiations, communications, ...
- 6. Social: can it be trusted?
 - Beliefs, expectations, culture, ...

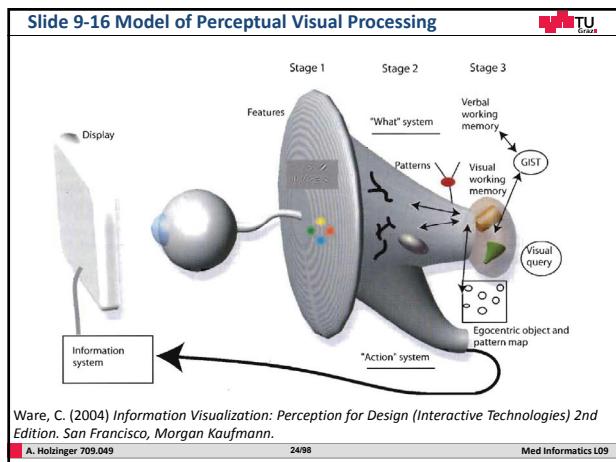
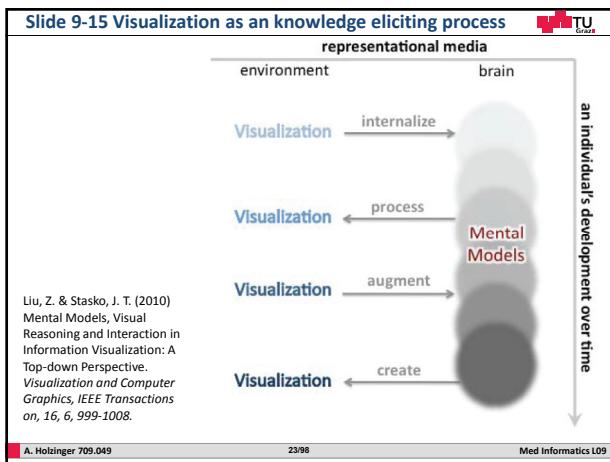
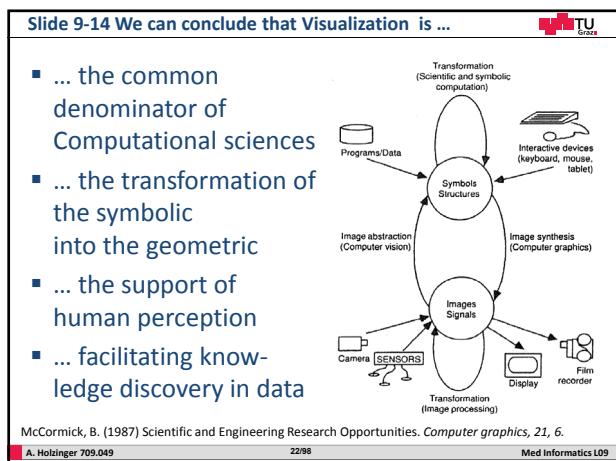
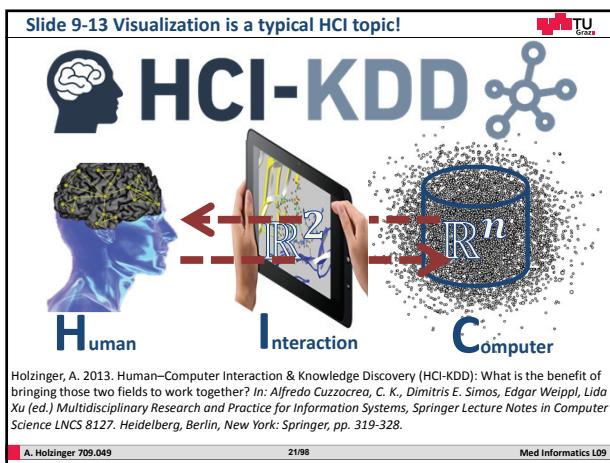
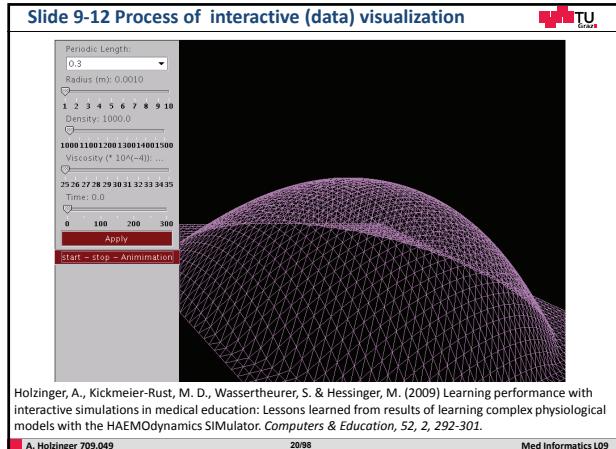
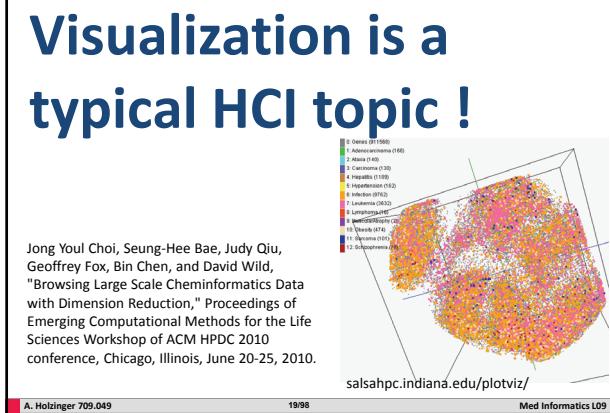
Burton-Jones, A., Storey, V. C., Sugumaran, V. & Ahluwalia, P. 2005. A semiotic metrics suite for assessing the quality of ontologies. *Data & Knowledge Engineering*, 55, (1), 84-102.

