

Artificial Intelligence in Medicine

Clustering

Nir Friedman and Tommy Kaplan 19/2/24

"In my defense, I was left unsupervised"

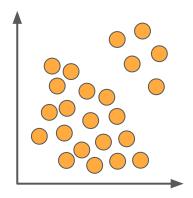
Lee St. John



Course outlook (so far)

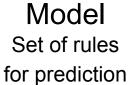
The basics of machine learning

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

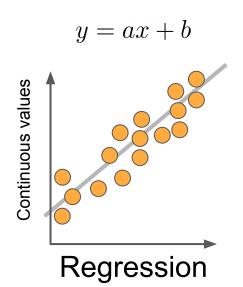


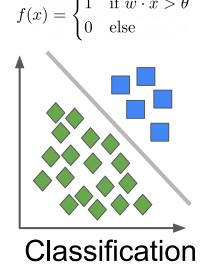
Learning Finding the optimal model

model type, parameters, simple, general, interpretable





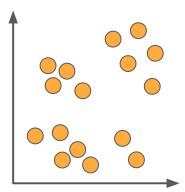




Unsupervised learning

In real life, data is often:

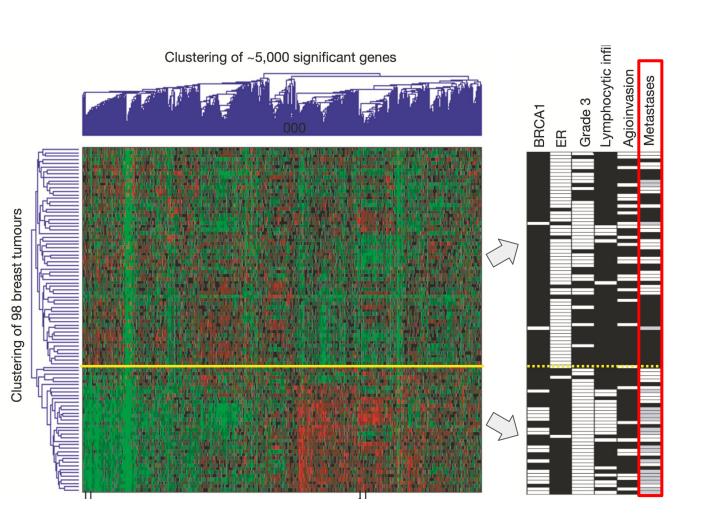
- Unlabeled
- High-dimensional
- Unorganized (missing/errors)
- Unfamiliar
- Unexpected

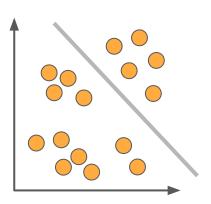




Previously, we predicted metastases presence

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 Quist^l, John C. Matese^c, Patrick O. Brown^m, David Botstein^c, Per Eystein Lønning^g, and Anne-Lise Børresen-Dale^{b,n}

>6 >4 >2 1:1 >2 >4 >6 >8

Luminal

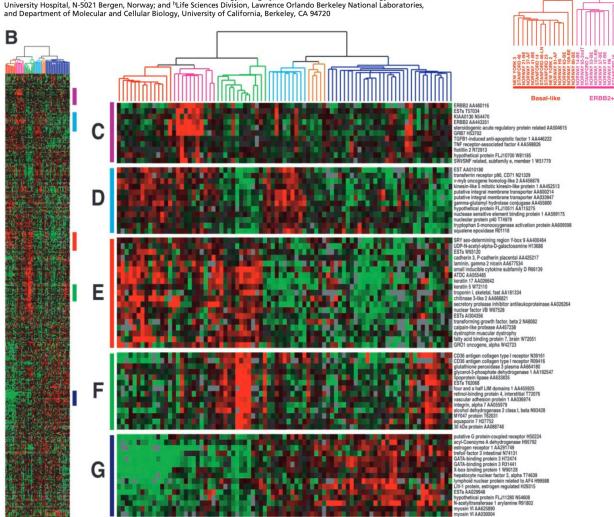
Subtype A

A

Breast-like

Subtype C Subtype B

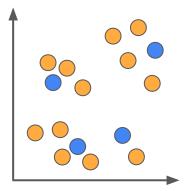
Departments of ^bGenetics and ^lSurgery, The Norwegian Radium Hospital, Montebello, N-0310 Oslo, Norway; ^dDepartment of Genetics and Lineberg Comprehensive Cancer Center, University of North Carolina, Chapel Hill, NC 27599; Departments of ^eHealth Research and Policy and Statistics, ^{(c}Genetics, ^lPathology, ^lSurgery, and ^mBiochemistry and Howard Hughes Medical Institute, Stanford University School of Medicine, Stanford, CA 94305; Departments of ^eMedicine (Section of Oncology), ^{(s}Surgery, and ^eBiochemical Endocrinology, Haukeland University Hospital, N-5021 Bergen, Norway; and ^hLife Sciences Division, Lawrence Orlando Berkeley National Laboratories, and Department of Medicine (Section of Oncology), ^(s) (Section



