

האוניברסיטה העברית בירושלים THE HEBREW UNIVERSITY **OF JERUSALEM**

Artificial Intelligence in Medicine

Clustering

Nir Friedman and Tommy Kaplan

5/12/22

"In my defense, I was left unsupervised"

Lee St. John



Course outlook (so far)

The basics of machine learning

- Classifiers (= rules, predictions)
- Parameters learning
- Model selection
- What if data unlabeled?

y = ax + b $f(x) = \begin{cases} 1 & \text{if } w \cdot x > \theta \\ 0 & \text{else} \end{cases}$ The second sec

Learning Finding the optimal model

model type, parameters, simple, general, interpretable

Model Set of rules for prediction

Data

Unsupervised learning

In real life, data is often:

- Unlabeled
- High-dimensional
- Unorganized (missing/errors)
- Unfamiliar
- Unexpected





Previously, we predicted metastases presence

Sub-types? Different treatment? Prognosis?



Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications

Therese Sørlie^{a,b,c}, Charles M. Perou^{a,d}, Robert Tibshirani^e, Turid Aas^f, Stephanie Geisler^g, Hilde Johnsen^b, Trevor Hastie^e, Michael B. Eisen^h, Matt van de Rijnⁱ, Stefanie S. Jeffrey^j, Thor Thorsen^k, Hanne Quist¹, John C. Matese^c, Patrick O. Brown^m, David Botstein^c, Per Eystein Lønning⁹, and Anne-Lise Børresen-Dale^{b,n}

Departments of ^bGenetics and ^lSurgery, The Norwegian Radium Hospital, Montebello, N-0310 Oslo, Norway; ^dDepartment of Genetics and Lineberg Comprehensive Cancer Center, University of North Carolina, Chapel Hill, NC 27599; Departments of eHealth Research and Policy and Statistics, ^cGenetics, ⁱPathology, ^jSurgery, and ^mBiochemistry and Howard Hughes Medical Institute, Stanford University School of Medicine, Stanford, CA 94305; Departments of 9Medicine (Section of Oncology), fSurgery, and Biochemical Endocrinology, Haukeland University Hospital, N-5021 Bergen, Norway; and ^hLife Sciences Division, Lawrence Orlando Berkeley National Laboratories, and Department of Molecular and Cellular Biology, University of California, Berkeley, CA 94720







Breast-like Subtype C Subtype B Subtype A

Unsupervised learning

Clustering allows:

- Grouping
- Qualitative find archetypes
- Quantitative 1 how many flavors
- Quantitative 2 common vs rare/exceptional cases
- Axes by which sample vary





Clustering in medicine / medical research

- Group patients by sub-types
- Single-cell data better understanding of disease-associated cells (qualitative, quantitative, cellular/disease dynamics)
- Aging / neurodegenerative examples
- Metagenomics assembly

ARTICLE

The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups

Christina Curtis^{1,2}†*, Sohrab P. Shah^{3,4}*, Suet-Feung Chin^{1,2}*, Gulisa Turashvili^{3,4}*, Oscar M. Rueda^{1,2}, Mark J. Dunning², Doug Speed^{2,5}†, Andy G. Lynch^{1,2}, Shamith Samarajiwa^{1,2}, Yinyin Yuan^{1,2}, Stefan Gräf^{1,2}, Gavin Ha³, Gholamreza Haffari³, Ali Bashashati³, Roslin Russell², Steven McKinney^{3,4}, METABRIC Group[‡], Anita Langerad⁶, Andrew Green⁷, Elena Provenzano⁸, Gordon Wishart⁸, Sarah Pinder⁹, Peter Watson^{3,4,10}, Florian Markowetz^{1,2}, Leigh Murphy¹⁰, Ian Ellis⁷, Arnie Purushotham^{9,11}, Anne-Lise Børresen-Dale^{6,12}, James D. Brenton^{2,13}, Simon Tavaré^{1,2,5,14}, Carlos Caldas^{1,2,8,13} & Samuel Aparicio^{3,4}



Figure 5: The integrative subgroups have distinct clinical outcomes.





