
3000788 Intro to Comp Molec Biol

Lecture 27: Principles of machine learning

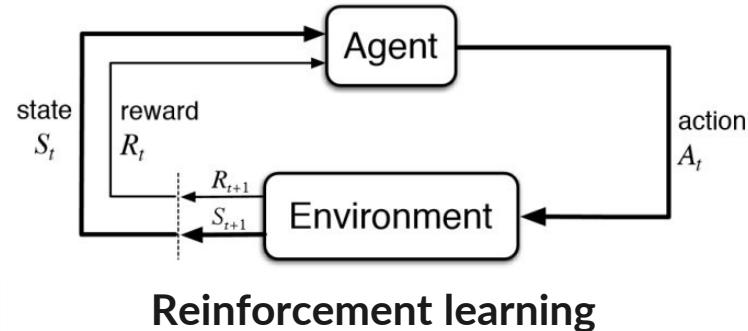
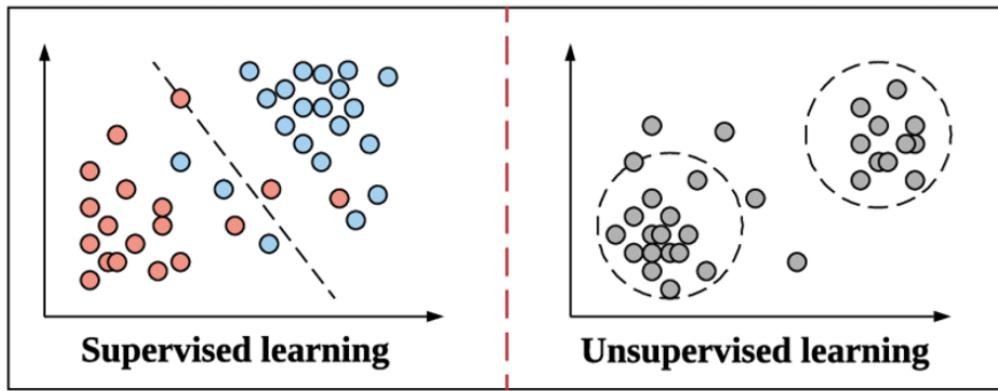
November 20, 2023



Sira Sriswasdi, PhD

- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
- Center for Artificial Intelligence in Medicine (CU-AIM)

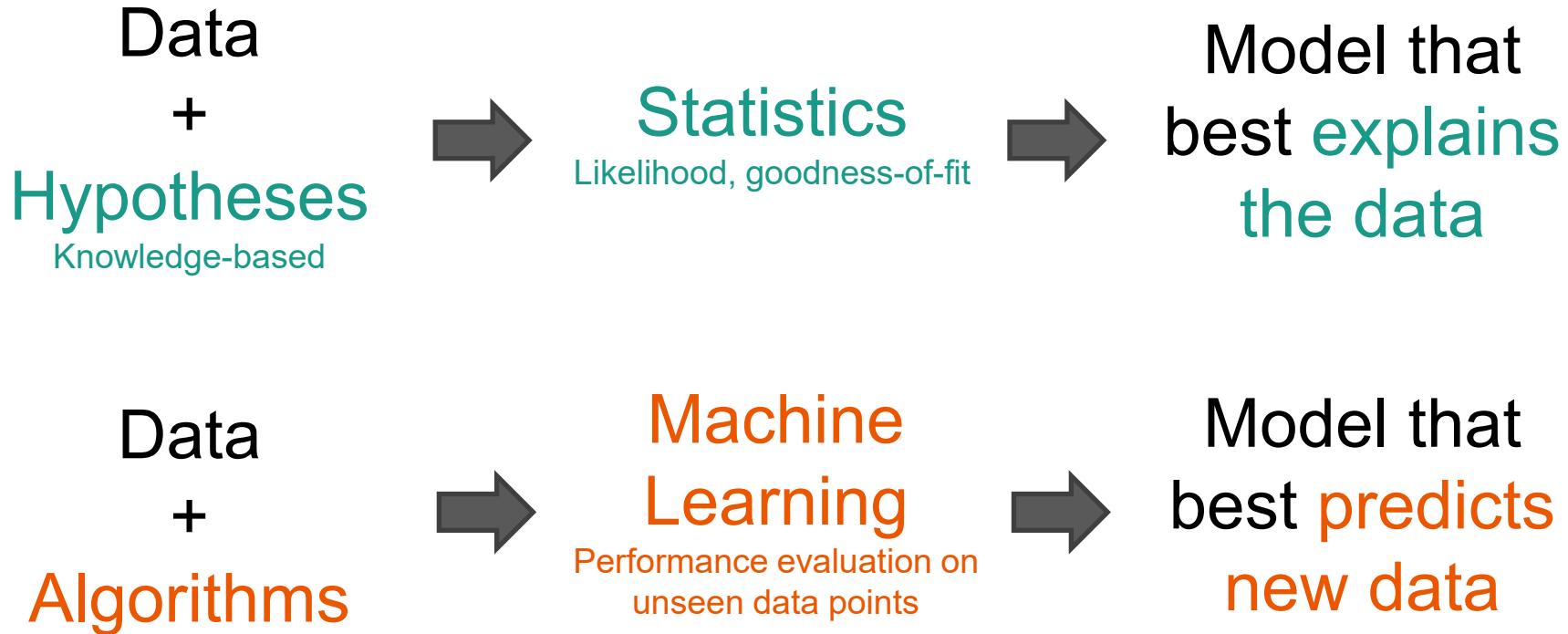
Machine learning paradigms



Qian, B. et al. "Orchestrating the Development Lifecycle of Machine Learning-Based IoT Applications: A Taxonomy and Survey"

- **Supervised:** Model learns from a dataset (x, y) to predict y from x
- **Unsupervised:** Pattern recognition with no target output (only x)
- **Reinforcement:** Model learns by interacting and receiving feedbacks from the environment (dynamic data)

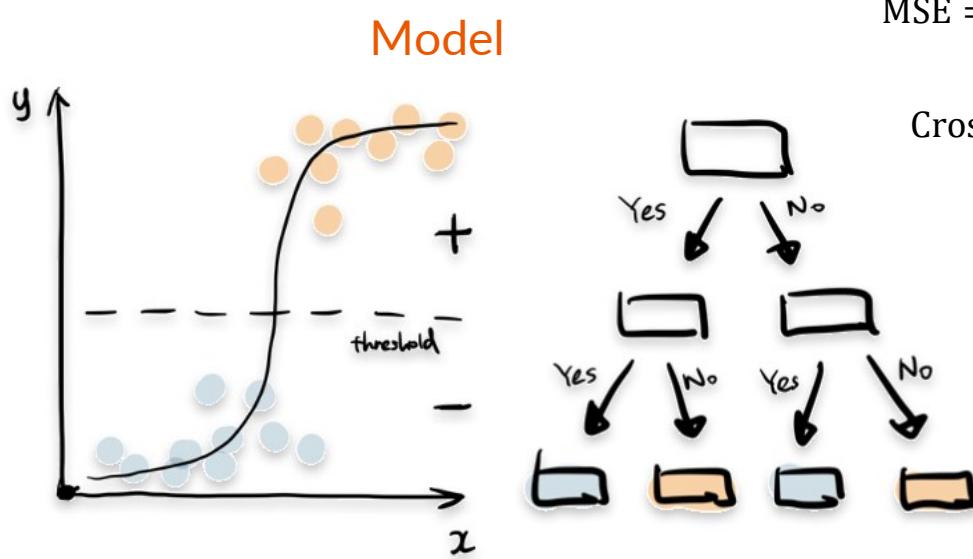
Machine learning versus human's way of thinking





Supervised learning

The cores of supervised learning

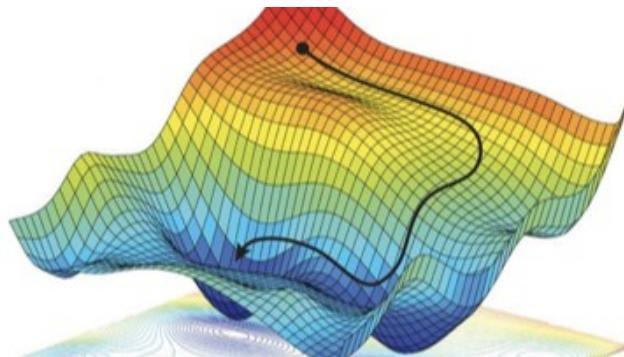


Objective / Loss Function

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{MAPE} = \frac{1}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{y_i} \times 100$$

$$\text{Crossentropy} = -\frac{1}{n} \sum_{i=1}^n y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$

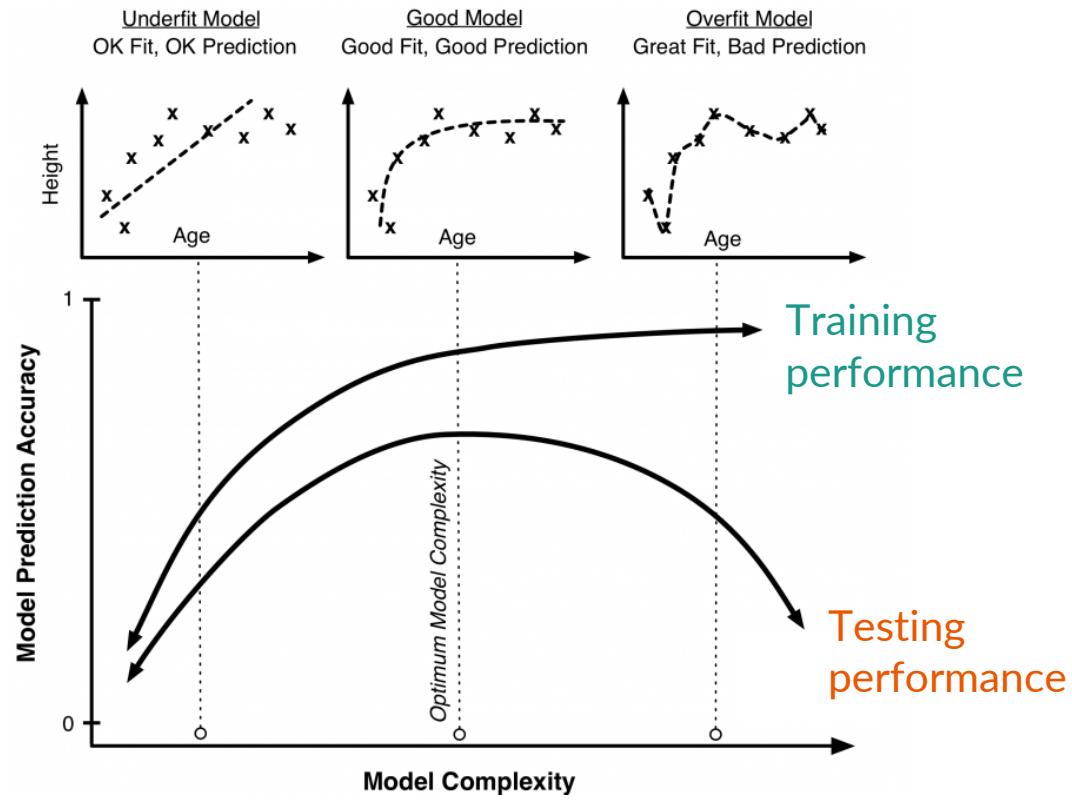
Optimization Algorithm



Supervised learning is all about control



https://en.wikipedia.org/wiki/Bull_riding



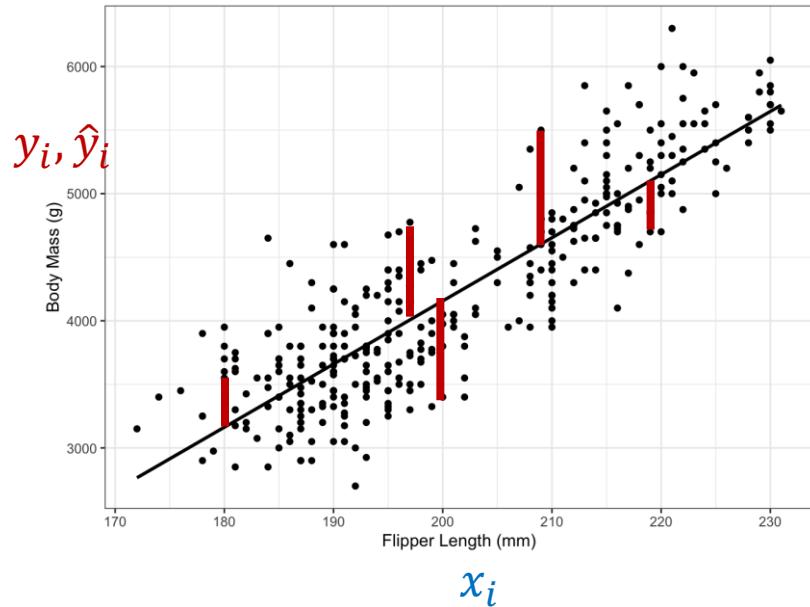
Statistical control of overfitting

- Better model achieves **higher likelihood**
- Complex model has **more parameters**
- **Information Criterion**
 - Akaike (AIC) = $2k - 2 \cdot \ln(\hat{L})$, where \hat{L} is the likelihood
 - Bayesian (BIC) = $\ln(n)k - 2 \cdot \ln(\hat{L})$, where n is the sample size
- **Nested model testing**
 - Simple model has n parameters, fit the data with likelihood \hat{L}_1
 - Complex model has $m > n$ parameters, fit the data with likelihood $\hat{L}_2 > \hat{L}_1$
 - Is the improvement $\frac{\hat{L}_2}{\hat{L}_1}$ worth the increase in $m - n$ parameters?



Linear and logistic regression

Linear regression (Ordinary Least Square)

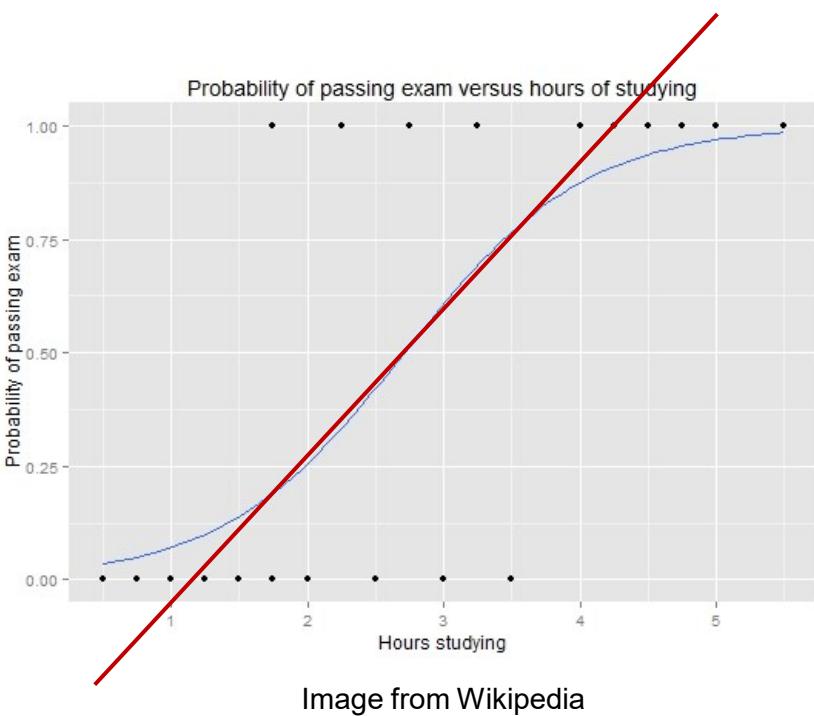


- Model: $\hat{y}_i = b_0 + b_1 x_i$
- Minimize MSE: $\frac{1}{n} \sum_i (y_i - [b_0 + b_1 x_i])^2$
- $\frac{\delta MSE}{\delta b_0} = -2 \sum_i y_i - 2b_1 \sum_i x_i - 2n b_0$
- $\frac{\delta MSE}{\delta b_1} = -2 \sum_i x_i y_i - 2b_1 \sum_i x_i^2 - 2b_0 \sum_i x_i$
- $b_0 = \frac{\sum xy - \sum x \sum y}{(\sum x)^2 - n \sum x^2}$
- $b_1 = \frac{\sum y \sum x - n \sum xy}{(\sum x)^2 - n \sum x^2}$

Ordinary Least Square interpretation

- Observed value = True value + Normally-distributed noise
- **Assumption:** Noises are identical and independent across samples
- Model: $(y_i - \hat{y}_i) \sim N(0, \sigma^2)$
- Density: $P(y_i - \hat{y}_i = \varepsilon_i | \sigma^2) \propto e^{\frac{-\varepsilon_i^2}{2\sigma^2}}$
- Likelihood: $\prod_i P(y_i - \hat{y}_i = \varepsilon_i | \sigma^2) \propto e^{\frac{-\sum_i \varepsilon_i^2}{2\sigma^2}}$
- MSE: $\frac{1}{n} \sum_i (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_i \varepsilon_i^2$
- Minimizing MSE is the same as maximizing likelihood

Logistic regression



- Classification output = 0 or 1
- Linear regression outputs $-\infty$ to ∞
- Probability of success p
- Log odd: $\ln\left(\frac{p}{1-p}\right)$
 - $\ln\left(\frac{p}{1-p}\right) \rightarrow -\infty$ as $p \rightarrow 0$
 - $\ln\left(\frac{p}{1-p}\right) \rightarrow \infty$ as $p \rightarrow 1$
- Transform linear regression output with log odd!