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



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 Langerød⁶, Andrew Green⁷, Elena Provenzano⁸, Gordon Wishart⁸, Sarah Pinder⁹, Peter Watson^{3,4,10}, Florian Markowetz^{1,2}, Leigh Murphyl⁰, 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}

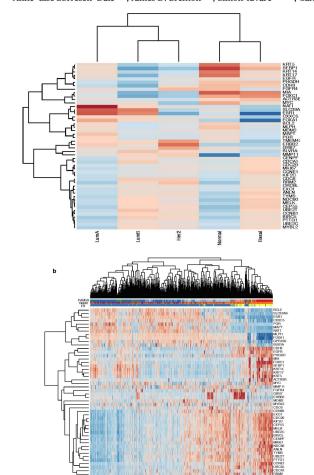
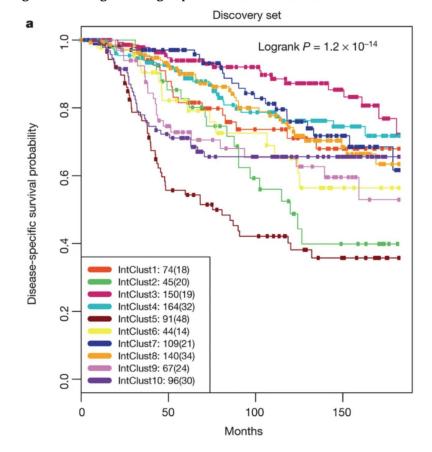


Figure 5: The integrative subgroups have distinct clinical outcomes.



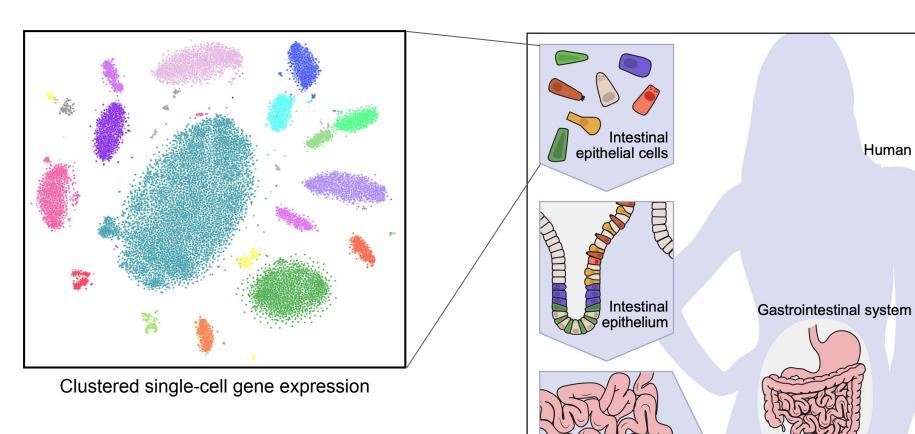
Human



SCIENCE FORUM

The Human Cell Atlas

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

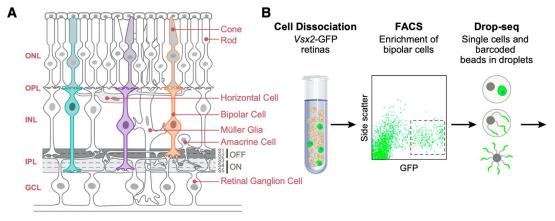


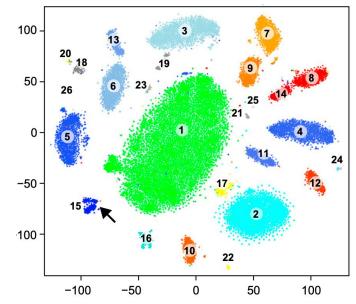
Small intestine

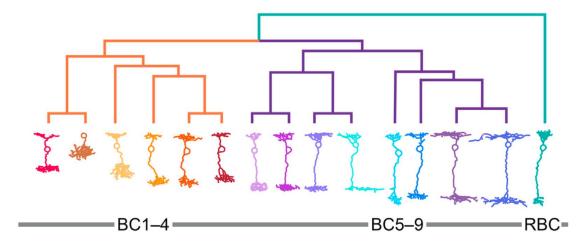
Comprehensive Classification of Retinal Bipolar Neurons by Single-Cell Transcriptomics