A. Holzinger 709.049 17/98 Med Informatics L09

Slide 9-11 Definitions of the term "Visualization"

- **Visualization** = generally a method of computer science to transform the symbolic into the geometric, to form a mental model and foster unexpected insights;
- **Information visualization** = the interdisciplinary study of the visual representation of large-scale collections of non-numerical data, such as files and software, databases, networks etc., to allow users to see, explore, and understand information at once;
- **Data visualization** = visual representation of complex data, to communicate information clearly and effectively, making data useful and usable;
- **Visual Analytics** = focuses on analytical reasoning of complex data facilitated by interactive visual interfaces;
- **Content Analytics** = a general term addressing so-called "unstructured" information – mainly text – by using mixed methods from visual analytics and business intelligence;

A. Holzinger 709.049 18/98 Med Informatics L09



Usefulness of Visualization Science

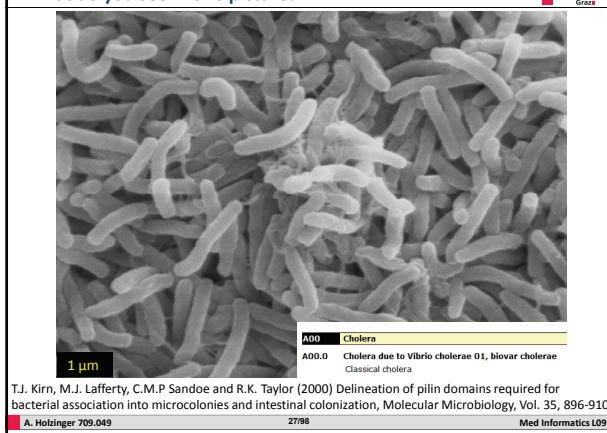
A. Holzinger 709.049 25/98 Med Informatics L09

Slide 9-17 A look back into history ...



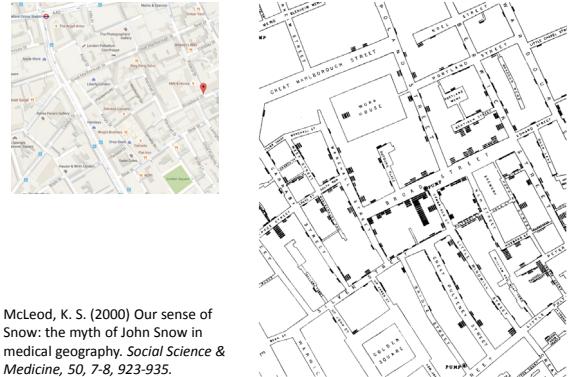
A. Holzinger 709.049 26/98 Med Informatics L09

What do you see in this picture?



A. Holzinger 709.049 27/98 Med Informatics L09

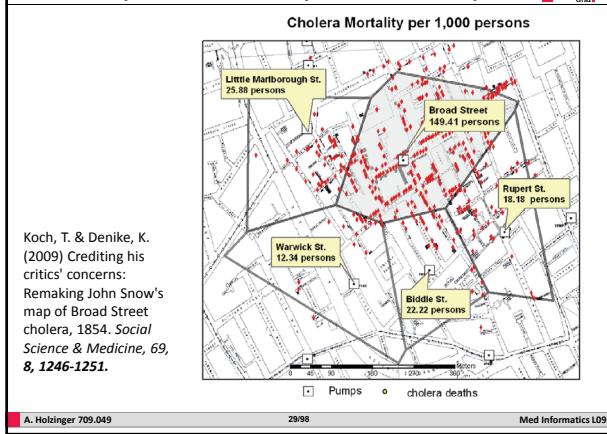
Slide 9-18 Medical Visualization by John Snow (1854)



McLeod, K. S. (2000) Our sense of Snow: the myth of John Snow in medical geography. *Social Science & Medicine*, 50, 7-8, 923-935.

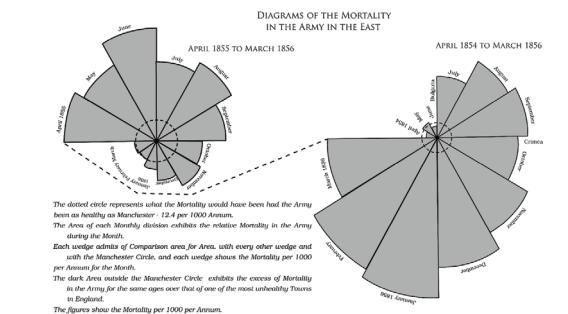
A. Holzinger 709.049 28/98 Med Informatics L09

Slide 9-19 Systematic Visual Analytics > Content Analytics



A. Holzinger 709.049 29/98 Med Informatics L09

Florence Nightingale – first medical quality manager



Meyer, B. C. & Bishop, D. S. (2007) Florence Nightingale: nineteenth century apostle of quality. *Journal of Management History*, 13, 3, 240-254.

A. Holzinger 709.049 30/98 Med Informatics L09

How many visualization methods do exist?

A. Holzinger 709.049

31/98

Med Informatics LOS

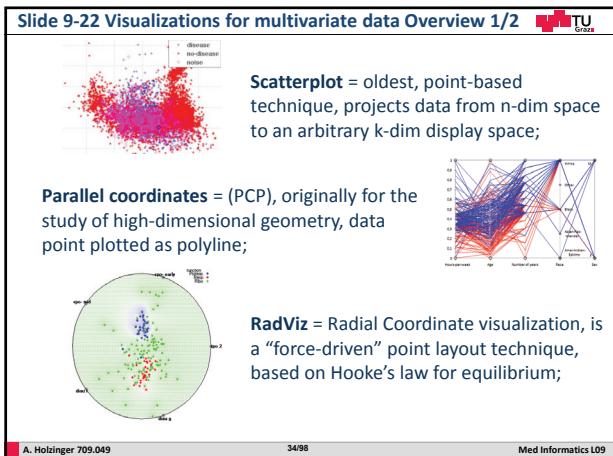
Slide 9-21: A taxonomy of Visualization Methods