Logistic regression

- Model: $\ln\left(\frac{\hat{y}_i}{1-\hat{y}_i}\right) = f(x_i) = b_0 + b_1x_{i,1} + \cdots + b_nx_{i,n}$
- $\hat{y}_i = \frac{e^{b_0+b_1x_{i,1}+\cdots+b_nx_{i,n}}}{1+e^{b_0+b_1x_{i,1}+\cdots+b_nx_{i,n}}}$
 - When $f(x_i) \rightarrow \infty$, $\hat{y}_i \rightarrow 1$
 - When $f(x_i) \rightarrow -\infty$, $\hat{y}_i \rightarrow 0$
- Can we keep using MSE as the loss function?
 - Brier score = $\frac{1}{N} \sum_i (\textcolor{brown}{y}_i - \hat{y}_i)^2$
 - But this does not interpret logistic output as probability

Likelihood for logistic regression

- Likelihood: $P(y_i | x_i) = \hat{y}_i^{y_i} (1 - \hat{y}_i)^{1-y_i}$
 - y_i is either 0 or 1
 - When y_i is 0, the likelihood is $1 - \hat{y}_i$
 - When y_i is 1, the likelihood is \hat{y}_i
- Log likelihood: $y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$
 - This is the cross-entropy loss function!
 - Maximizing likelihood is the same as minimizing cross-entropy

Regularization of linear model

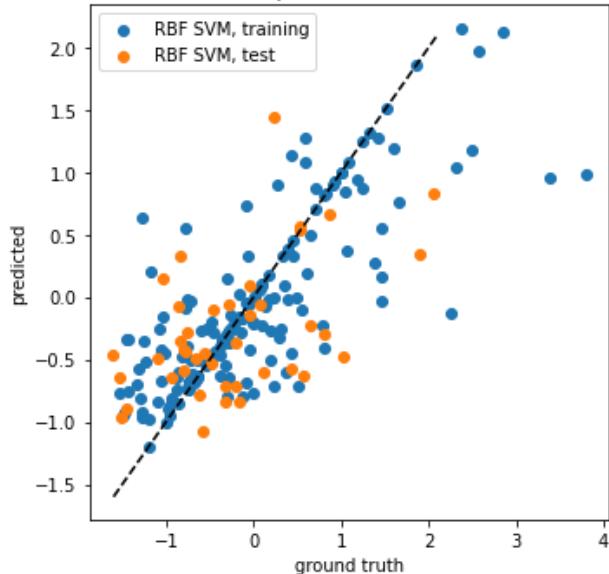
- L1 regularization (LASSO): $\text{MSE} + \alpha \sum_k |b_k|$
- L2 regularization (Ridge): $\text{MSE} + \alpha \sum_k b_k^2$
- α is the **hyperparameter** that controls the regularization strength
- Hyperparameter must be tuned for every dataset!

Tuning regularization strength



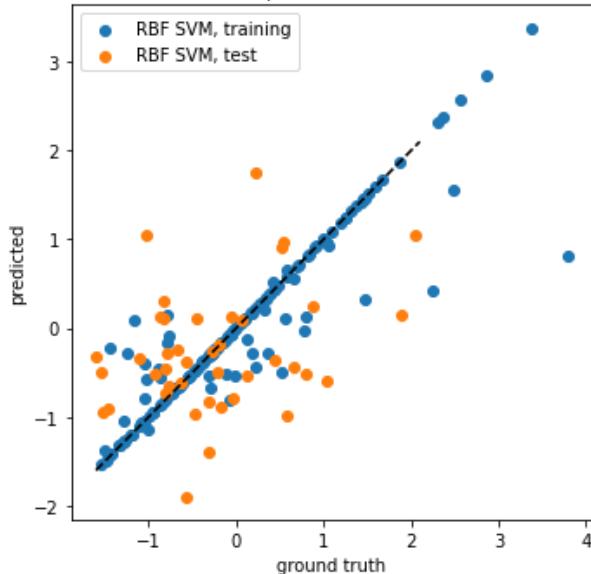
Just right

C: 1.0, epsilon: 0, MSE: 0.5305



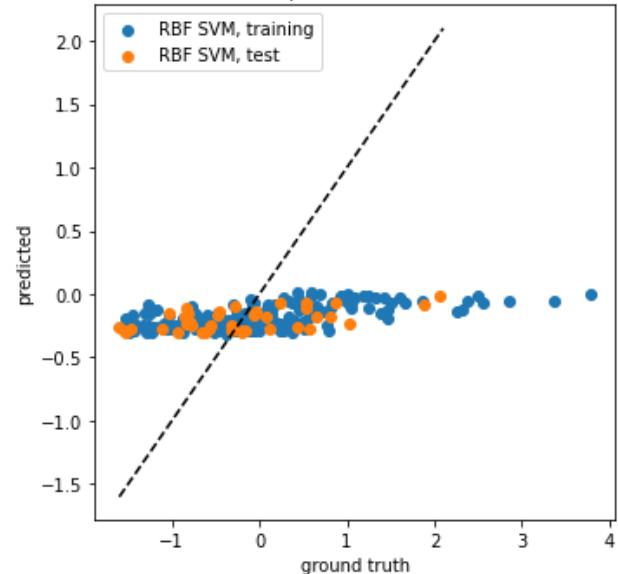
Too weak

C: 10.0, epsilon: 0, MSE: 0.8129



Too strong

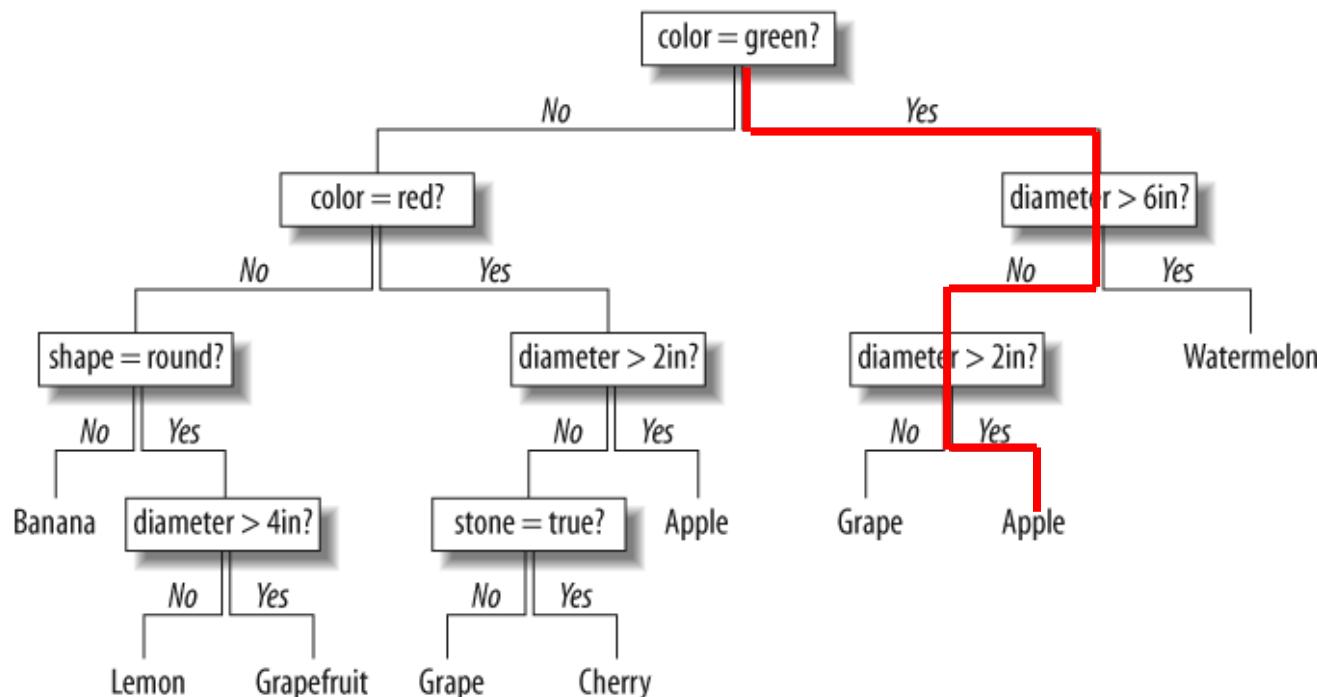
C: 0.01, epsilon: 0, MSE: 0.6388





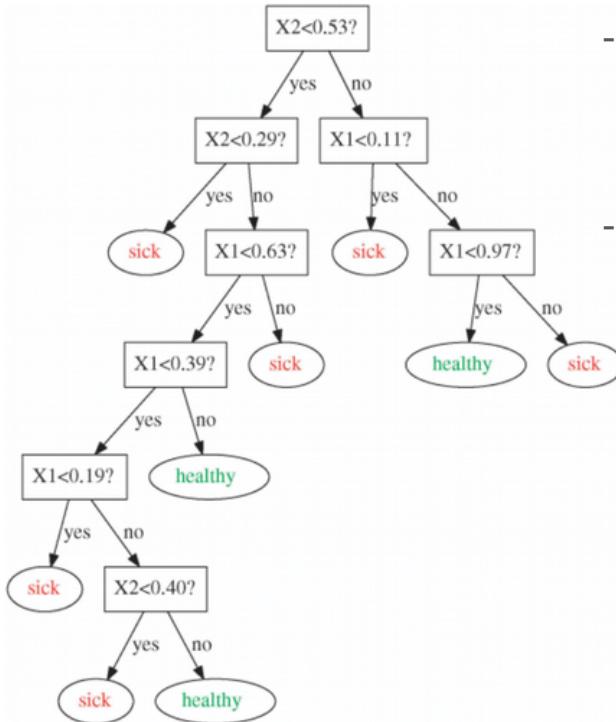
Decision tree

Decision tree

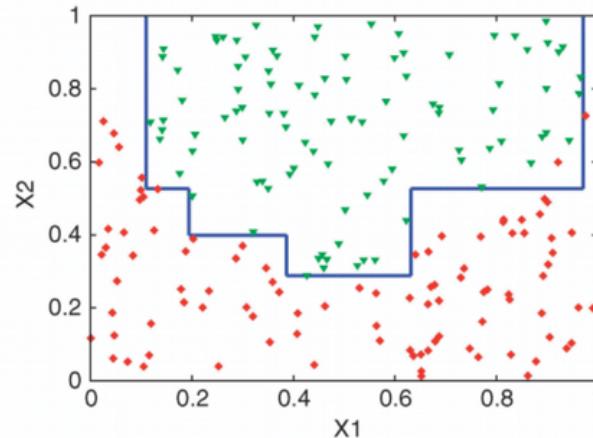


Source: Programming collective intelligence by Toby Segaran

Decision tree behaviors



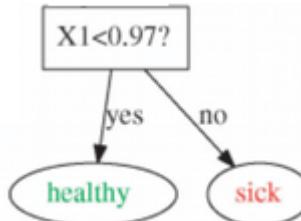
- Each decision is a threshold on each feature
 - Piecewise linear
 - Parallel to an axis
- Good for criteria-based classification



Splitting quality

100 healthy, 100 sick

85 healthy, 5 sick
45 healthy, 51 sick

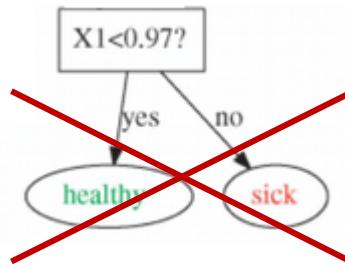


15 healthy, 95 sick
55 healthy, 49 sick

- Gini impurity: $\sum p(1 - p)$
- Entropy: $-\sum p \ln(p)$
 - Minimal at $p = 0$ or $1 \rightarrow$ Perfect split
 - Maximal at $p = 0.5 \rightarrow$ 50-50 split
- Search for feature and cutoff that yield lowest impurity or entropy

Control mechanisms for tree building

1. Too few samples to make a split

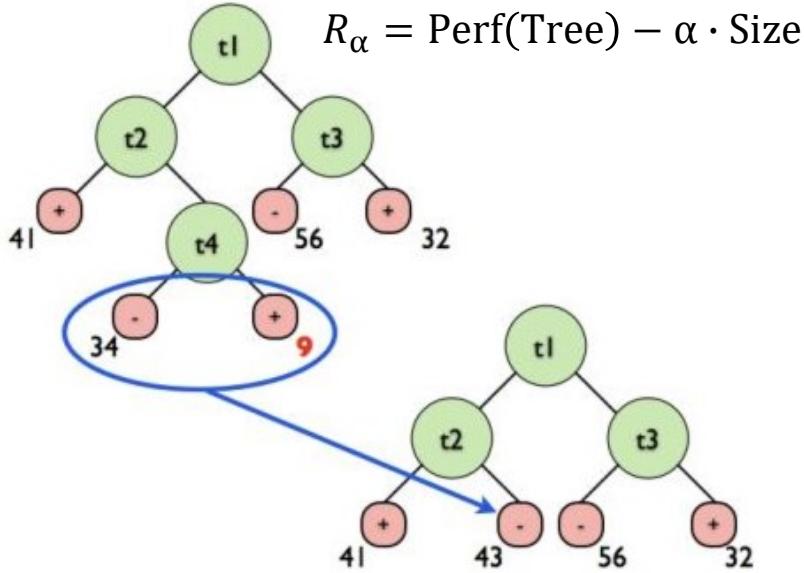


2. Too few samples on either branch

3. Impurity or entropy does not change much after the split

- Limit the tree size
- Limit the improvement in quality
- Limit the number of samples that support a split

Tree pruning (post-processing)



Patel, N. and Upadhyay, S. "Study of Various Decision Tree Pruning Methods with their Empirical Comparison in WEKA" IJCA 2012

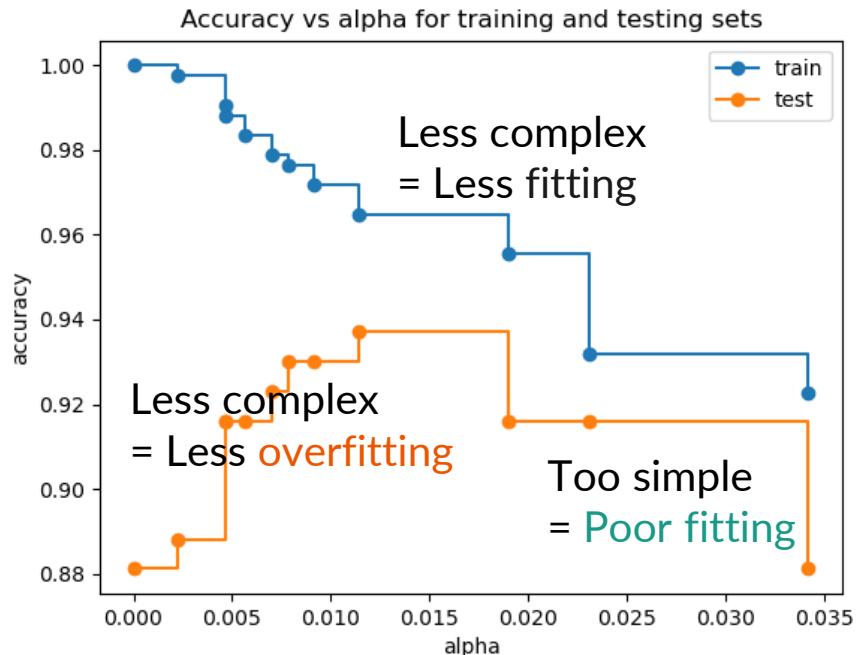
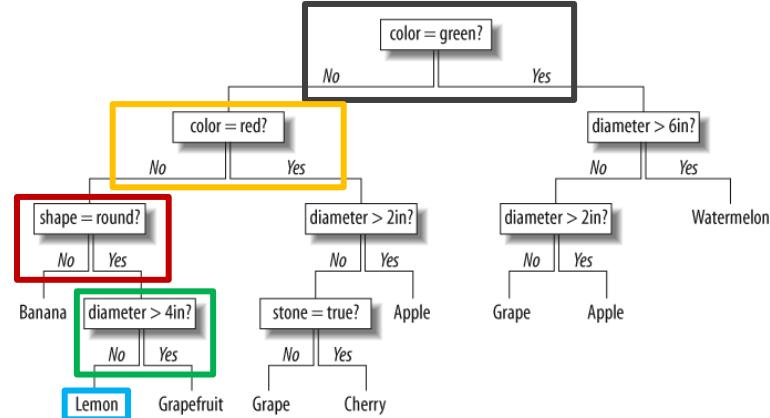