Karthik Shekhar,^{1,9} Sylvain W. Lapan,^{2,7,9} Irene E. Whitney,^{4,9} Nicholas M. Tran,⁴ Evan Z. Macosko,^{2,5,6} Monika Kowalczyk,¹ Xian Adiconis,^{1,5} Joshua Z. Levin,^{1,5} James Nemesh,^{2,5,6} Melissa Goldman,^{2,5} Steven A. McCarroll,^{2,5,6} Constance L. Cepko,^{2,3,7,*} Aviv Regev,^{1,7,8,*} and Joshua R. Sanes^{4,10,*}

⁴Center for Brain Science and Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA 02130, USA







¹Broad Institute of Harvard and MIT, Cambridge, MA 02142, USA

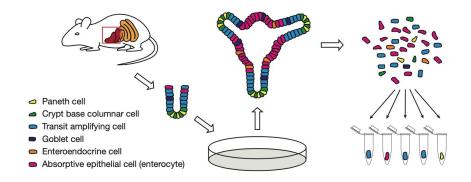
²Department of Genetics, Harvard Medical School, Boston, MA 02115, USA

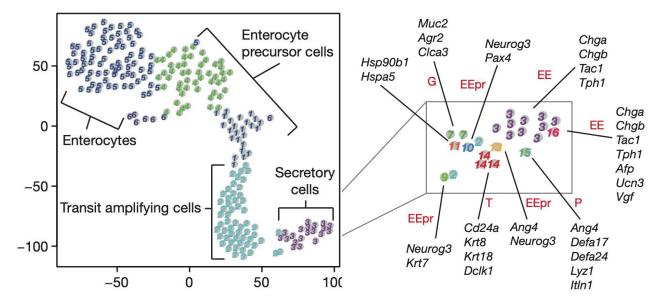
³Department of Ophthalmology, Harvard Medical School, Boston, MA 02115, USA



Single-cell messenger RNA sequencing reveals rare intestinal cell types

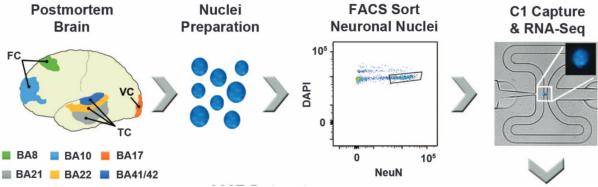
Dominic Grün 1,2* , Anna Lyubimova 1,2* , Lennart Kester 1,2 , Kay Wiebrands 1,2 , Onur Basak 1,2 , Nobuo Sasaki 1,2 , Hans Clevers 1,2 & Alexander van Oudenaarden 1,2

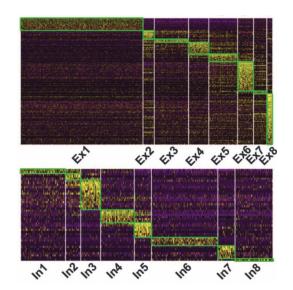


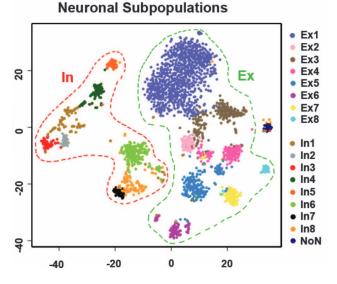


Neuronal subtypes and diversity revealed by single-nucleus RNA sequencing of the human brain