- 1) Data Visualization (Pie Charts, Area Charts or Line Graphs, ...)
 - 2) Information Visualization (Semantic networks, tree-maps, radar-chart, ...)
 - 3) Concept Visualization (Concept map, Gantt chart, PERT diagram, ...)
 - 3) Metaphor Visualization (Metro maps, story template, iceberg, ...)
 - 4) Strategy Visualization (Strategy Canvas, roadmap, morpho box,...)
 - 5) Compound Visualization

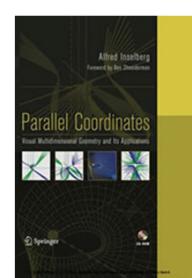
A. Holzinger 709.049

33/98

Med Informatics L09



Parallel Coordinates



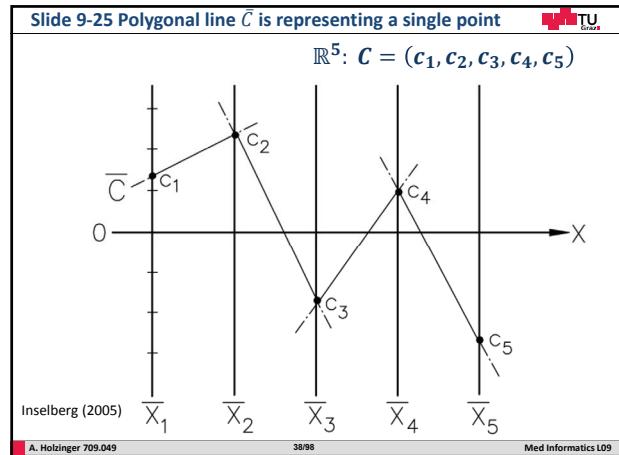
Slide 9-24 Parallel Coordinates – multidim. Visualization



- On the plane with Cartesian-coords, a vertical line, labeled \bar{X}_i is placed at each $x = i - 1$ for $i = 1, 2, \dots, N$.
- These are the axes of the parallel coordinate system for \mathbb{R}^N .
- A point $C = (c_1, c_2, \dots, c_N) \in \mathbb{R}^N$ is mapped into the polygonal line \bar{C}
- the N -vertices with xy -coords $(i - 1, c_i)$ are now on the parallel axes.
- In \bar{C} the full lines and not only the segments between the axes are included.

Inselberg, A. (2005) Visualization of concept formation and learning. *Kybernetes: The International Journal of Systems and Cybernetics*, 34, 1/2, 151-166.

A. Holzinger 709.049 37/98 Med Informatics L09

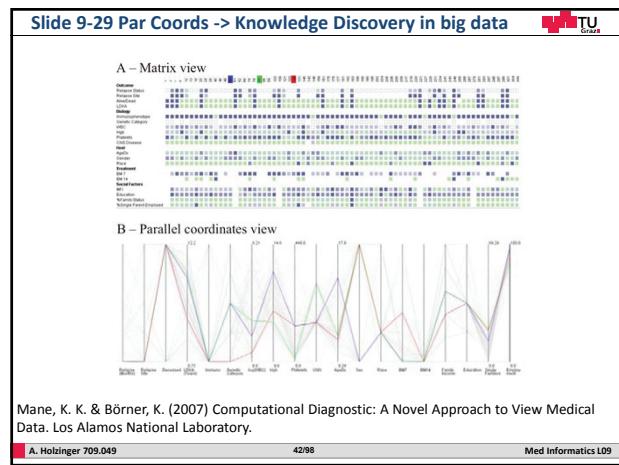
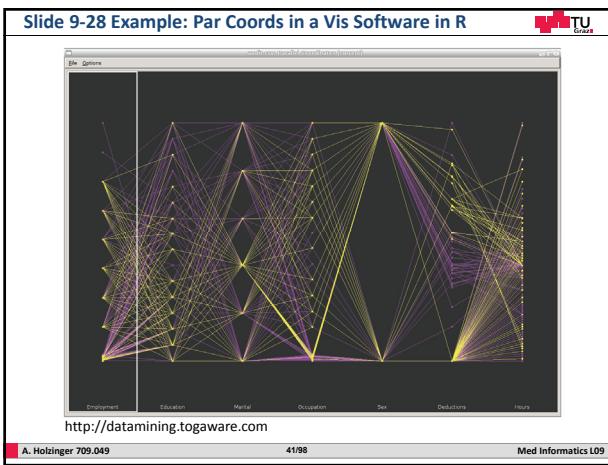
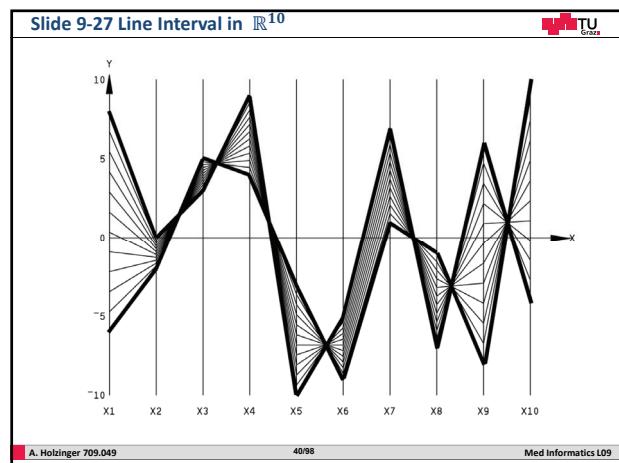


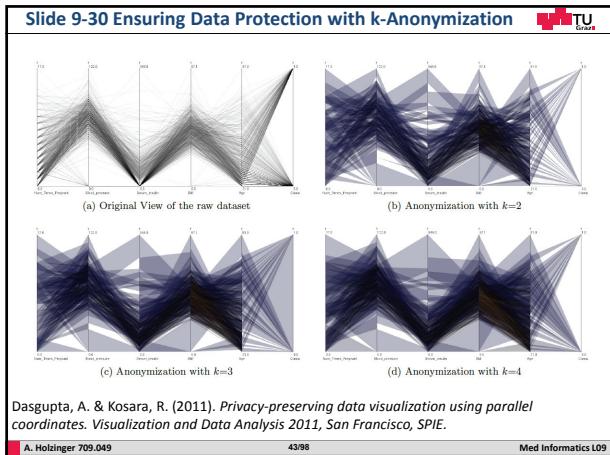
Slide 9-26 Heavier polygonal lines represent end-points



- A polygonal line \bar{P} on the $N - 1$ points represents a point
- $P = (p_1, \dots, p_{i-1}, p_i, \dots, p_N) \in \ell$
- since the pair of values \dots, p_{i-1}, p_i marked on the \bar{X}_{i-1} and \bar{X}_i axes.
- In the following slide we see several polygonal lines, intersecting at $\ell_{(i-1),i}$
- representing data points on a line $\ell \subset \mathbb{R}^{10}$.
- Note: The indexing is essential and is important for the visualization of proximity properties such as the minimum distance between a pair of lines.