Image from scikit-learn.org

Regularization on features

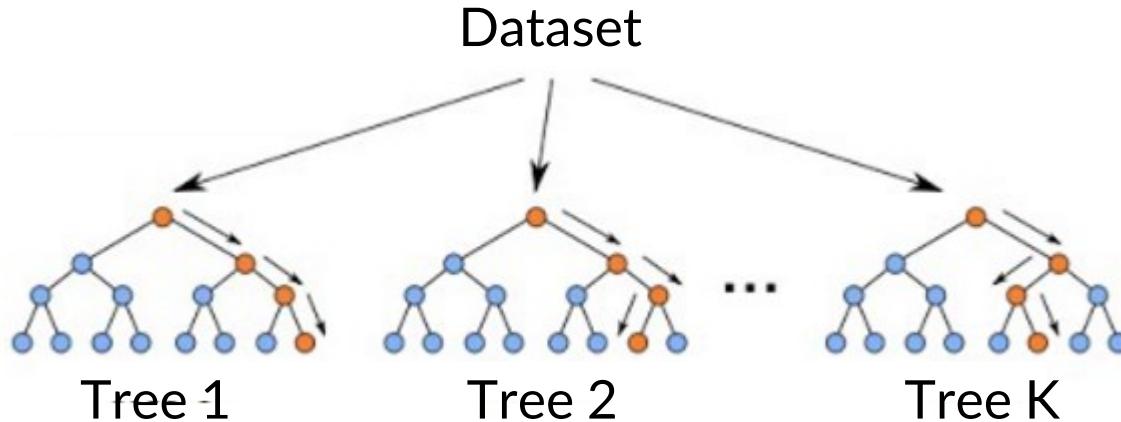
- Linear model: $\hat{y}_i = b_0 + b_1 x_{i,1} + \cdots + b_n x_{i,n}$
 - LASSO
- Tree model:
 - Repeatedly using the same feature
 - Early decision affects the rest
- Feature bagging
 - Look at only N features at each step
 - Force model to use diverse features





Ensemble approaches

Bagging: Random forest



- Sample 80% of the dataset to train each decision tree
- Each tree may overfit to different part of the dataset
- But the consensus should be correct

Boosting

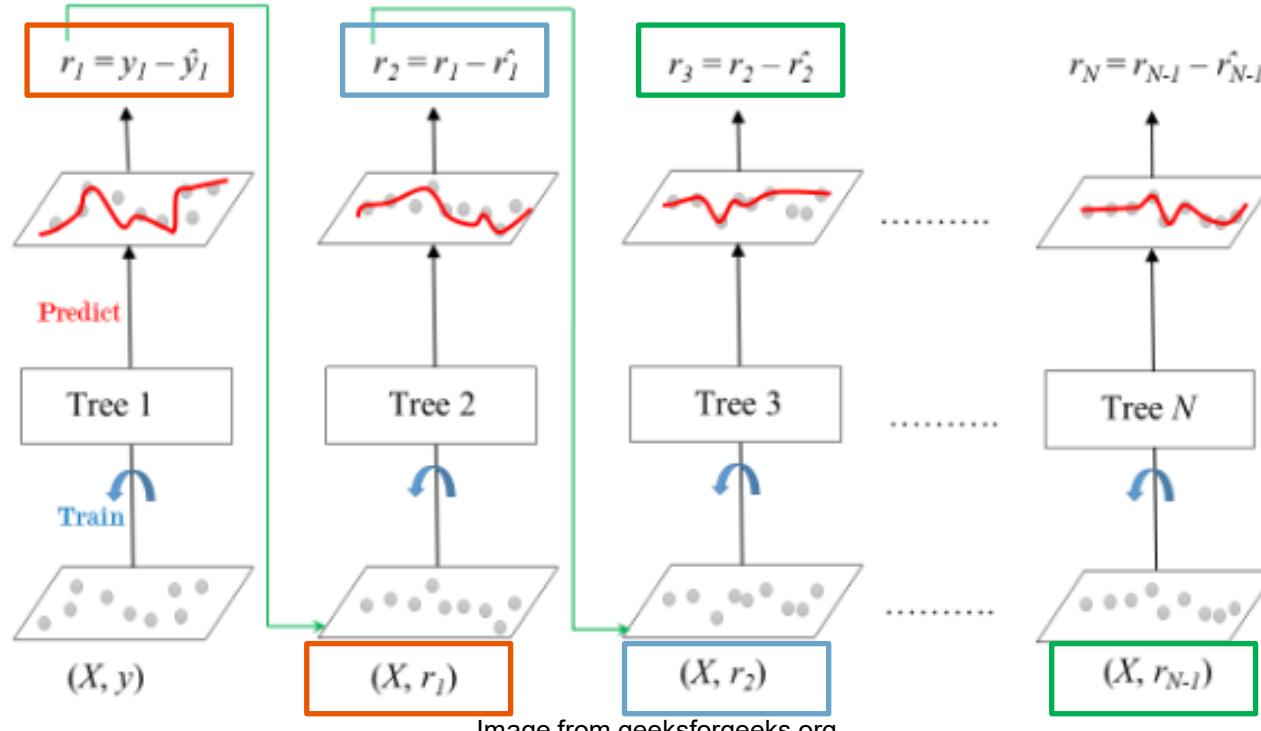
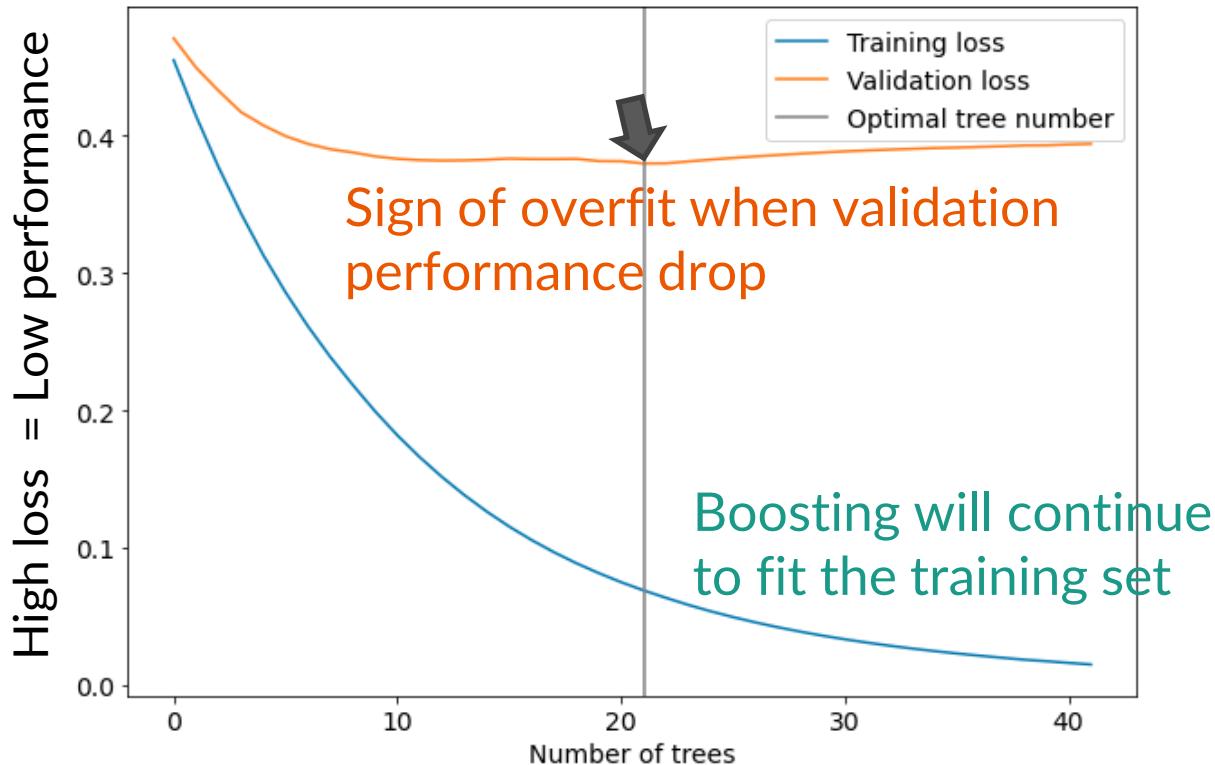


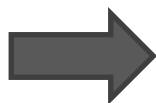
Image from geeksforgeeks.org

Controlling the boosting process



Impact of ensemble

Boosting solves
underfitting



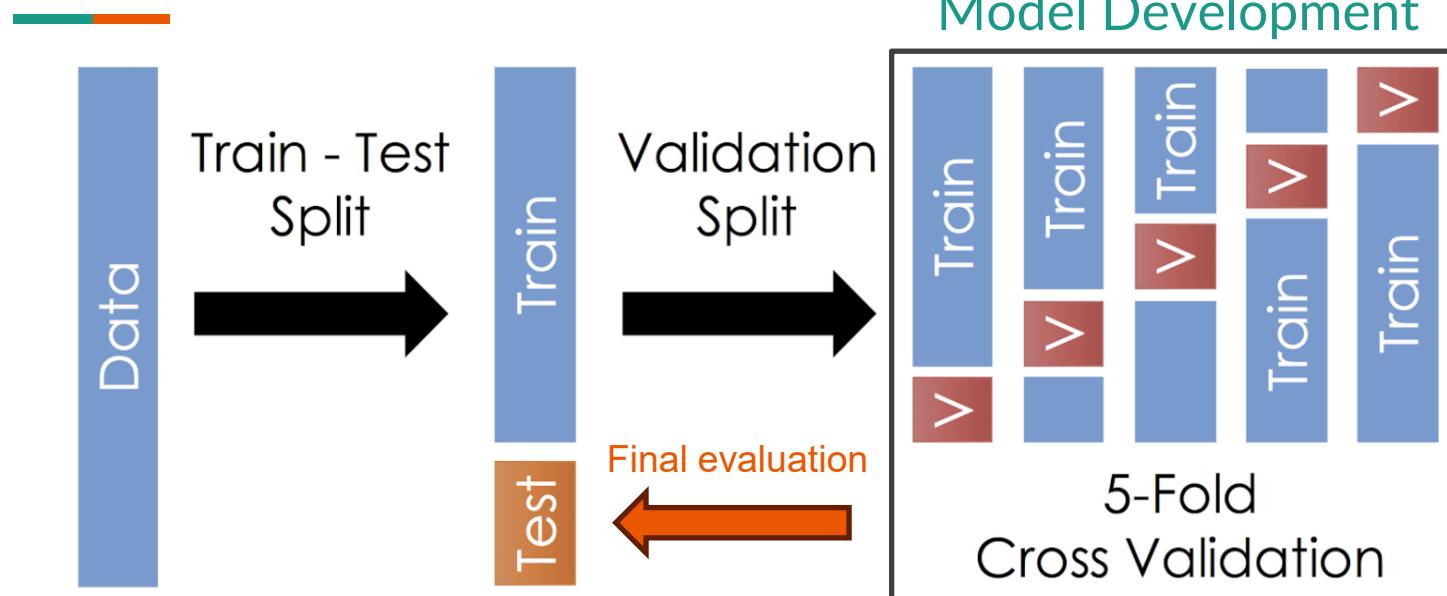
Bagging prevent
overfitting





Model evaluation

Train-Val-Test

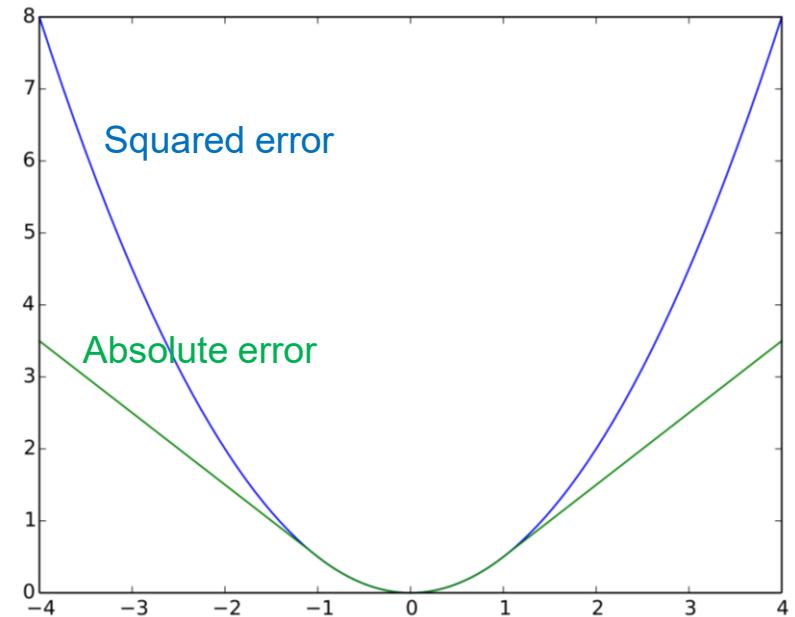


- **Training** data determines the best coefficients / weights
- **Validation** data determine the best hyperparameters
- **Test** data determine performance on new datasets

Source: medium.com

Performance metrics: regression

- Mean Square Error
- Mean Absolute Error
- Mean Absolute Percentage Error
- R^2 (Coefficient of Determination)



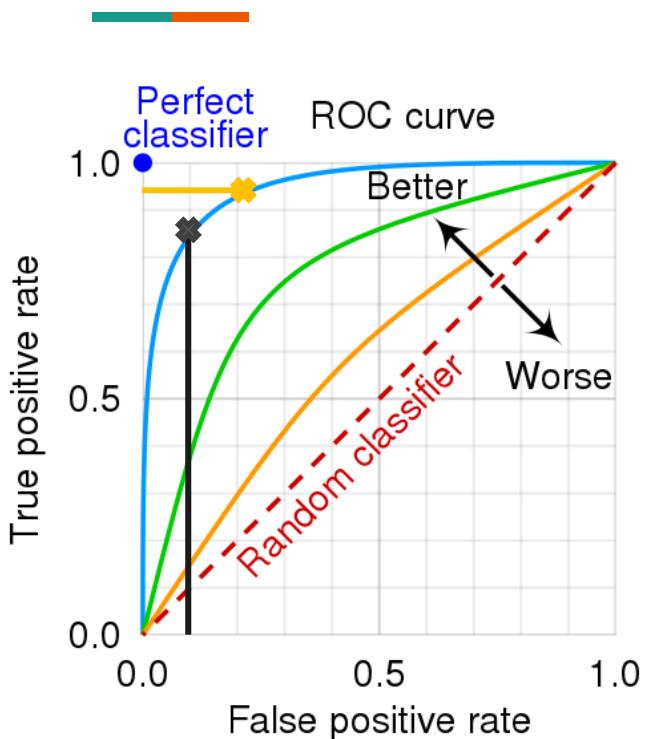
Performance metrics: classification



		Predicted	
		Negative	Positive
Actual	Negative	True Negative	False Positive
	Positive	False Negative	True Positive
		Predicted < 0.5	Predicted > 0.5

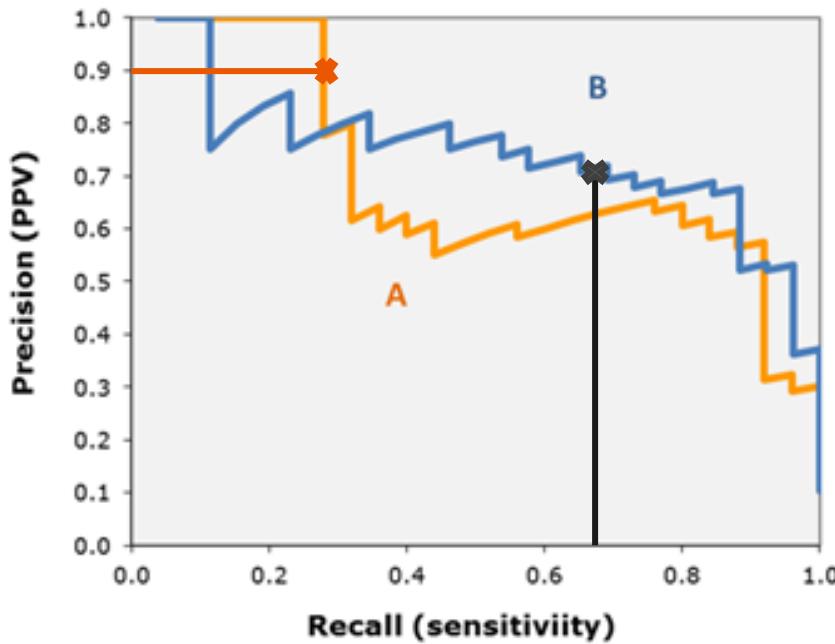
- Accuracy = $(TN + TP) / \text{total}$
- Precision = $TP / (TP + \text{FP})$ = Positive predictive value
- Recall = $TP / (TP + \text{FN})$ = Sensitivity
- Specificity = $TN / (TN + \text{FP})$