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 Zhang ¹†



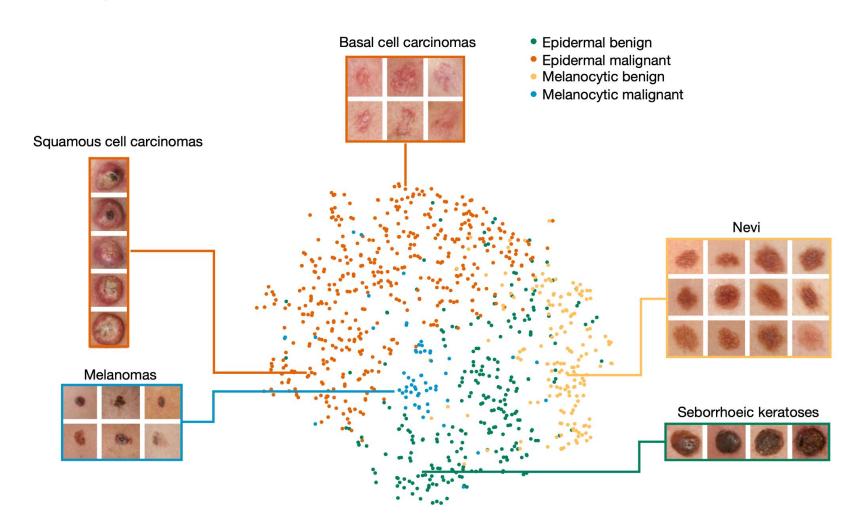




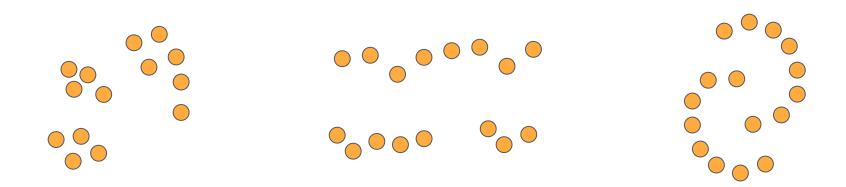
LETTER

Dermatologist-level classification of skin cancer with deep neural networks

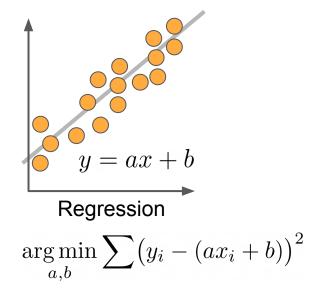
Andre Esteva¹*, Brett Kuprel¹*, Roberto A. Novoa^{2,3}, Justin Ko², Susan M. Swetter^{2,4}, Helen M. Blau⁵ & Sebastian Thrun⁶



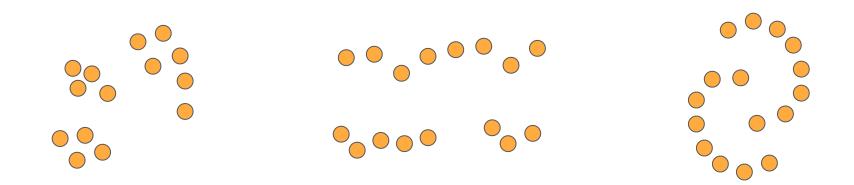
People are good in low-dimensional clustering



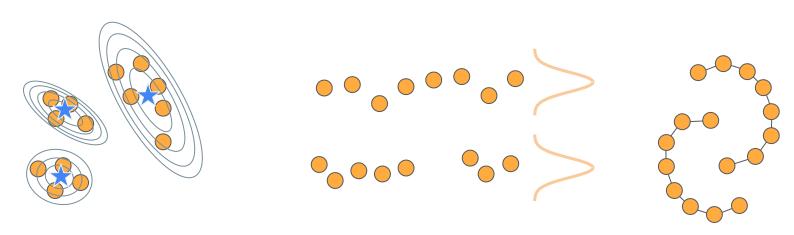
But how would we formally define a good clustering?



People are good in low-dimensional clustering



But how would we formally define a good clustering?



Top-down approach

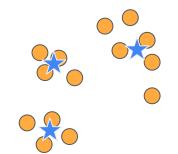
Bottom-up approach

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:

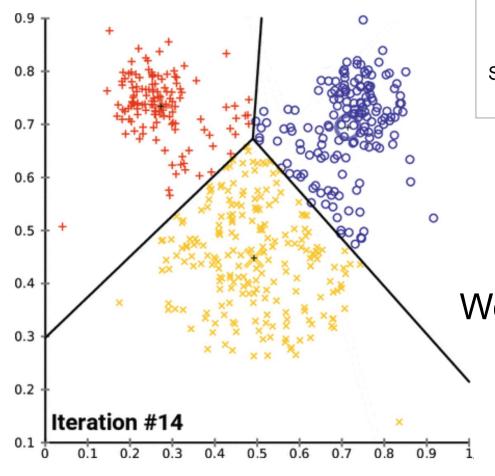
Choose at random K data points, as centroids

Loop:

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

Stop:

Stochasticity



Init:

Choose at random K data points, as centroids

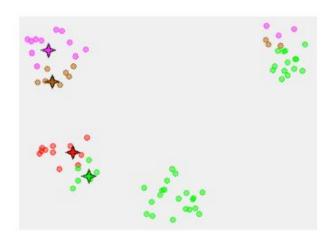
Loop:

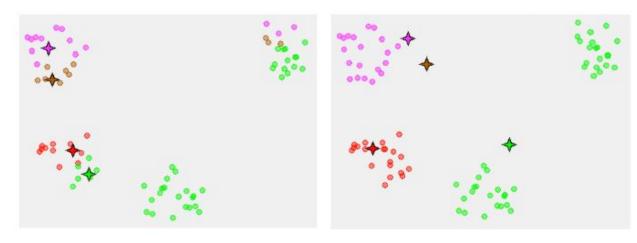
- 1. Re-assign each point to nearest centroid
- Move each centroid to mean (or "center of mass") of assigned data points

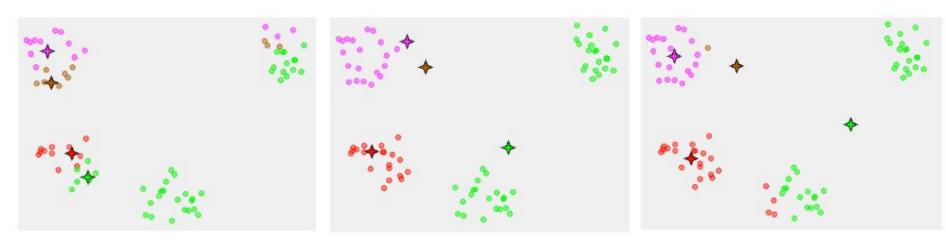
Stop:

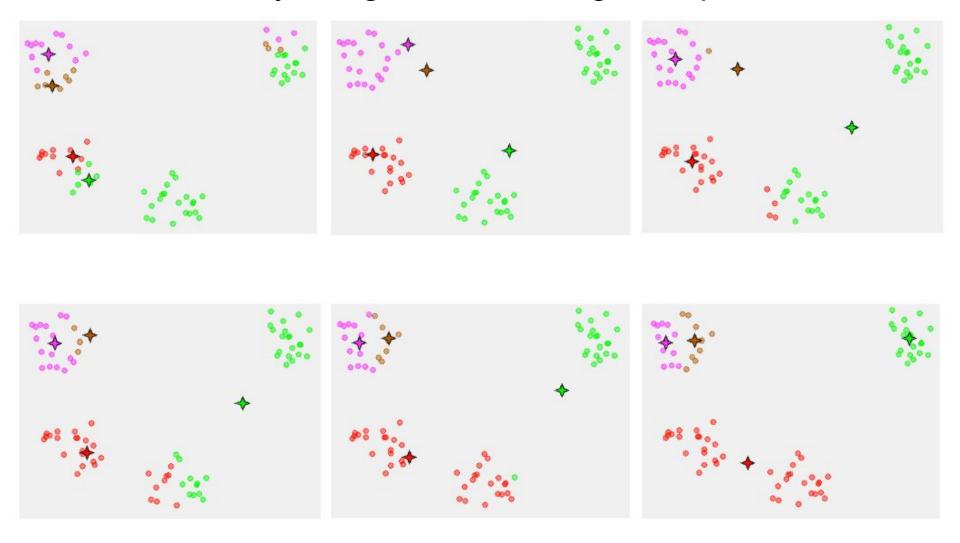
Upon convergence (no assignment changes)

Would you cluster differently?









- Stochasticity
 - Convergence to local optimum?

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)

o YES!

Step 1 assigns each point to closest centroid (better or equal)

Step 2 minimizes sum-of-squared-distances for all clusters

- Stochasticity
- Random restarts

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:

- Stochasticity
- Random restarts
- Initialization
 - K-means++
 - Cluster subset of data

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:

- 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?