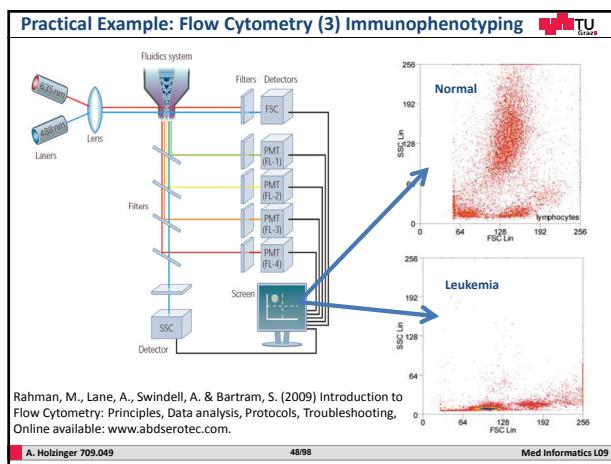
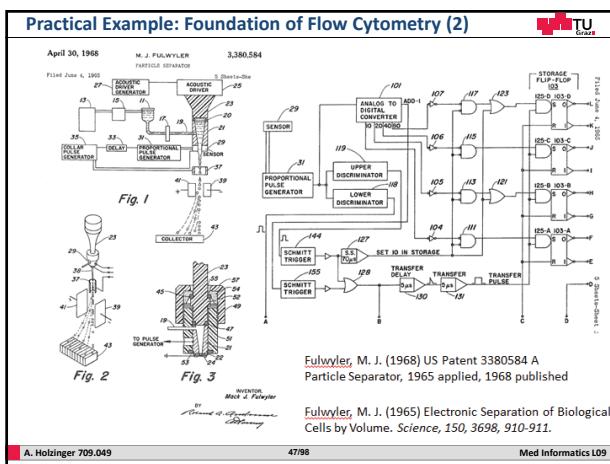
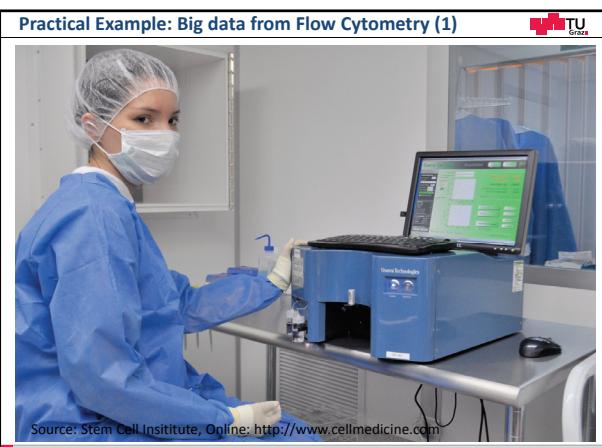
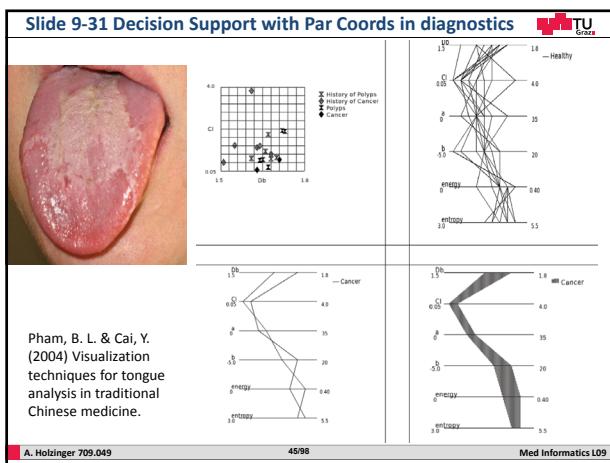
A. Holzinger 709.049 39/98 Med Informatics L09





Why are such approaches not used in enterprise hospital information systems?

A. Holzinger 709.049 44/98 Med Informatics L09



Practical Example: Flow Cytometry (4) Immunophenotyping

■ Forward scatter channel (FSC) intensity equates to the particle's size and can also be used to distinguish between cellular debris and living cells.

■ Side scatter channel (SSC) provides information about the granular content within a particle.

■ Both FSC and SSC are unique for every particle, and a combination of the two may be used to differentiate different cell types in a heterogeneous sample.

Rahman et al. (2009)

A. Holzinger 709.049 49/98 Med Informatics L09

Example: 2D Parallel Coordinates in Cytometry

Streit, M., Ecker, R. C., Österreicher, K., Steiner, G. E., Bischof, H., Bangert, C., Kopp, T. & Rogojanu, R. (2006) 3D parallel coordinate systems—A new data visualization method in the context of microscopy-based multicolor tissue cytometry. *Cytometry Part A*, 69A, 7, 601–611.

A. Holzinger 709.049 50/98 Med Informatics L09

Example: Limitations of 2D Parallel Coordinates

Streit et al. (2006)

A. Holzinger 709.049 51/98 Med Informatics L09

Parallel Coordinates in 3D

Streit et al. (2006)

A. Holzinger 709.049 52/98 Med Informatics L09

Slide 9-32 RadViz – Idea based on Hooke's Law

Demšar, J., Curk, T., & Erjavec, A. Orange: Data Mining Toolbox in Python; Journal of Machine Learning Research 14:2349–2353, 2013.