Threshold-free metrics



- Sensitivity-specificity at every output threshold
- Area under the ROC curve (AUROC, AUC)
 - Random guess = 0.5
 - Perfect model = 1.0
- Pick threshold based on use case
 - Specificity >0.9
 - Sensitivity >0.9

Precision-Recall curve



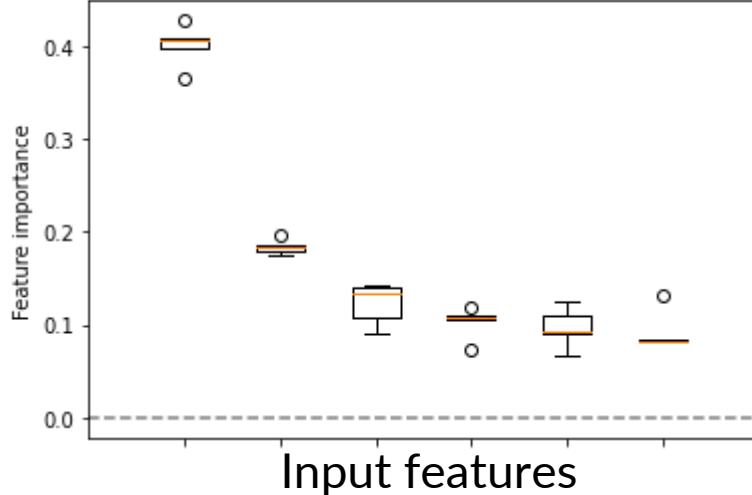
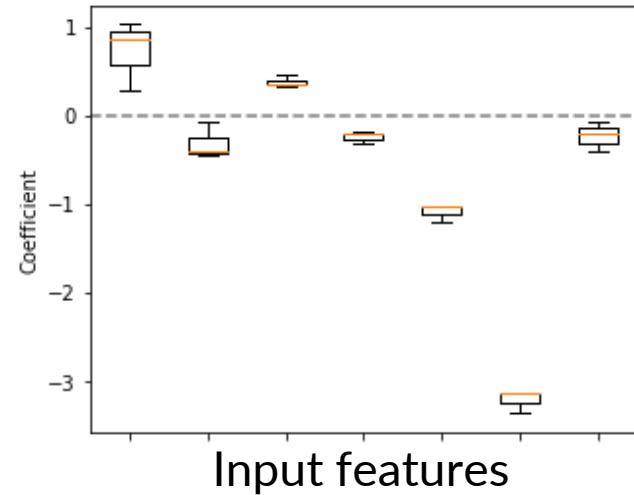
<https://acute caretesting.org/en/articles/precision-recall-curves-what-are-they-and-how-are-they-used>

- The best model can depend on use case



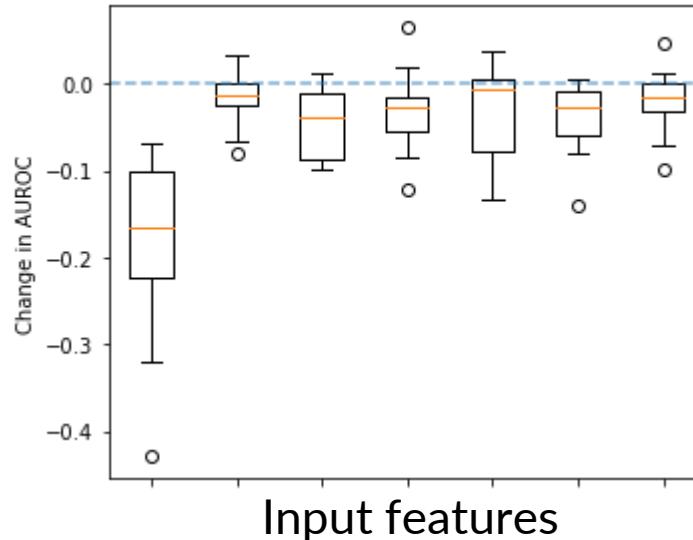
Explainability

Feature importance



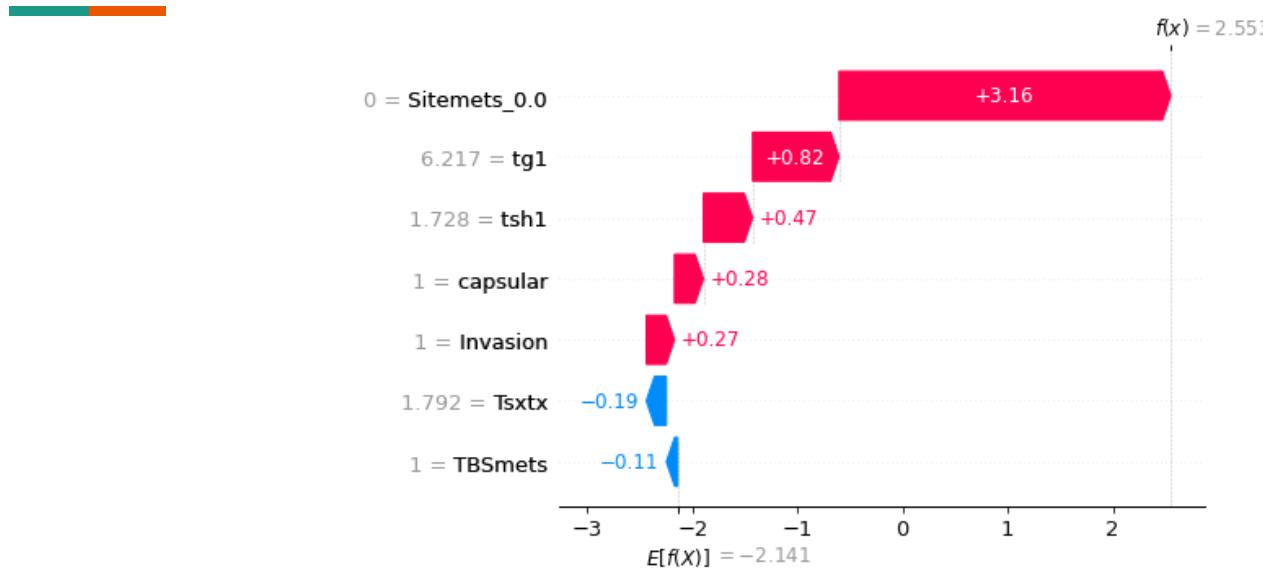
- Coefficients of linear, logistic, and SVM models
- Average improvement in impurity or entropy in tree models
- Model-level explanation

Change in performance after dropping a feature



- Compare performance with and without each input feature
- Big drop = important

Shapley value



- Change in predicted value due to the addition of a feature i
 - $\varphi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(n-|S|-1)!}{n!} [v(S \cup \{i\}) - v(S)]$
- Sample-level explanation

Key points

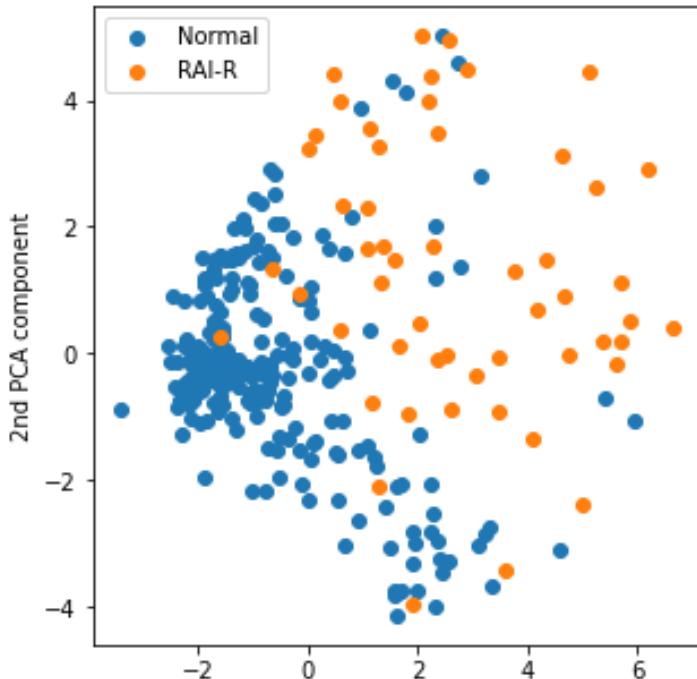
- Machine learning vs statistics = **data driven** vs **hypothesis driven**
- ML components = **model architecture + objective + learning algorithm**
- Heart of ML = balancing between data fitting and generalization
- Explainability ensures that the model learns the right knowledge



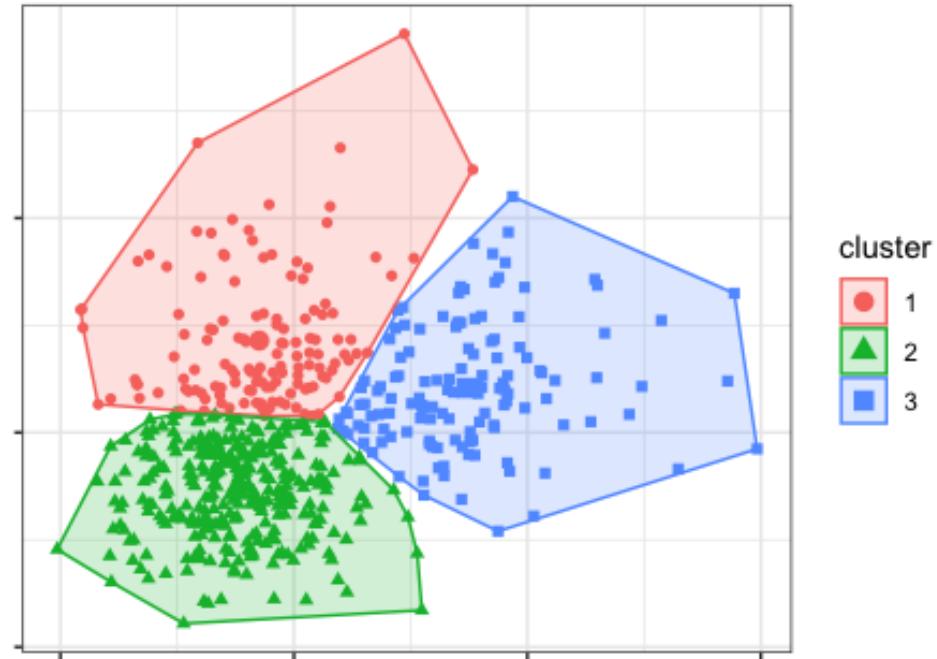
Unsupervised learning

Two primary branches of unsupervised learning

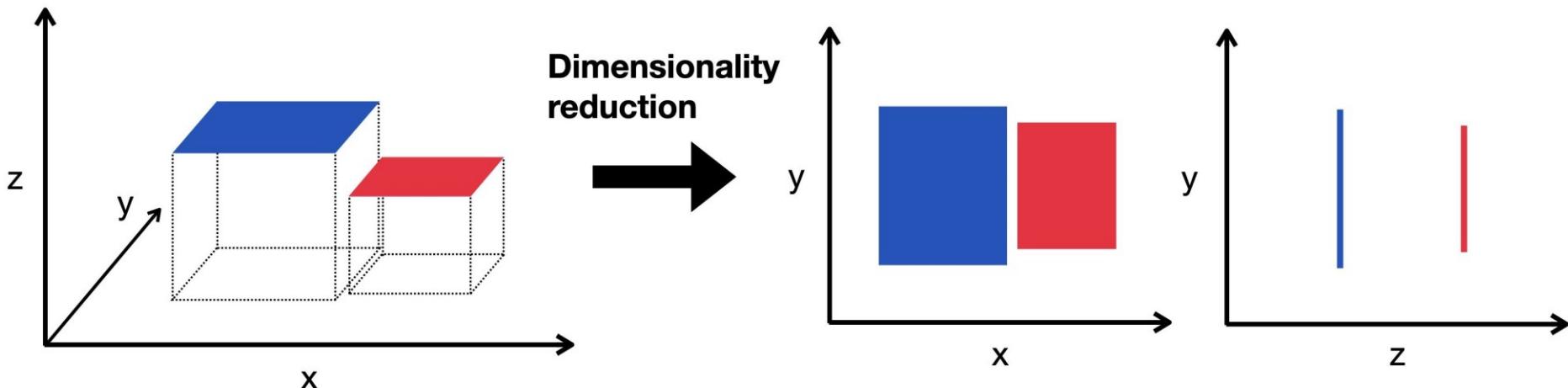
Dimensionality Reduction



Clustering



Dimensionality reduction



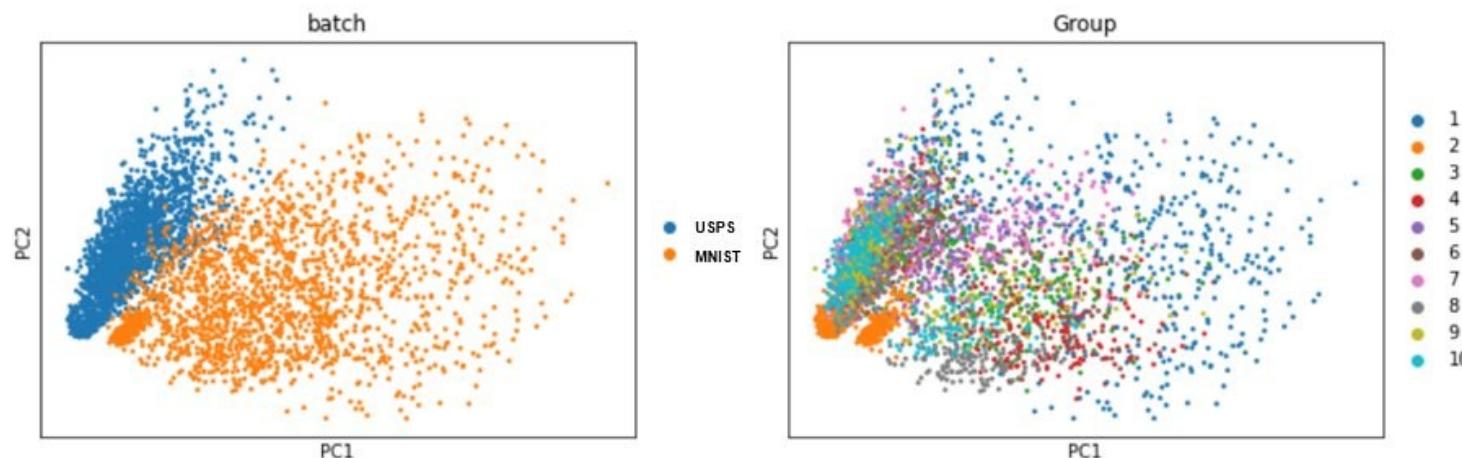
https://www.sc-best-practices.org/preprocessing_visualization/dimensionality_reduction.html

- Reduce dimension (number of features) while maintaining information
- Patient with similar symptoms also exhibit similar lab tests or have similar demographics or similar medical history