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:

- 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:

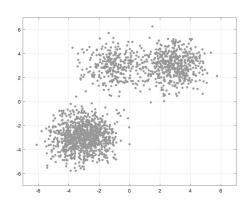
Choose at random K data points, as centroids

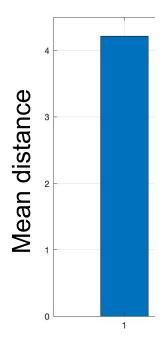
Loop:

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

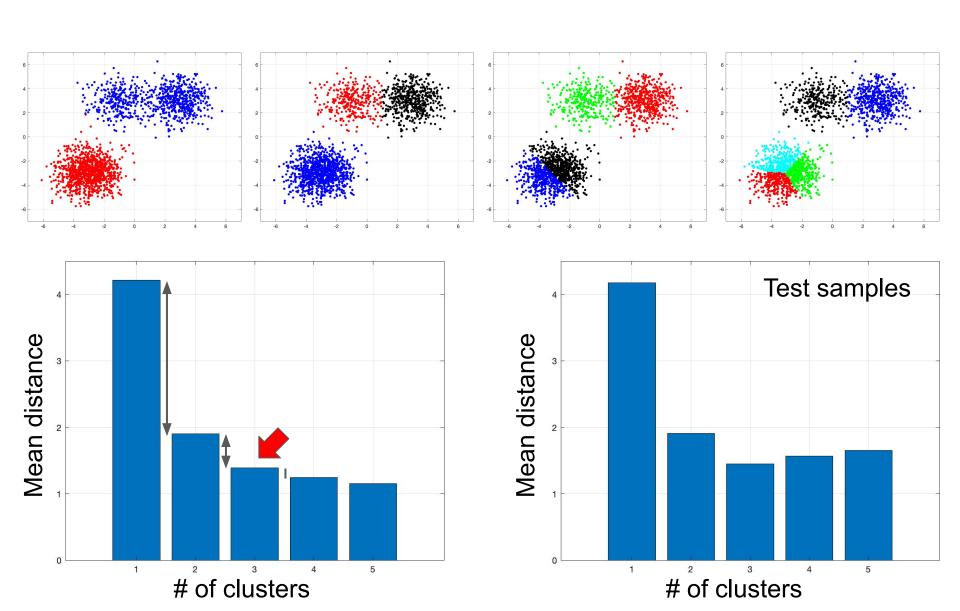
Stop:

Choosing K using the Elbow method

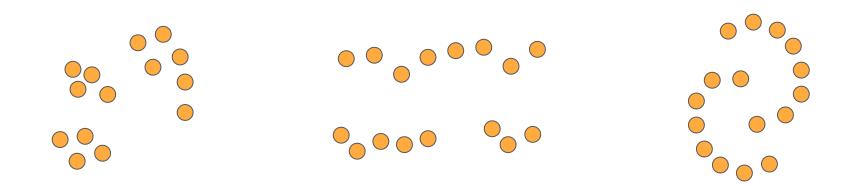




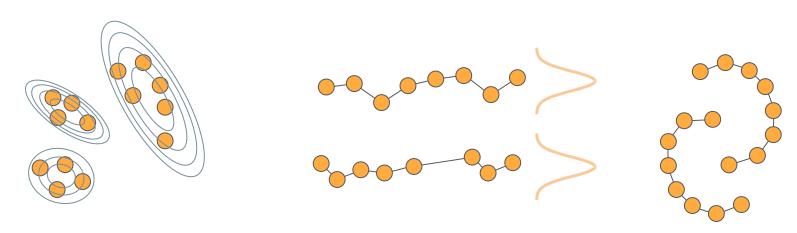
Choosing K using the Elbow method



People are good in low-dimensional clustering

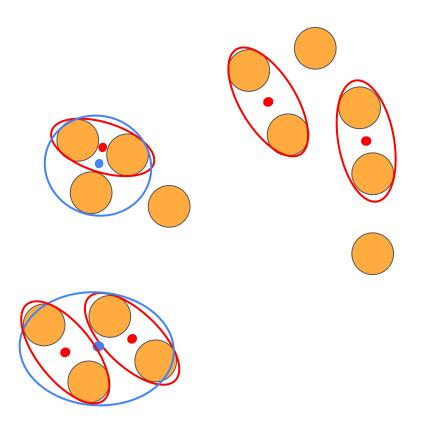


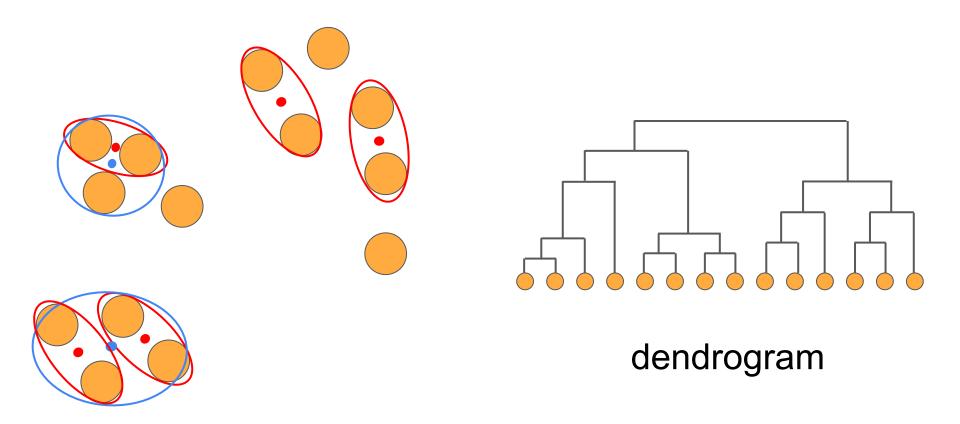
But how would we formally define a good clustering?

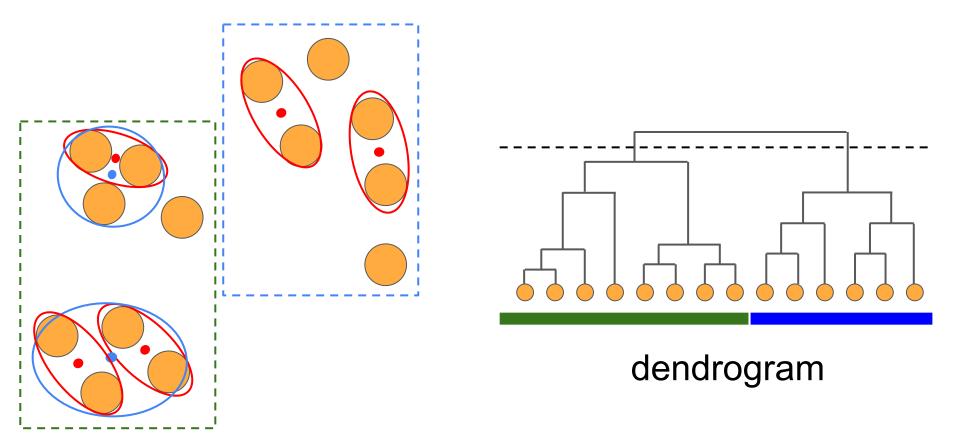


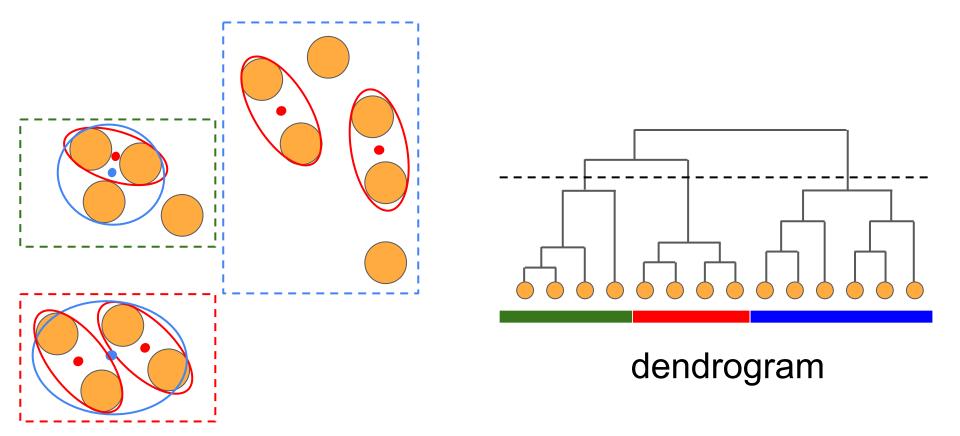
Top-down approach

Bottom-up approach





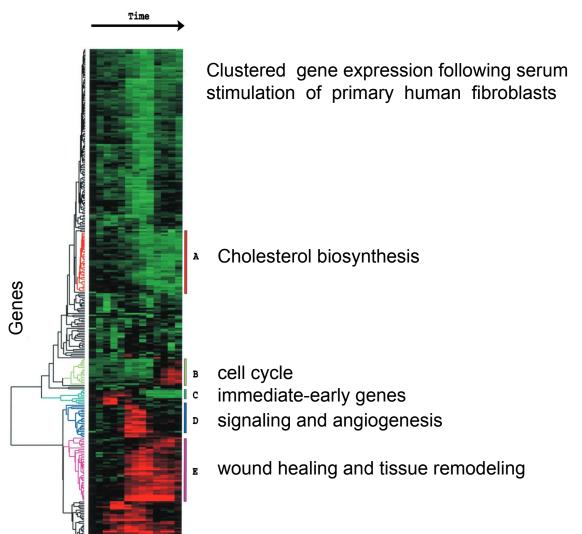




Cluster analysis and display of genome-wide expression patterns

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



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. Jeffreyⁱ, Thor Thorsen^k, Hanne Quistⁱ, John C. Matese^c, Patrick O. Brown^m, David Botstein^c, Per Eystein Lønning^g, and Anne-Lise Børresen-Dale^{b,n}