Source: <http://orange.biolab.si/>

A. Holzinger 709.049 53/98 Med Informatics L09

Slide 9-33 RadViz Principle

- Let us consider a point $y_i = (y_{i1}, y_{i2}, \dots, y_{in})$ from the n -dimensional space
- This point is now mapped into a single point u in the plane of anchors: for each anchor j the stiffness of its spring is set to y_j
- Now the Hooke's law is used to find the point u , where all the spring forces reach equilibrium (means they sum to 0). The position of $u = [u_1, u_2]$ is now derived by:

$$\sum_{j=1}^n (\vec{s}_j - \vec{u}) y_j = 0 \quad \sum_{j=1}^n \vec{s}_j y_j = \vec{u} \sum_{j=1}^n y_j$$

$$\vec{u} = \frac{\sum_{j=1}^n \vec{s}_j y_j}{\sum_{j=1}^n y_j} \quad u_1 = \frac{\sum_{j=1}^n y_j \cos(\alpha_j)}{\sum_{j=1}^n y_j} \quad u_2 = \frac{\sum_{j=1}^n y_j \sin(\alpha_j)}{\sum_{j=1}^n y_j}$$

Novakova, L. & Stepankova, O. (2009). *RadViz and Identification of Clusters in Multidimensional Data*. 13th International Conference on Information Visualisation, 104–109.

A. Holzinger 709.049 54/98 Med Informatics L09

Slide 9-34 RadViz mapping principle and algorithm

1. Normalize the data to the interval $(0, 1)$

$$\bar{x}_{ij} = \frac{x_{ij} - \min_j}{\max_j - \min_j}$$

2. Now place the dimensional anchors

3. Now calculate the point to place each record and to draw it:

$$y_i = \sum_{j=1}^n \bar{x}_{ij}$$

$$\vec{u}_i = \frac{\sum_{j=1}^n \vec{s}_j \bar{x}_{ij}}{y_i}$$

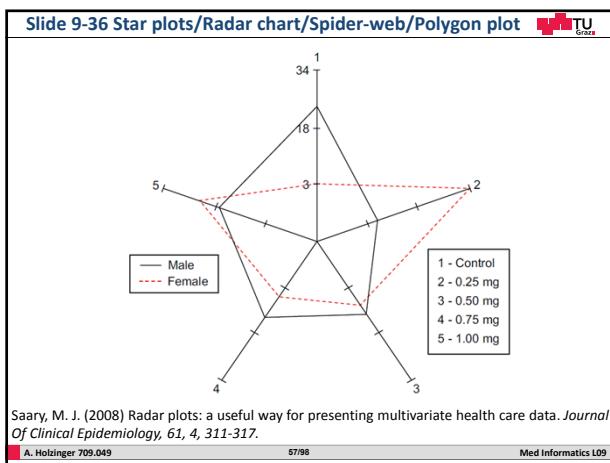
Novakova, L. & Stepankova, O. (2009). *RadViz and Identification of Clusters in Multidimensional Data*. 13th International Conference on Information Visualisation, 104-109.

A. Holzinger 709.049 55/98 Med Informatics L09

Slide 9-35 RadViz for showing the existence of clusters

Novakova, L. & Stepankova, O. (2009). *RadViz and Identification of Clusters in Multidimensional Data*. 13th International Conference on Information Visualisation, 104-109.

A. Holzinger 709.049 56/98 Med Informatics L09



- Slide 9-37 Star Plot production**
- Arrange N axes on a circle in \mathbb{R}^2
 - $3 \leq N \leq N_{max}$
Note: An amount of $N_{max} \leq 20$ is just useful, according to Lanzenberger et al. (2005)
 - Map coordinate vectors $P \in \mathbb{R}^N$ from $\mathbb{R}^N \rightarrow \mathbb{R}^2$
 - $P = \{p_1, p_2, \dots, p_N\} \in \mathbb{R}^N$ where each p_i represents a different attribute with a different physical unit
 - Each axis represents one attribute of data
 - Each data record, or data point P is visualized by a line along the data points
 - A Line is perceived better than points on the axes
- A. Holzinger 709.049 58/98 Med Informatics L09