FEATURE ARTICLE

SCIENCE FORUM

The Human Cell Atlas

AVIV REGEV*, SARAH A TEICHMANN*, ERIC S LANDER*



Comprehensive Classification of Retinal Bipolar Neurons by Single-Cell Transcriptomics





Neuronal subtypes and diversity revealed by single-nucleus RNA sequencing of the human brain

Blue B. Lake,^{1*} Rizi Ai,^{2*} Gwendolyn E. Kaeser,^{3,4*} Neeraj S. Salathia,^{5*} Yun C. Yung,³ Rui Liu,¹ Andre Wildberg,² Derek Gao,¹ Ho-Lim Fung,¹ Song Chen,¹ Raakhee Vijayaraghavan,⁵ Julian Wong,³ Allison Chen,³ Xiaoyan Sheng,³ Fiona Kaper,⁵ Richard Shen,⁵ Mostafa Ronaghi,⁵ Jian-Bing Fan,⁵† Wei Wang,²† Jerold Chun,³† Kun Zhang¹†









People are good in low-dimensional clustering



But how would we formally define a good clustering?



People are good in low-dimensional clustering



But how would we formally define a good clustering?



Formal definition - top down

- Density-based clustering
- Minimize distances from "centroid"

• Each point x "belongs" to class C_i whose center is at μ_i

$$\underset{C,\mu}{\operatorname{arg\,min}} \sum_{i=1}^{k} \sum_{x \in C_i} |x - \mu_i|^2$$



How to find C,μ ?

Developed by Lloyd, 1957

Init:

Choose at random K data points, as centroids

Loop:

- 1. Re-assign each point to nearest centroid
- 2. Move each centroid to mean (or "center of mass") of assigned data points

Stop:

Upon convergence (no assignment changes)

• Stochasticity



Init:

Choose at random K data points, as centroids Loop:

- 1. Re-assign each point to nearest centroid
- 2. Move each centroid to mean (or "center of mass") of assigned data points

Stop:

Upon convergence (no assignment changes)

Would you cluster differently?

• Stochasticity - might not converge to optimal solution





- Stochasticity
- Random restarts
- Initialization
 - K-means++
 - Cluster subset of data
- Running time
 - \circ Each iteration: $\mathcal{O}(KN)$
 - But how many iterations are typically needed?
- Possible data transformations
 - Feature selection
 - Data transformation

Init:

Choose at random K data points, as centroids Loop:

- 1. Re-assign each point to nearest centroid
- Move each centroid to mean (or "center of mass") of assigned data points

Stop:

Upon convergence (no assignment changes)

Choosing K using the Elbow method



People are good in low-dimensional clustering



But how would we formally define a good clustering?



Top-down approach

Bottom-up approach

Hierarchical Agglomerative Clustering

• Iteratively, merge similar points / sub-groups





dendrogram

Hierarchical Agglomerative Clustering

• Iteratively, merge similar points / sub-groups





dendrogram

Proc. Natl. Acad. Sci. USA Vol. 95, pp. 14863–14868, December 1998 Genetics

Cluster analysis and display of genome-wide expression patterns

MICHAEL B. EISEN*, PAUL T. SPELLMAN*, PATRICK O. BROWN[†], AND DAVID BOTSTEIN*[‡]

*Department of Genetics and [†]Department of Biochemistry and Howard Hughes Medical Institute, Stanford University School of Medicine, 300 Pasteur Avenue, Stanford, CA 94305



Hierarchical Agglomerative Clustering

How to define distance between clusters?



Running time: $\mathcal{O}(N^3)$

What have we learned?

- Unsupervised data
- How to approach?
- Cluster to find typical samples (archetypes)
- Top-down (divisive) clustering
 - K-means (+ useful tricks)
- Bottom-up (agglomerative) clustering
 - Hierarchical clustering
 - Interpretation
 - Variations

Syllabus

- 1. Introduction
- 2. Classification
- 3. Learning 1
- 4. Al in ophthalmology (Prof. Itay Chowers)
- 5. Learning 2
- 6. Regression
- 7. Clustering
- 8. Visualization (and dimensionality reduction)
- 9. Deep learning in image analysis (Prof. Leo Joskowicz)
- 10. Missing data, statistical dependencies
- 11. Natural language in medicine (Dr. Gabi Stanovsky)
- 12. Decisions (utility)
- 13. Longitudinal Data / Project