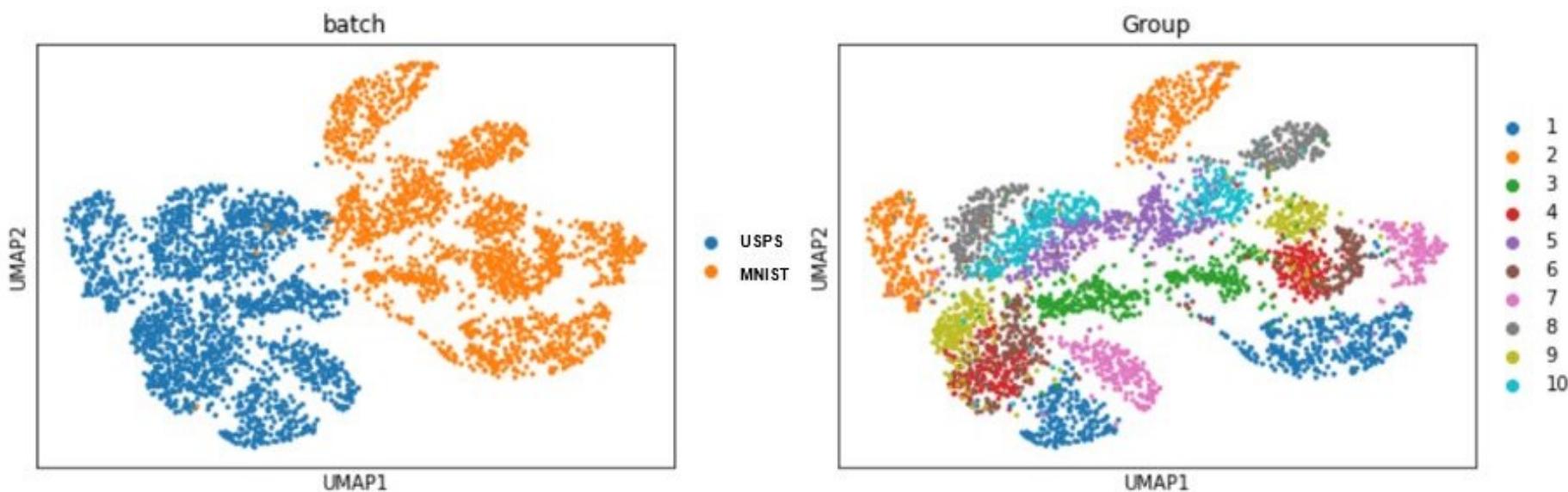
Digit datasets



*adapted from doi:10.1109/TKDE.2017.2669193



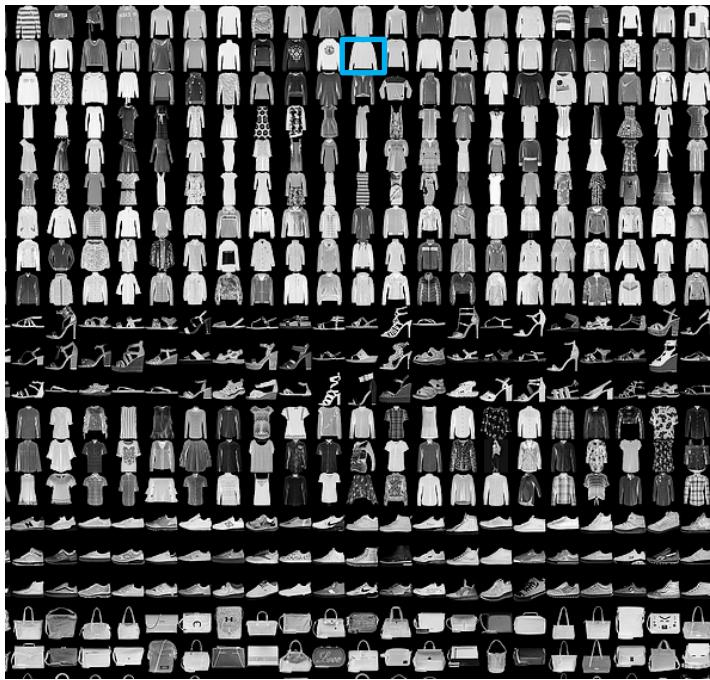
UMAP captures every group in 2D



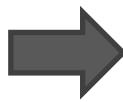
<https://twitter.com/lkmklsmn/status/1436357177887895555>

- Both data source and digit identity can be distinguished

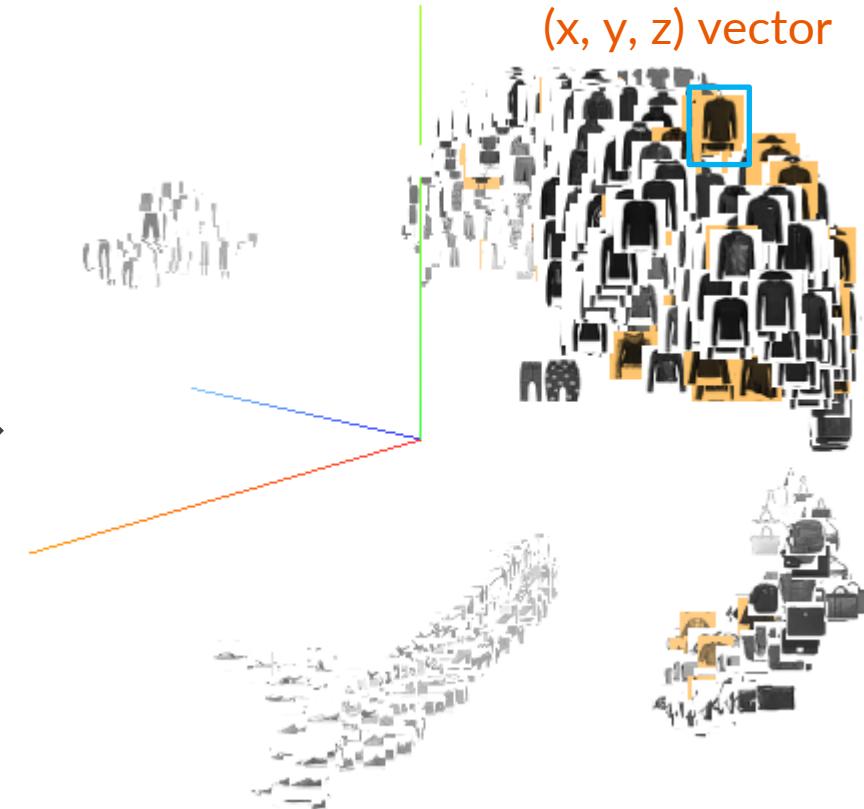
Fashion MNIST



28x28 pixel image



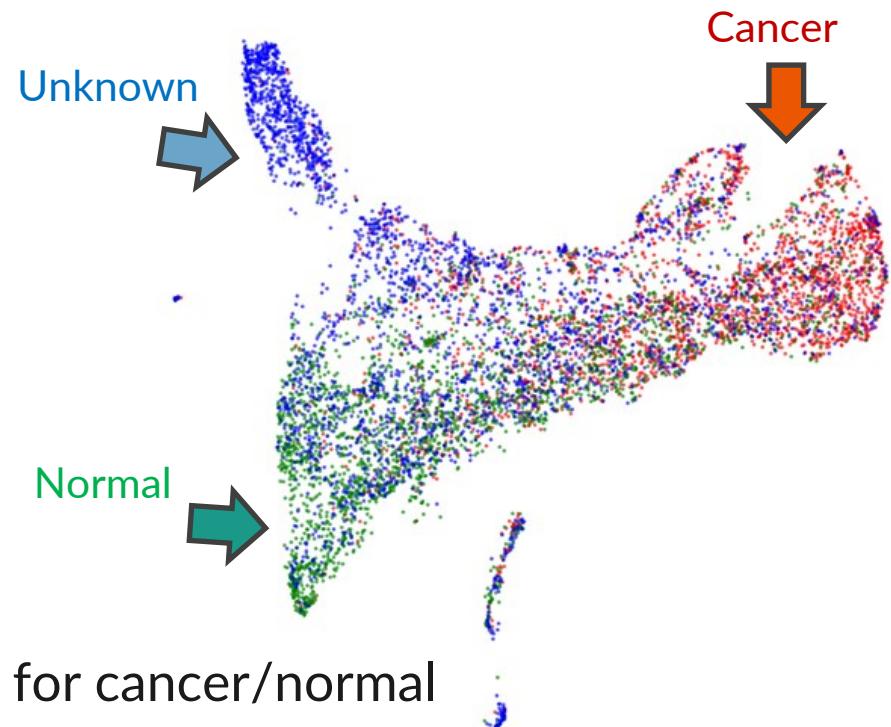
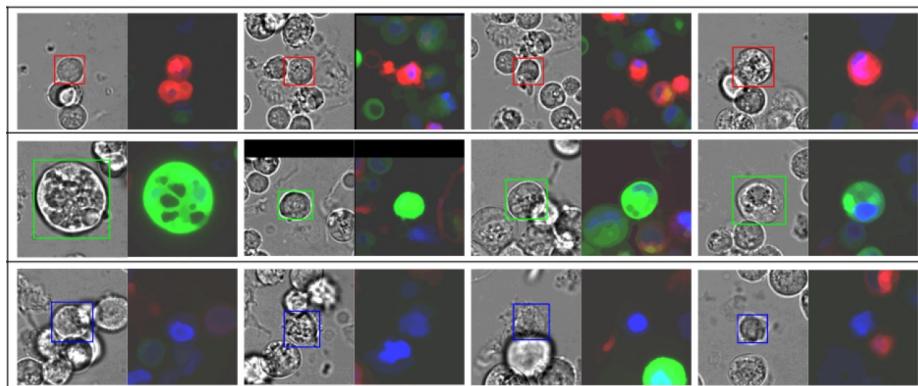
(x , y , z) vector



<https://github.com/zalandoresearch/fashion-mnist>

<https://pair-code.github.io/understanding-umap/>

2D visualization for cell images

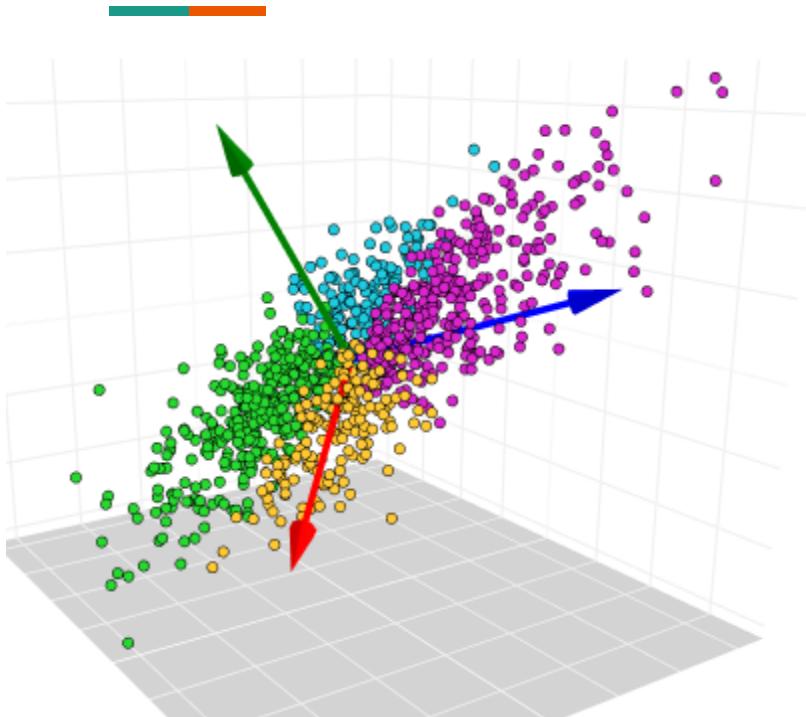


- Extracted from **deep learning model** for cancer/normal cell type classification

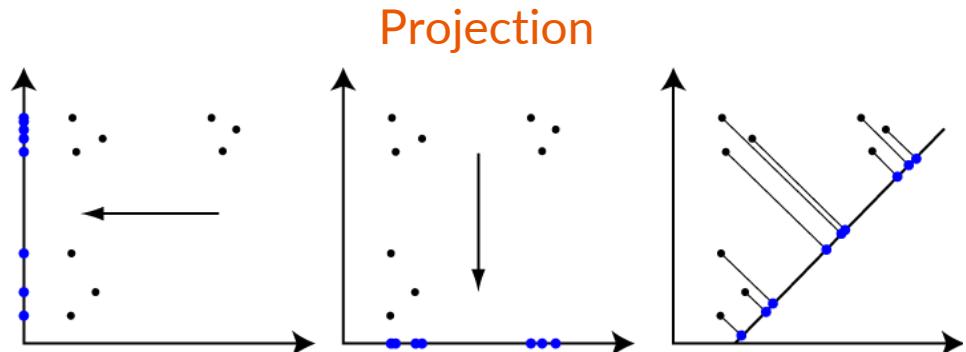


Principal component analysis (PCA)

Variance is information



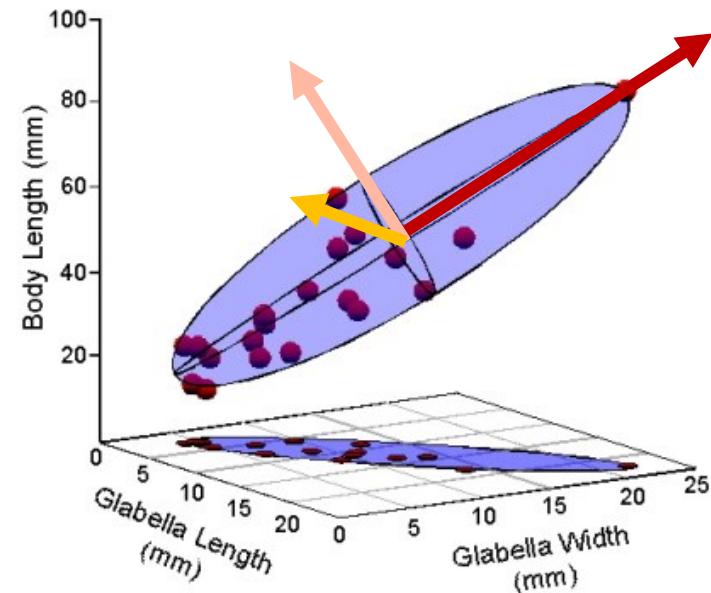
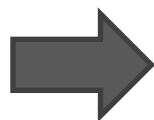
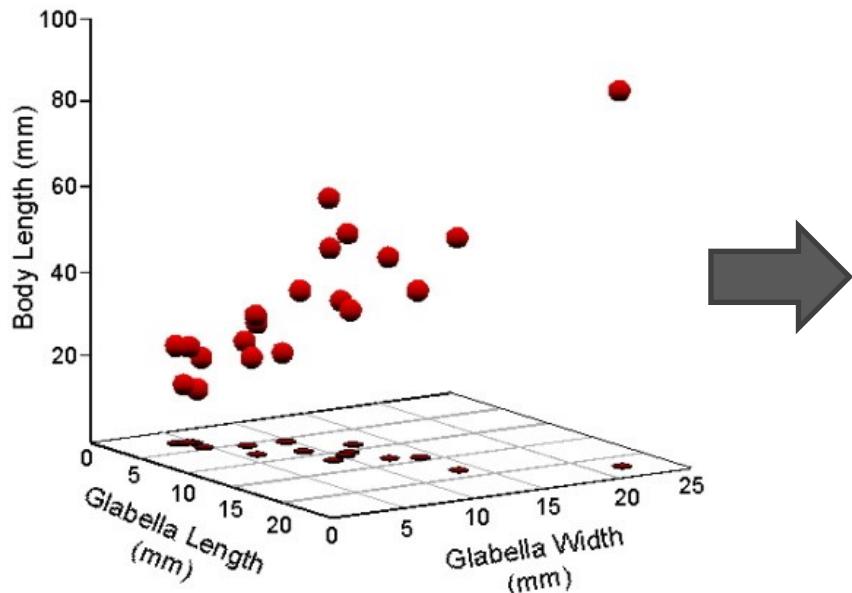
<https://towardsdatascience.com/principal-component-analysis-pca-explained-visually-with-zero-math-1cbf392b9e7d>



<https://shapeofdata.wordpress.com/2013/04/16/visualization-and-projection/>

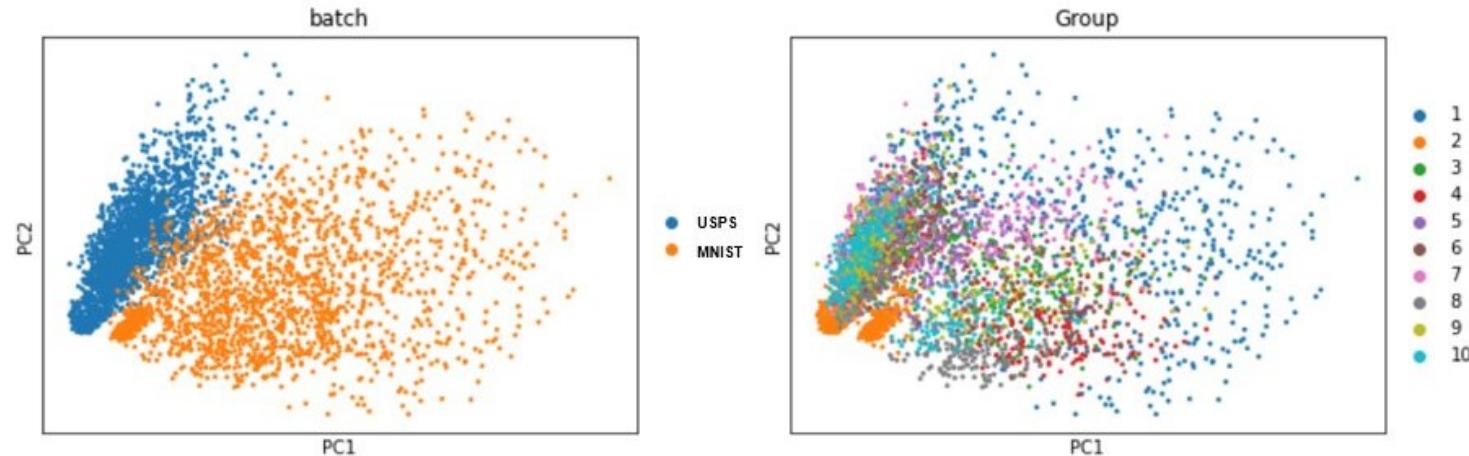
- High variances = more power to distinguish groups of data points