Slide 9-38 Algorithm for drawing the axes and the lines

```

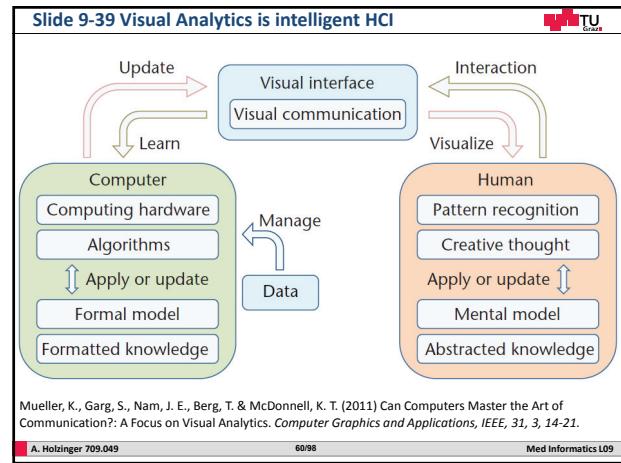
anglesector = 2 * pi / N
for each ai from axes[]
{
    anglei = i * anglesector
    xi = mid.x + r * cos(anglei)
    yi = mid.y + r * sin(anglei)
    DrawLine(midpoint.x, midpoint.y, xi, yi)

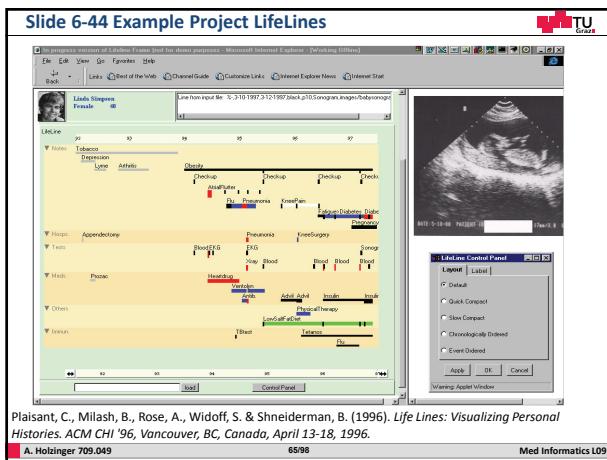
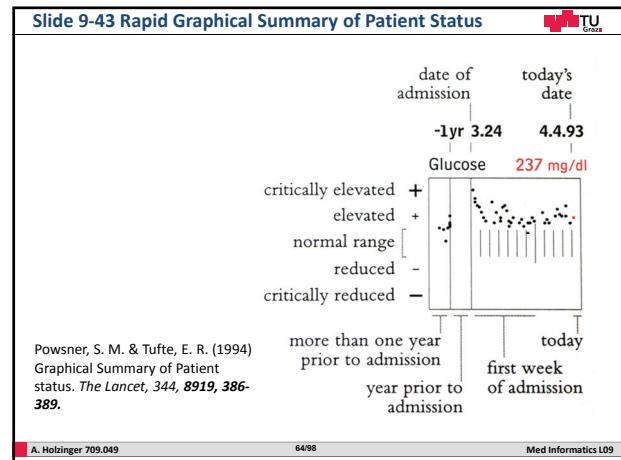
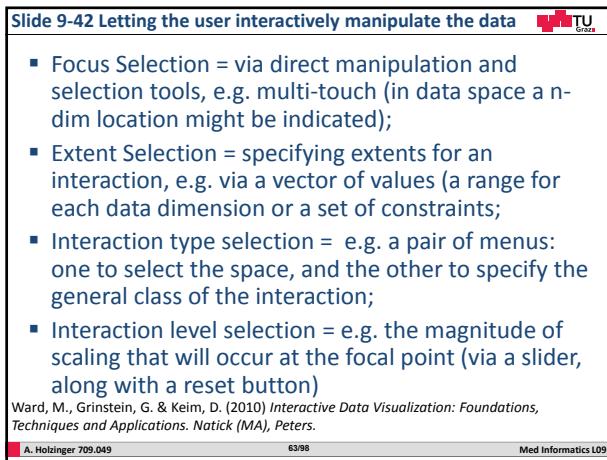
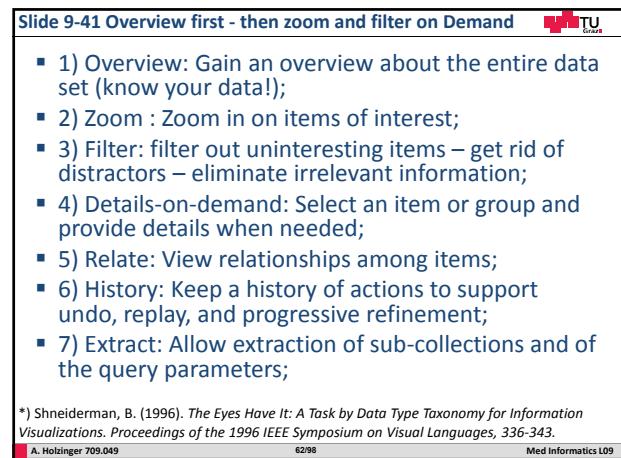
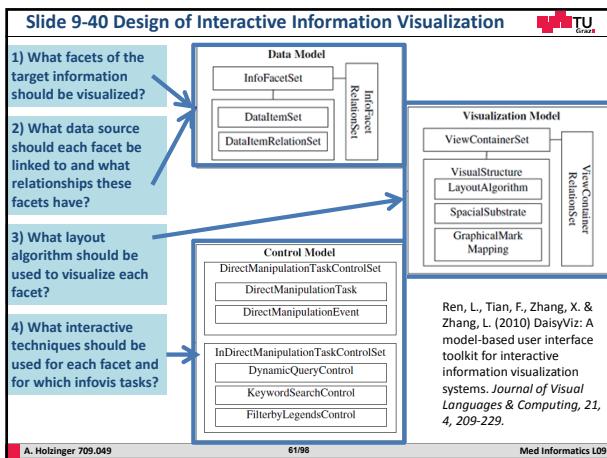
    maxi = ai.upperBound()
    scaled_vali = ai.value() * r / maxi
    x_vali = mid.x + scaled_vali * cos(anglei)
    y_vali = mid.y + scaled_vali * sin(anglei)
    DrawLine(x_vali, y_vali, x_vali-1, y_vali-1)
}

```

Dim 6 Dim 1
Dim 5 Dim 2
Dim 4 Dim 3

A. Holzinger 709.049 59/98 Med Informatics L09





What are temporal analysis tasks?

A. Holzinger 709.049 66/98 Med Informatics L09

Slide 6-45 Temporal analysis tasks

Classification = given a set of classes: the aim is to determine which class the dataset belongs to; a classification is often necessary as pre-processing;

Clustering = grouping data into clusters based on similarity; the similarity measure is the key aspect of the clustering process;

Search/Retrieval = look for a priori specified queries in large data sets (query-by-example), can be exact matched or approximate matched (similarity measures are needed that define the degree of exactness);

Pattern discovery = automatically discovering relevant patterns in the data, e.g. local structures in the data or combinations thereof;

Prediction = foresee likely future behaviour of data – to infer from the data collected in the past and present how the data will evolve in the future (e.g. autoregressive models, rule-based models etc.)

Aigner, W., Miksch, S., Schumann, H. & Tominski, C. (2011) *Visualization of Time-Oriented Data*. Human-Computer Interaction Series. London, Springer.