>6 >4 >2 1:1 >2 >4 >6 >8

Luminal

Subtype A

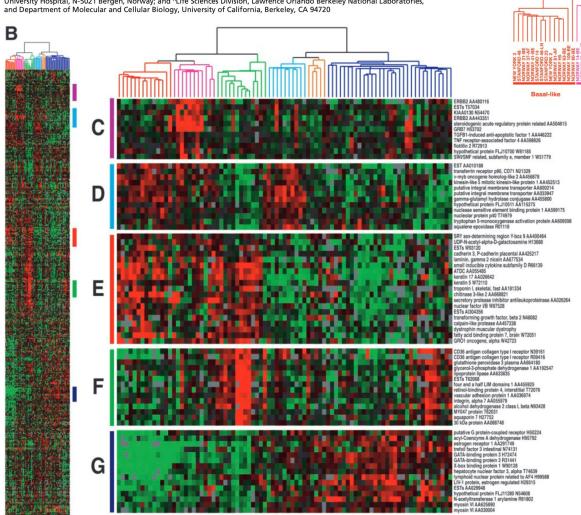
A

ERBB2+

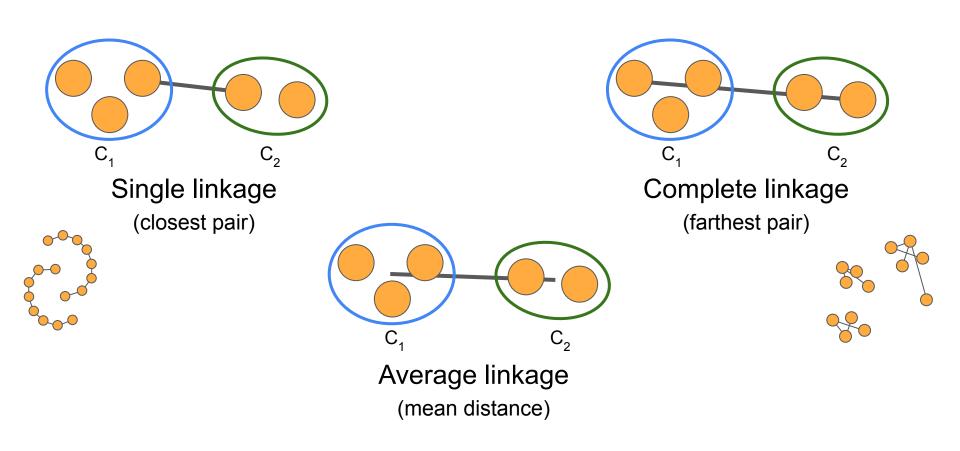
Breast-like

Subtype C Subtype B

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, ^{(c}Genetics, ^lPathology, ^lSurgery, and ^mBiochemistry and Howard Hughes Medical Institute, Stanford University School of Medicine, Stanford, CA 94305; Departments of ^eMedicine (Section of Oncology), ^{(s}Surgery, and ^(e)Biochemical Endocrinology, Haukeland University Hospital, N-5021 Bergen, Norway; and ^(e)Life Sciences Division, Lawrence Orlando Berkeley National Laboratories, and Department of Medicine (Section of Oncology), ^(e)Life Sciences Division, Lawrence Orlando Berkeley National Laboratories, ^(e)



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

Often there is more than one "good" solution

"Clustering: science or art?"

Syllabus

1	1/1	Al in ophthalmology (Prof. Itay Chowers)
2	8/1	Classification
3	15/1	Learning 1
4	22/1	Learning 2
5	7/2	Regression (Wed.)
6	12/2	Deep learning in image analysis (Prof. Leo Joskowicz)
7	19/2	Clustering
8	26/2	Dimensionality reduction and visualization
9	28/2	Deep learning, Missing data (Wed.)
10	4/3	Natural language in medicine (Dr. Gabi Stanovsky)
11	11/3	?