PCA prioritizes directions with high variances



Source: the paleontological association

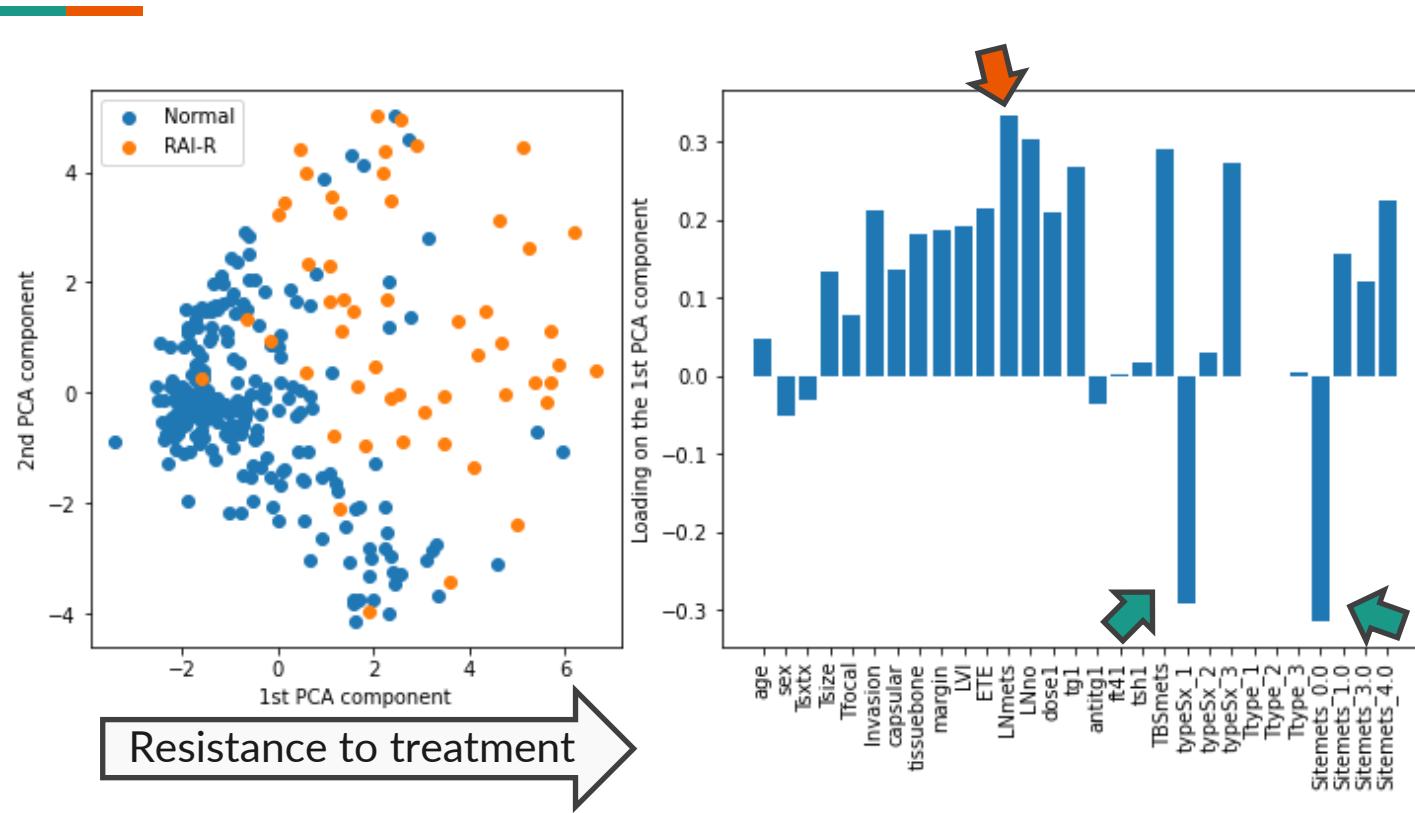
Interpretation of PCA result



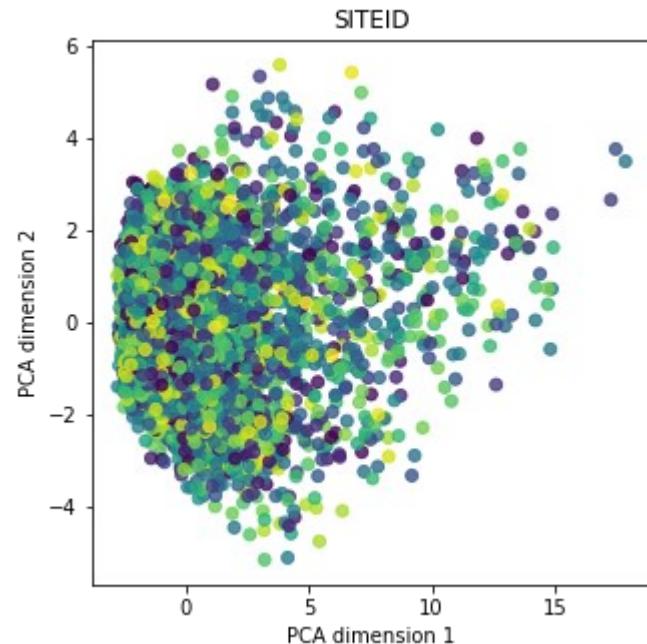
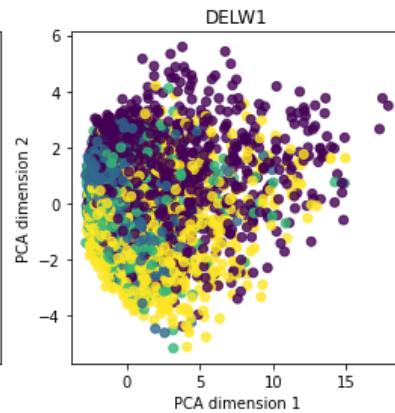
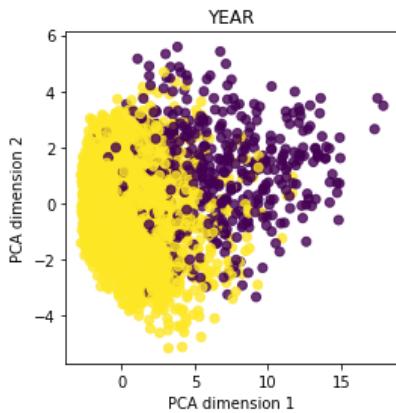
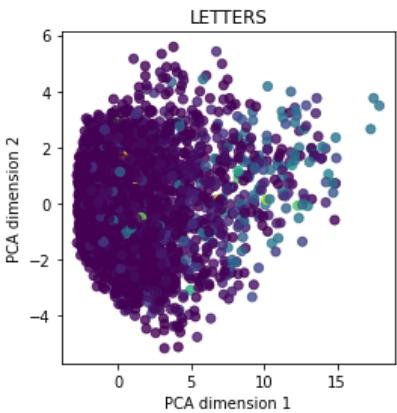
<https://twitter.com/lkmklsmn/status/1436357177887895555>

- PC1 captures the variance between data sources
- PC2 somewhat captures the variance between digit identity

Interpreting loadings on individual PC

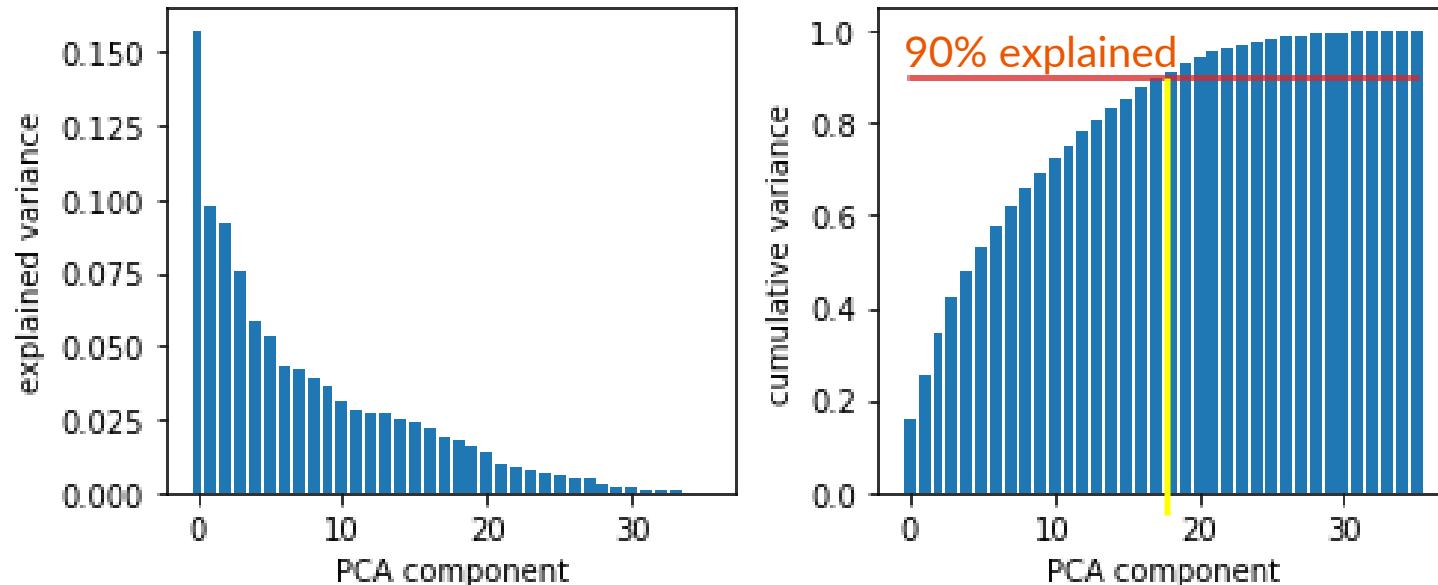


Exploring PCA results



- Color by feature values to understand how PCA group data points
- Color by potential confounding factors

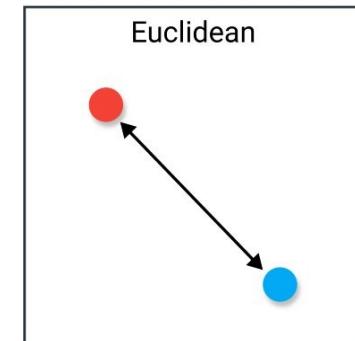
PCA for dimensionality reduction



- By default, PCA retains the number of dimensions
- We can **select only the first k PC** for downstream analyses

Pros and cons of PCA

- Each PC can be interpreted from the loadings
 - Highly correlated features tend to be grouped into the same PC
 - PCA is a good initial dimensionality reduction step
-
- PCA strictly preserves Euclidean distance
 - But some datasets require different distance metric!

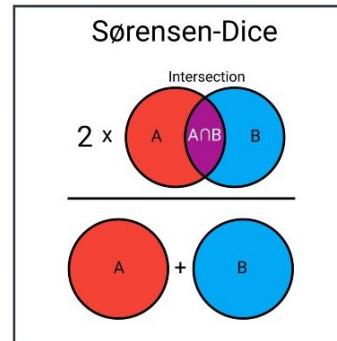
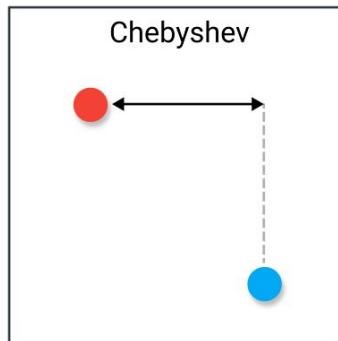
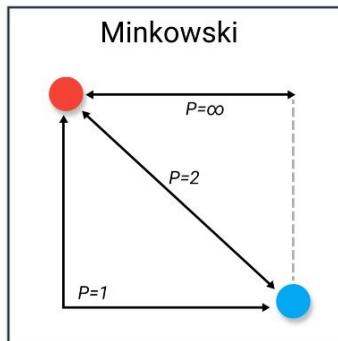
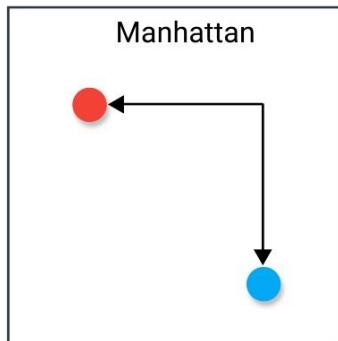
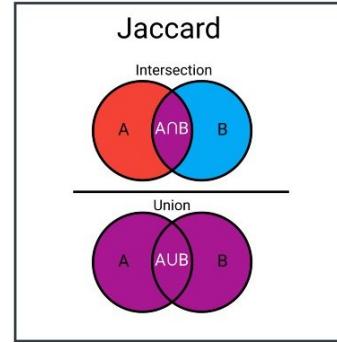
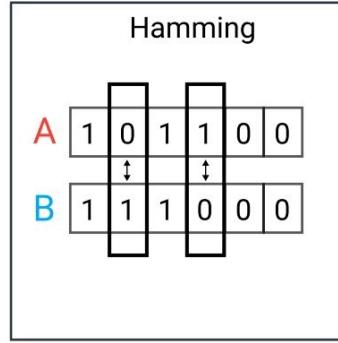
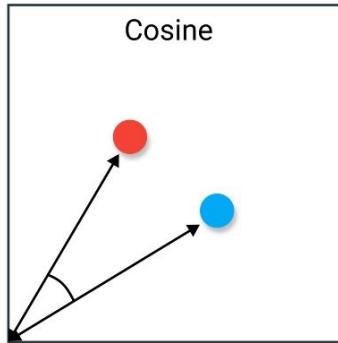
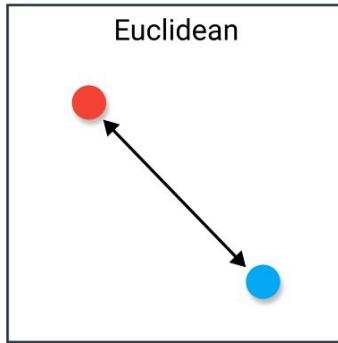


<https://towardsdatascience.com/9-distance-measures-in-data-science-918109d069fa>



Multidimensional Scaling (MDS)

Distances



Pairwise distance matrix

	A	B	C	D	E	F	G
A							
B	19.00						
C	27.00	31.00					
D	8.00	18.00	26.00				
E	33.00	36.00	41.00	31.00			
F	18.00	1.00	32.00	17.00	35.00		
G	13.00	13.00	29.00	14.00	28.00	12.00	

<http://www.slimsuite.unsw.edu.au/teaching/upgma/>

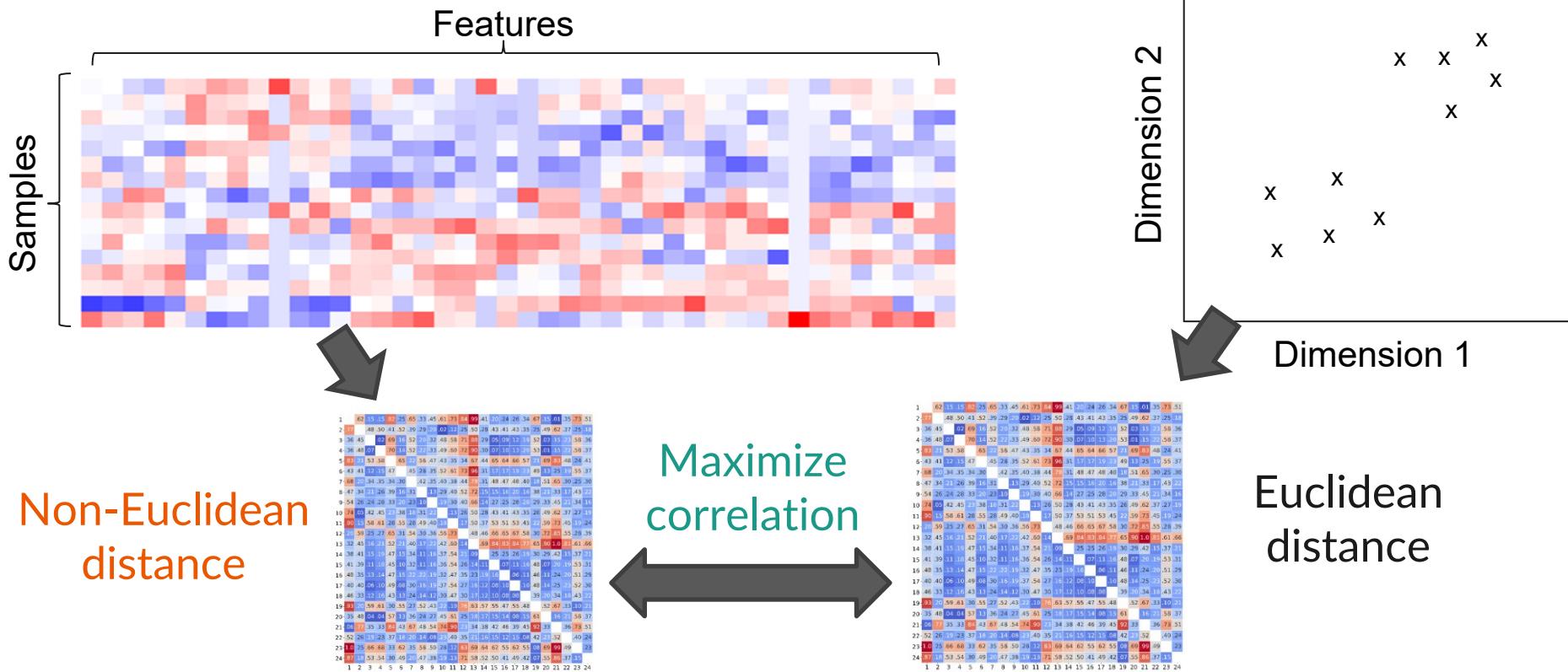
Metric properties

- $D(i, j)$ = distance between sample i and sample j
- $D(i, i) = 0$
- $D(i, j) = 0$ iff $i = j$
- $D(i, j) = D(j, i)$
- $D(i, j) + D(j, k) \geq D(i, k)$