A. Holzinger 709.049 67/98 Med Informatics L09

Example: Subspace Clustering

Remember: The curse of dimensionality

(a) 11 Objects in One Unit Bin

(b) 6 Objects in One Unit Bin

(c) 4 Objects in One Unit Bin

A. Holzinger 709.049 Data Mining Concepts and Techniques 69 Med Informatics L09

Repeat some definitions

- Dataset - consists of a matrix of data values, rows represent individual instances and columns represent dimensions.
- Instance - refers to a vector of d measurements.
- Cluster - group of instances in a dataset that are more similar to each other than to other instances. Often, similarity is measured using a distance metric over some or all of the dimensions in the dataset.
- Subspace - is a subset of the d dimensions of a given dataset.
- Subspace Clustering – seek to find clusters in a dataset by selecting the most *relevant* dimensions for each cluster separately.
- Feature Selection - process of determining and selecting the dimensions (features) that are most relevant to the data mining task.

A. Holzinger 709.049 70/98 Med Informatics L09

Parsons et al. SIGKDD Explorations 2004

Parsons, L., Haque, E. & Liu, H. 2004. Subspace clustering for high dimensional data: a review. SIGKDD Explorations 6, (1), 90-105.

(a) Dimension a

(b) Dimension b

(c) Dimension c

(a) Dims a & b

(b) Dims b & c

(c) Dims a & c

A. Holzinger 709.049

Similar: Principal Component Analysis (PCA)

Curse of dimensionality (Bellman, 1957): As more dimensions are available, data becomes more sparse and distance measures are less meaningful.

A. Holzinger 709.049 72/98 Med Informatics L09

MMDS 2012
Workshop on Algorithms for Modern Massive Data Sets

CONTEXT – Recent related progress

Matrix completion: Given $y = P_{\Omega}[\mathbf{A}_0]$, $\Omega \subset [m] \times [n]$, recover \mathbf{A}_0 .

Impossible in general ($|\Omega| \ll mn$)
Well-posed if \mathbf{A}_0 is structured (low-rank), but still NP-hard
Tractable via convex optimization: $\min \|\mathbf{A}\|_*$ s.t. $y = P_Q(\mathbf{A})$
... if Ω is "nice" (random subset) ...
... and \mathbf{A}_0 interacts "nicely" with P_Q (\mathbf{A}_0 incoherent – not "spiky").

Hugely active area: Candes+Recht '08, Keshavan+Oh+Montanari '09, Candes+Tao '09, Gross '10, Recht '10, Negahban+Wainwright '10

Stanford University July 10-13, 2012

A. Holzinger 709.049 73/98 Med Informatics L09

012
Workshop on Algorithms for Data Sets

CONTEXT – Recent related progress

Subspace Clustering: Given $Y : [y_1, \dots, y_n] \subset S_1, \dots, S_k$, recover the subspaces.

Impossible in general (solutions highly ambiguous)
Well-posed if $\{S_i\}$ are few and structured (low-dim), but still combinatorial
Tractable via convex optimization: $\min \|\mathbf{X}\|_0 + \|\mathbf{E}\|_1$ s.t. $Y = Y\mathbf{X} + \mathbf{E}$.
... for random samples Y
... \mathbf{X} and outliers \mathbf{E} are sparse (or low-rank, column-wise sparse).

Hugely active area: Rao, Tron, Ma, Vidal '08, Elhamifar and Vidal '2010, Liu, Lin, Sun, Yan, Ma et. al. '2011, Soltanolkotabi and Candes '2011

See Rene Vidal's Talk CC BY-SA

A. Holzinger 709.049 74/98 Med Informatics L09

Slide 6-46 Future Outlook

TU Graz

- Time (e.g. entropy) and Space (e.g. topology)
- Knowledge Discovery from “unstructured” ;-)
(Forrester: >80%) data and applications of structured components as methods to index and organize data -> Content Analytics
- Open data, Big data, sometimes: small data
- Integration in “real-world” (e.g. Hospital), mobile
- How can we measure the benefits of visual analysis as compared to traditional methods?
- Can (and how can) we develop powerful visual analytics tools for the non-expert end user?

A. Holzinger 709.049 75/98 Med Informatics L09

TU Graz

HCI-KDD

Thank you!

A. Holzinger 709.049 76/98 Med Informatics L09

Sample Questions (1)

TU Graz

- What is semiotic engineering?
- Please explain the process of intelligent interactive information visualization!
- What is the difference between visualization and visual analytics?
- Explain the model of perceptual visual processing according to Ware (2004)!
- What was the historical start of systematic visual analytics? Why is this an important example?
- Please describe very shortly 6 of the most important visualization techniques!
- Transform five given data points into parallel coordinates!
- How can you ensure data protection in using parallel coordinates?
- What is the basic idea of RadViz?
- For which problem would you use a star-plot visualization?

A. Holzinger 709.049 77/98 Med Informatics L09

Sample Questions (2)

TU Graz

- What are the basic design principles of interactive intelligent visualization?
- What is the visual information seeking mantra of Shneiderman (1996)?
- Which concepts are important to let the end user interactively manipulate the data?
- What is the problem involved in looking at neonatal polysomnographic recordings?
- Why is time very important in medical informatics?
- What was the goal of LifeLines by Plaisant et al (1996)?
- Which temporal analysis tasks can you determine?
- Why is pattern discovery in medical informatics so important?
- What is the aim of foreseeing the future behaviour of medical data?

A. Holzinger 709.049 78/98 Med Informatics L09

Some useful links

- <http://vis.lbl.gov/Events/SC07/Drosophila/> (some really cool examples of high-dimensional data)
- <http://people.cs.uchicago.edu/~wiseman/chernoff.html> (Chernoff Faces in Java)
- <http://lib.stat.emu.edu> (Iris sample data set)
- <http://graphics.stanford.edu/data/voldata> (113-slice MRI data set of CT studies of cadaver heads)

A. Holzinger 709.049 79/98 Med Informatics L09

Appendix: Parallel Coordinates in a Vis Software in R

http://datamining.togaware.com

A. Holzinger 709.049 80/98 Med Informatics L09

Korrelation = +1 Korrelation = -1 Zwei Cluster Kreis Normalverteilung

Zur Visualisierung von hochdimensionalen Daten in der Statistik müssen drei wichtige Aspekte beachtet werden:

- die Anordnung der Achsen
- Die Anordnung der Achsen ist entscheidend für die Suche nach Strukturen in den Daten. In einer typischen Datenanalyse werden meist viele Anordnungen ausprobiert. Es wurden Anordnungsheuristiken entwickelt, die Einblicke in interessante Strukturen erlauben.
- die Rotation der Achsen (Daten)
- Da die 1-te Koordinate durch die Ecke auf der 1ten Achse bestimmt wird, kann eine Rotation der Achsen (= Rotation der Daten) ein anderes Bild ergeben. Die beiden linken Grafiken können als Rotation der Achsen (oder Daten) um 90 Grad aufgefasst werden. Trotz gleicher Struktur ergeben sich unterschiedliche Strukturen in den parallelen Koordinaten.
- die Skalierung der Achsen
- Die parallelen Koordinaten sind im Wesentlichen eine Aneinanderreihung von Linien zwischen Paaren von Koordinatenachsen. Daher sollten die Variablen auf einen ähnlichen Maßstab skaliert sein. Verschiedene Skalierungen können ebenfalls interessante Einsichten in die Daten geben.

A. Holzinger 709.049 81/98 Med Informatics L09

Visual Multidimensional Geometry and its Applications (1)

A. Holzinger 709.049 82/98 Med Informatics L09

Appendix: Node-link graphs to visualize biological networks

Viau, C., McGuffin, M. J., Chiricota, Y. & Jurisica, I. (2010) The FlowVizMenu and Parallel Scatterplot Matrix: Hybrid Multidimensional Visualizations for Network Exploration. *Visualization and Computer Graphics, IEEE Transactions on*, 16, 6, 1100-1108.

A. Holzinger 709.049 83/98 Med Informatics L09

Appendix: Deep View Working Environment - Swiss PDB

Main windows

Specific windows

http://www.expasy.org

A. Holzinger 709.049 84/98 Med Informatics L09

Remember: Data – Information (it is a visualization task!)

Each multivariate observation can be seen as a data point in an n -dimensional vector space

$$x_i = [x_{i1}, \dots, x_{in}]$$

- “Look at your data”
- transfer data into information
- By use of human intelligence ...
- to transfer information into knowledge ($C \rightarrow P$)
- Challenge: To reduce the dimensionality of the data ...
- ... it is an information retrieval task!

Remember: The quality can be measured by two measures:

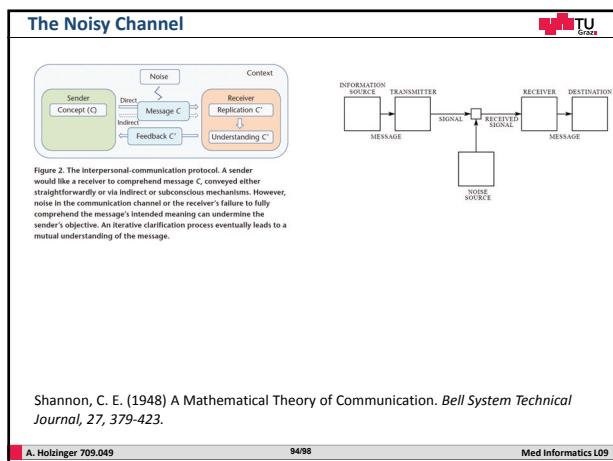
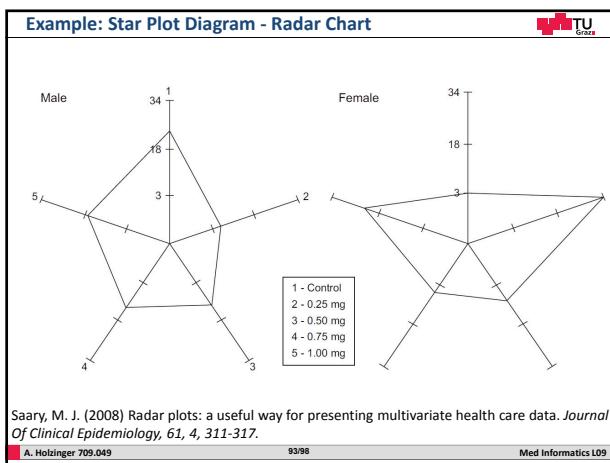
- Recall
- Precision

A. Holzinger 709.049 91/98 Med Informatics L09

Typical Problems in the Medical Clinical Domain

Holzinger, A., Hoeller, M., Bloice, M. & Urlesberger, B. (2008). Typical Problems with developing mobile applications for health care: Some lessons learned from developing user-centered mobile applications in a hospital environment. International Conference on E-Business (ICE-B 2008), Porto (PT), IEEE, 235-240.

A. Holzinger 709.049 92/98 Med Informatics L09



Slide 9-45 Example Algorithms for Selection

```

    ▪ Scatterplot-Select (xDim, yDim, xMin, xMax, yMin, yMax)
    ▪ 1 s ← 0▷ Initialize the set of records
    ▪ 2 for each record i▷ For each record,
    ▪ 3   do x ← NORMALIZE(i,xDim)▷ derive the location,
    ▪ 4   y ← NORMALIZE(i,yDim)
    ▪ 5   if xMin < x < xMax and yMin < y < yMax
    ▪ 6     do s ← s ∪ I▷ select points within rectangle
    ▪ 7 return s
    ▪ Point-in-Polygon(xs, ys, numPoints, x, y)
    ▪ 1 j ← numPoints - 1
    ▪ 2 oddNodes ← false
    ▪ 3 for i<0 to numPoints -1
    ▪ 4   do if ys[i]<y and ys[j]>y or ys[j]<y and ys[i]>y
    ▪ 5     do if xs[i]+(ys[i]/(ys[j]-ys[i]))*(xs[j]-xs[i])<x
    ▪ 6       do oddNodes ← not oddNodes
    ▪ 7       j ← I
    ▪ 8 return oddNodes
  
```

Ward, M., Grinstein, G. & Keim, D. (2010) *Interactive Data Visualization: Foundations, Techniques and Applications*. Natick (MA), Peters.

A. Holzinger 709.049 95/98 Med Informatics L09