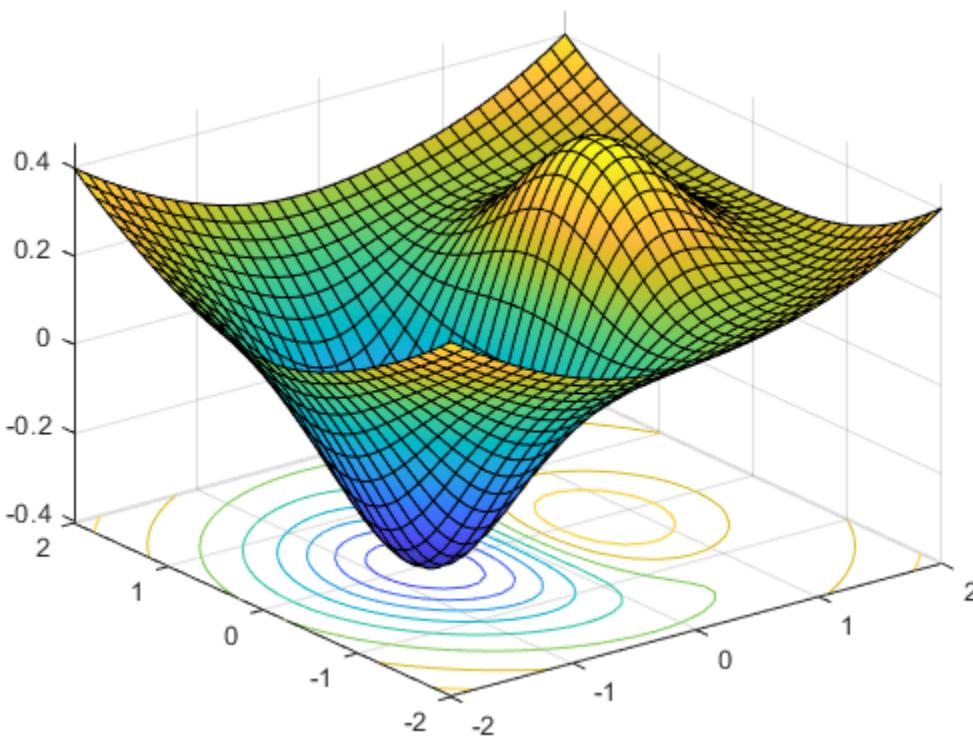
Non-metric

- Any user-defined dissimilarity
- $D(i, j) \neq D(j, i)$

Principal Coordinate Analysis (PCoA)

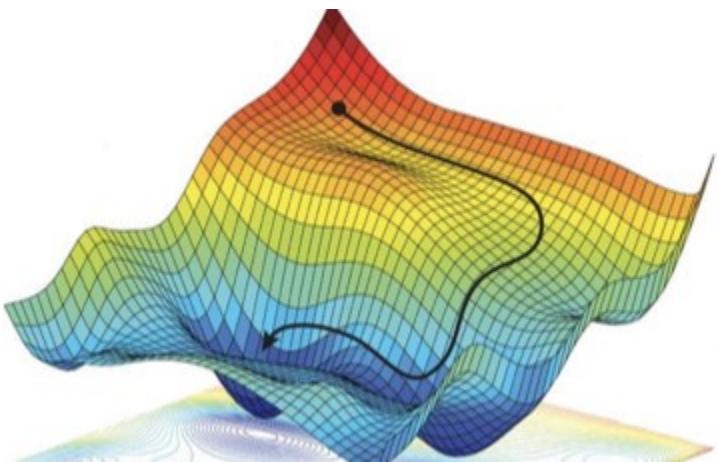


How to optimize a function?

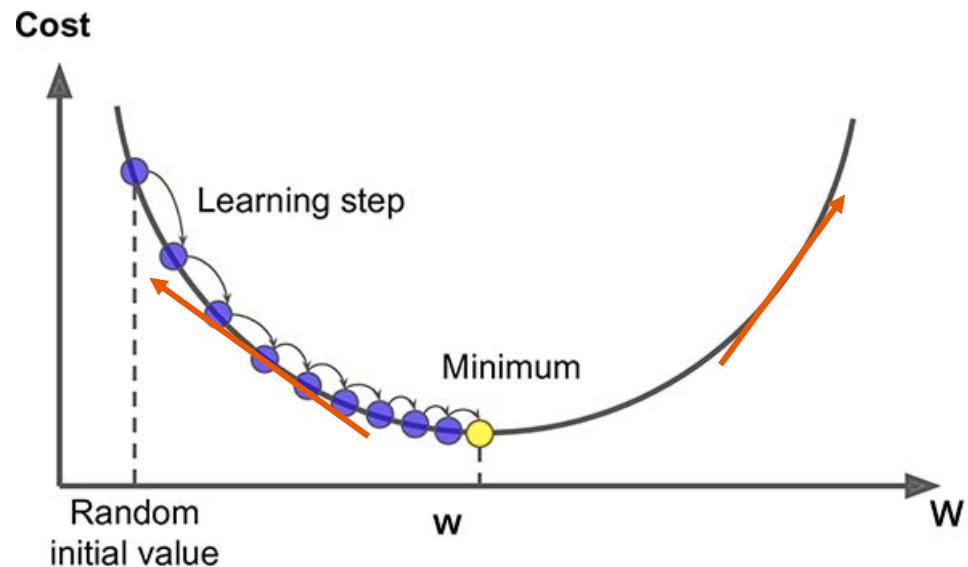


- Find (x_1, x_2, \dots, x_n) that minimize $f(x_1, x_2, \dots, x_n)$
- At minimum, the slope is zero in all directions
- Take derivative of each variable and set to zero
 - $\frac{\delta f}{\delta x_1} = 0$
 - $\frac{\delta f}{\delta x_2} = 0$
 - n equations with n variables

Gradient descent

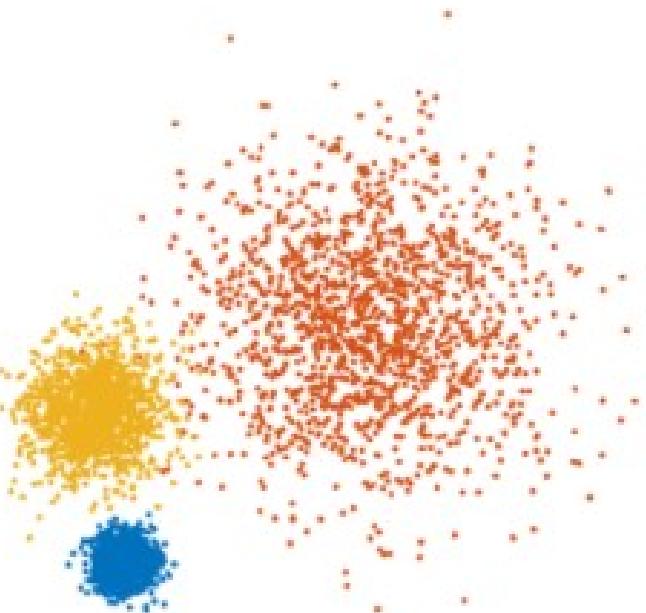


<https://medium.com/analytics-vidhya/gradient-descent-b0dc1af33517>



- Slope tells us if the function is increasing or decreasing if we increase x_i
 - So, we can update x_i accordingly

Limitation of PCA and MDS

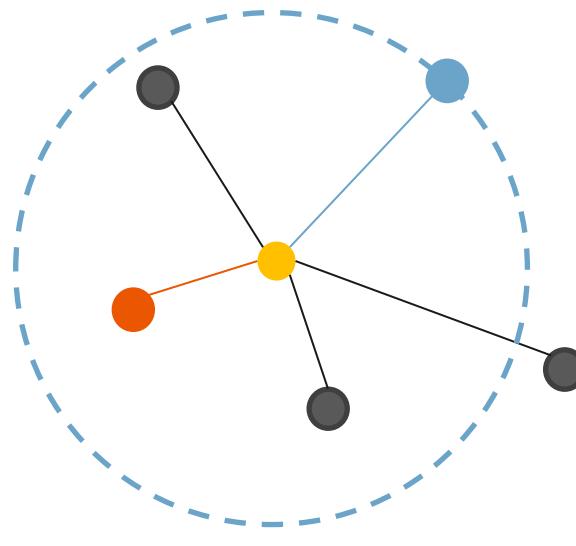
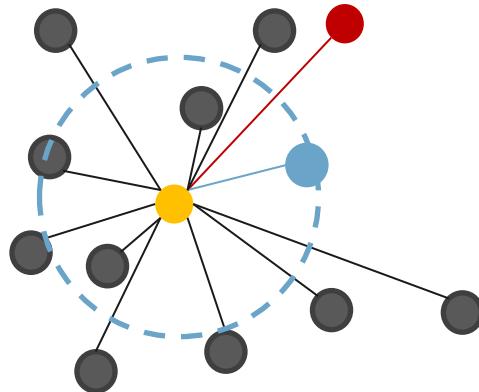


- A single definition of distance is used throughout the data space
- What if some data groups are noisier than the others?
 - Difference in data density



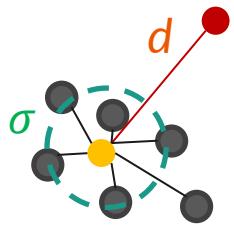
***t*-distributed stochastic neighbor embedding (*t*-SNE)**

Measuring data density



- Distance to the k -th nearest neighbor reflects data density
 - Small distance in dense area
 - Large distance in sparse area

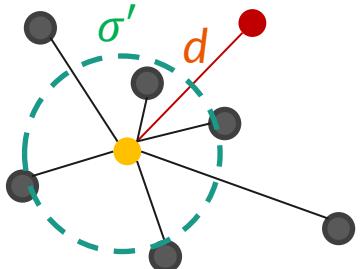
Probability of being a neighbor



$\text{score}(o | o) = \text{probability that } o \text{ would pick } o \text{ as neighbor}$
under a **normal distribution center at } o\text{**

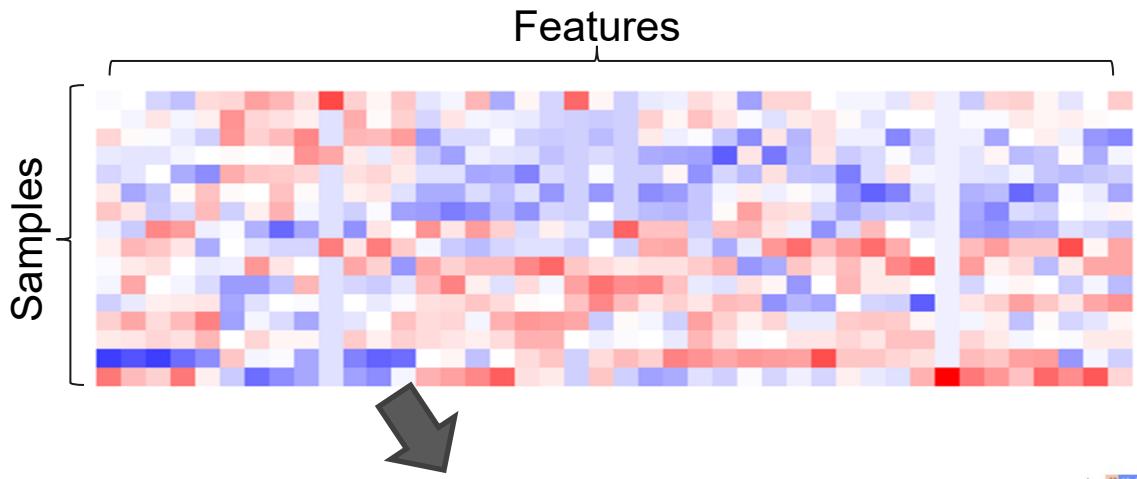
$$= \frac{e^{-\frac{d^2}{2\sigma^2}}/\sigma}{\sum e^{-\frac{(\text{dist}(o, o))^2}{2\sigma^2}}/\sigma}$$

o = other data points



- Same distance d normalized against density σ and distances to other nearby data points o

Finding the optimal projection for t -SNE



Probability of
being a neighbor
(Normal)
(σ depend on
density)

1	62	11	11	22	25	33	45	61	71	84	99	41	50	24	26	34	67	11	31	35	73	51		
2	52	48	30	41	50	22	29	02	11	29	50	28	43	41	47	33	29	50	28	41	41	37	29	33
3	38	45	32	69	16	52	22	37	48	58	71	89	29	55	10	18	52	03	13	33	58	36		
4	36	48	27	70	13	52	22	37	49	69	72	90	07	01	10	39	33	11	31	35	23	58	37	
5	41	73	53	47	52	22	47	43	35	34	67	44	52	06	41	46	52	21	69	81	48	24	41	
6	41	11	12	15	47	05	28	35	52	61	79	96	31	17	19	31	49	13	25	34	55	37		
7	41	34	31	35	42	02	42	35	52	61	79	96	31	17	19	31	49	13	25	34	55	37		
8	41	34	31	35	42	02	42	35	52	61	79	96	31	17	19	31	49	13	25	34	55	37		
9	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
10	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
11	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
12	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
13	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
14	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
15	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
16	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
17	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
18	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
19	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
20	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
21	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
22	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
23	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
24	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
25	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
26	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
27	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
28	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
29	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
30	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
31	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
32	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
33	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
34	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
35	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
36	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
37	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
38	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
39	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
40	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
41	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
42	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
43	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
44	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
45	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
46	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
47	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
48	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
49	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
50	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
51	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
52	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
53	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
54	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
55	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
56	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
57	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
58	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
59	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
60	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
61	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
62	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
63	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
64	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
65	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
66	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
67	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
68	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
69	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
70	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
71	54	28	28	33	20	11	11	35	40	66	50	37	28	29	35	45	52	31	35	34	36	38		
72	54	28	28	3																				

Why t -distribution for the projection?

- t -distribution has **fatter tails**
- Allow data points to be projected far away from each other

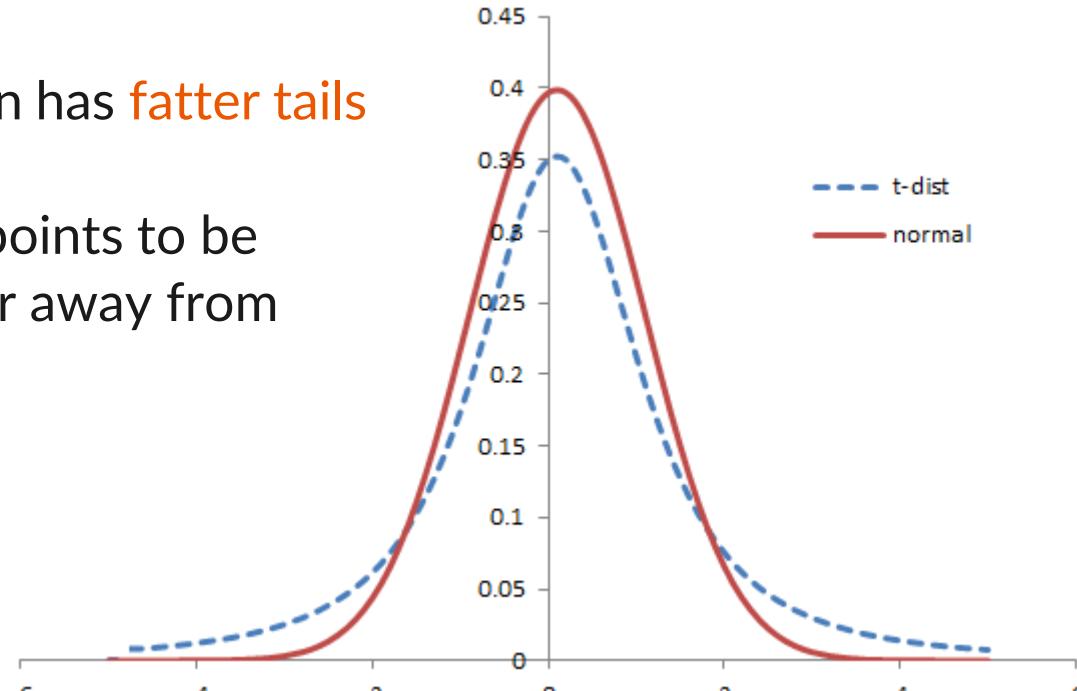
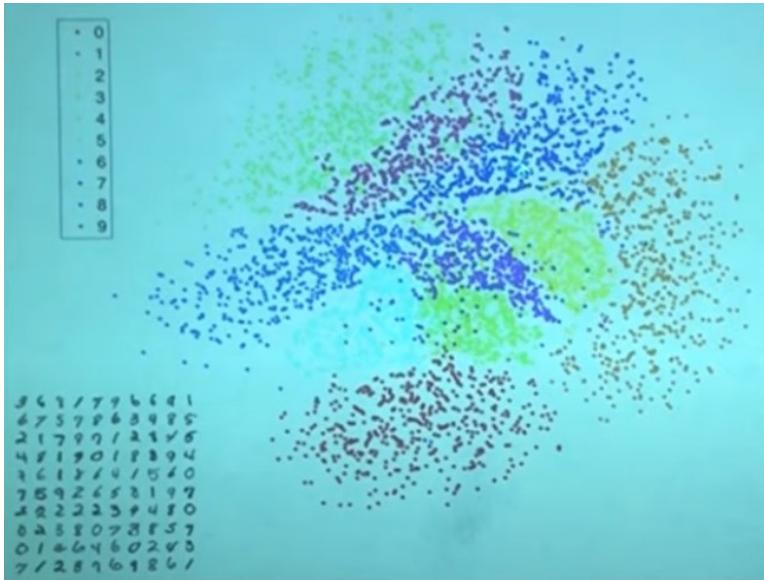
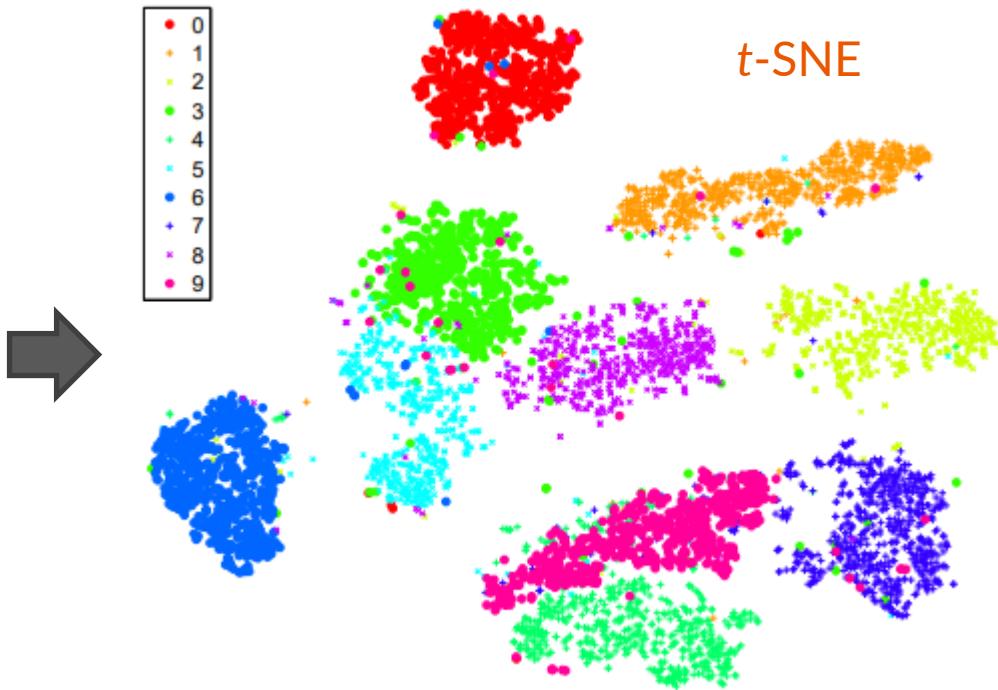


Image from riskprep.com

Impact of t -distribution

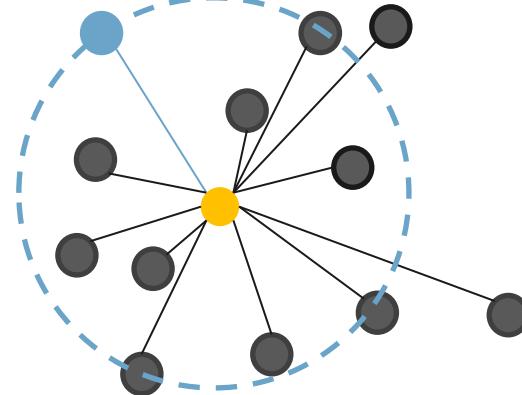
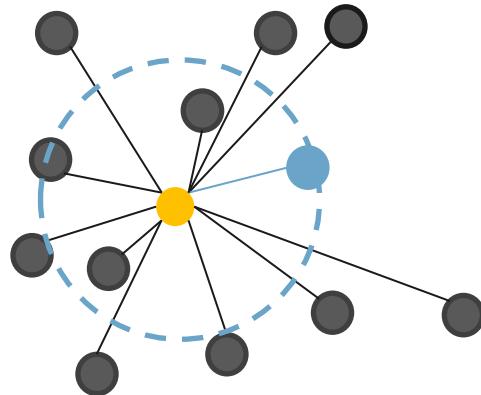


SNE (Normal → Normal)



Maaten, L. and Hinton, G. J of Machine Learning Research 9:2579-2605 (2008)

Perplexity

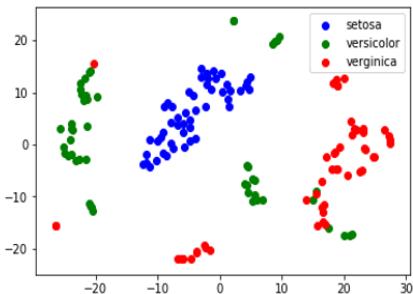


- How many nearest neighbors to consider to normalize data density?
 - Perplexity parameters

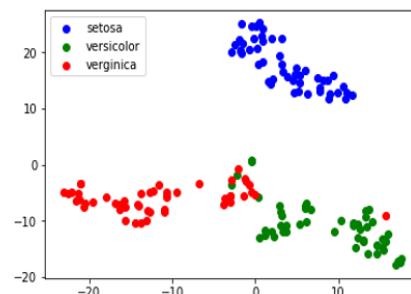
Impact of perplexity

—

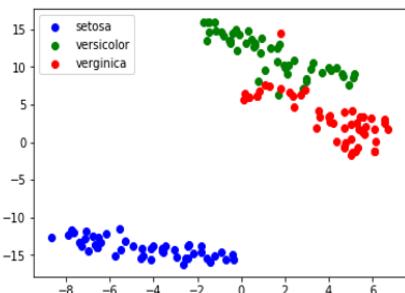
K = 5



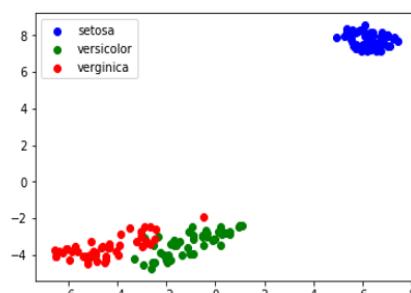
K = 10



K = 25



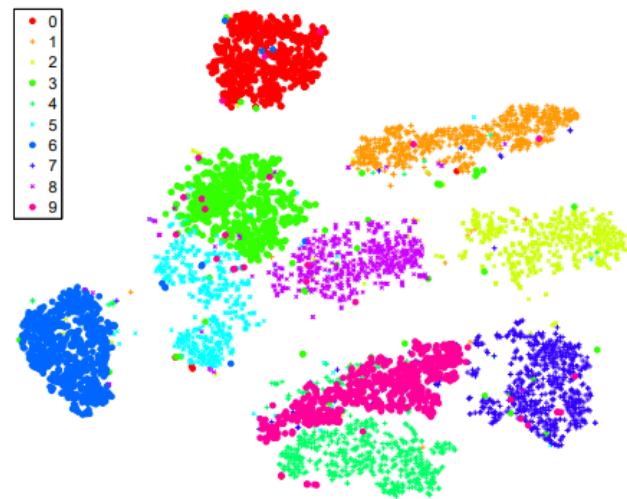
K = 50



- Too small perplexity = a lot of scattered data groups
- Try varying the perplexity and identify patterns that **consistently** appear

Pros and cons of *t*-SNE

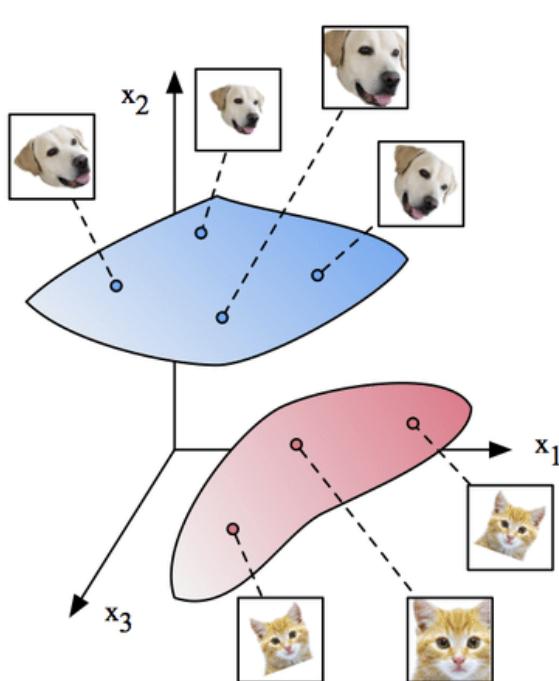
- Capture neighbor relationship
- Normalize data density
- Recompute every time new data is added
- Lose long-range relationship
- Axes of the resulting projection have no meaning
 - Don't use *t*-SNE coordinates for clustering or interpretation



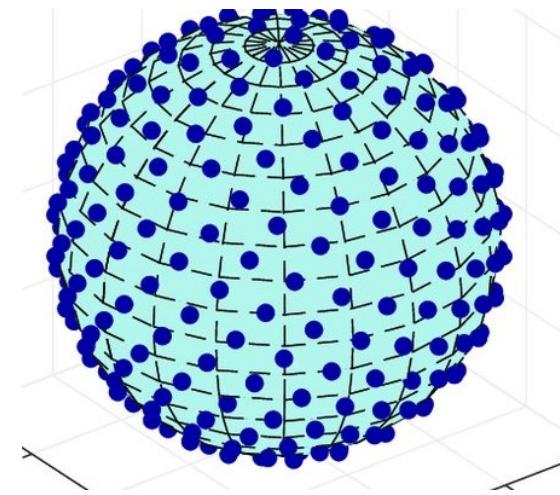


Uniform manifold approximation and projection (UMAP)

Two key assumptions



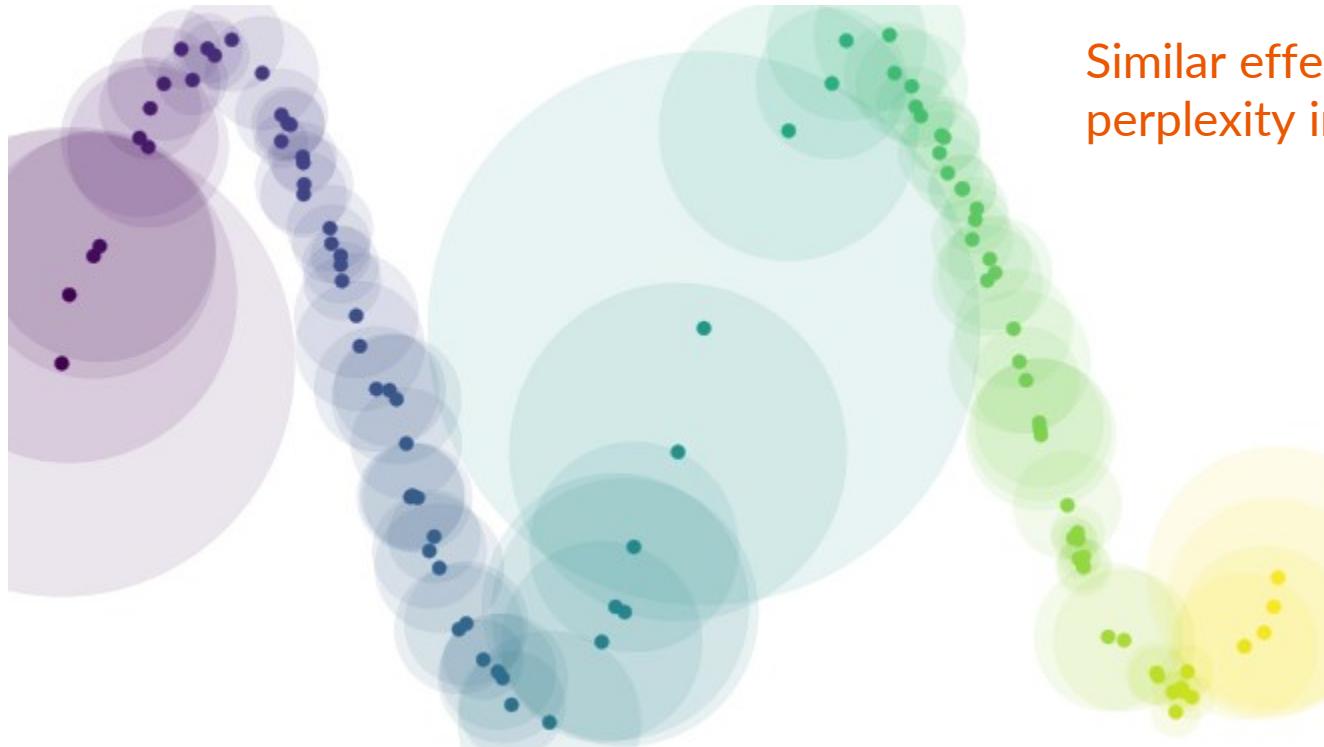
Chung, S. et al. "Classification and Geometry of General Perceptual Manifolds"



Ali, A. et al. IEEE Access PP(99):1 (2021)

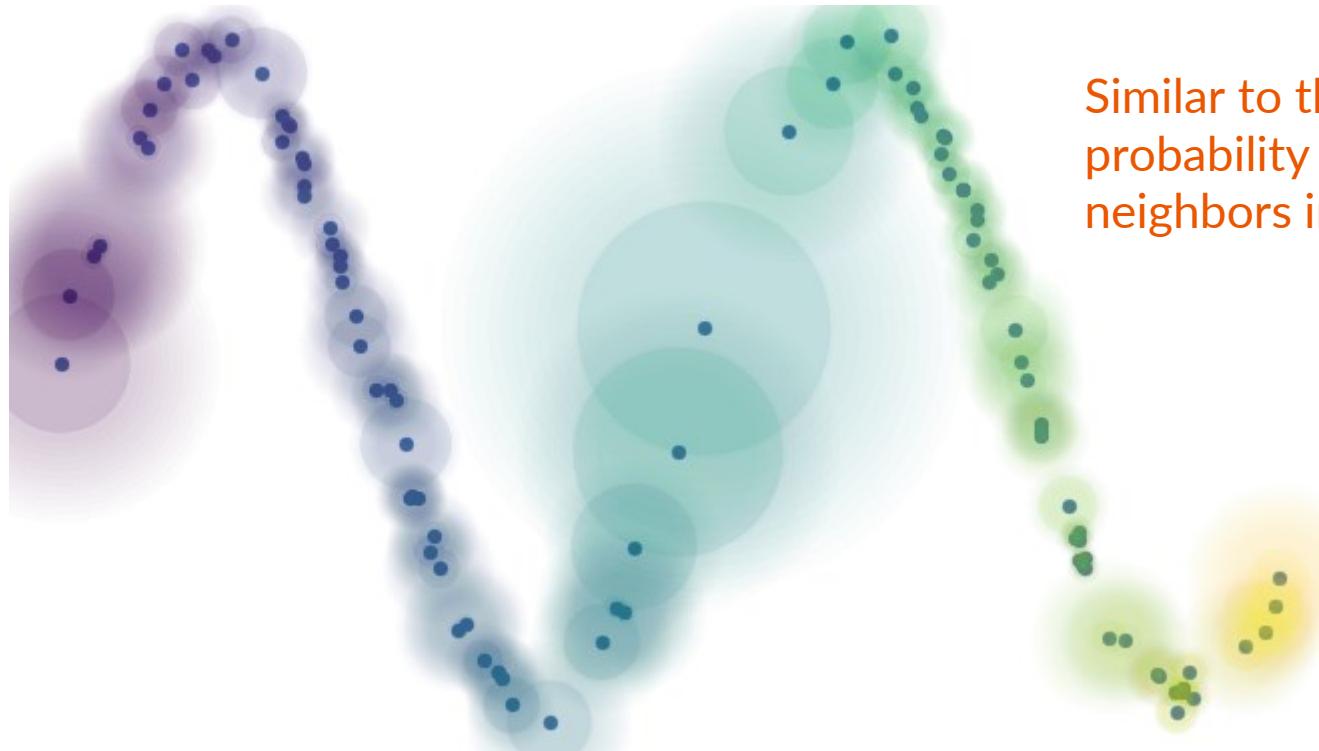
- Data came from **multiple manifolds**
- Data points were **sampled uniformly**

Uniform sampling = similar distance to k -th neighbor

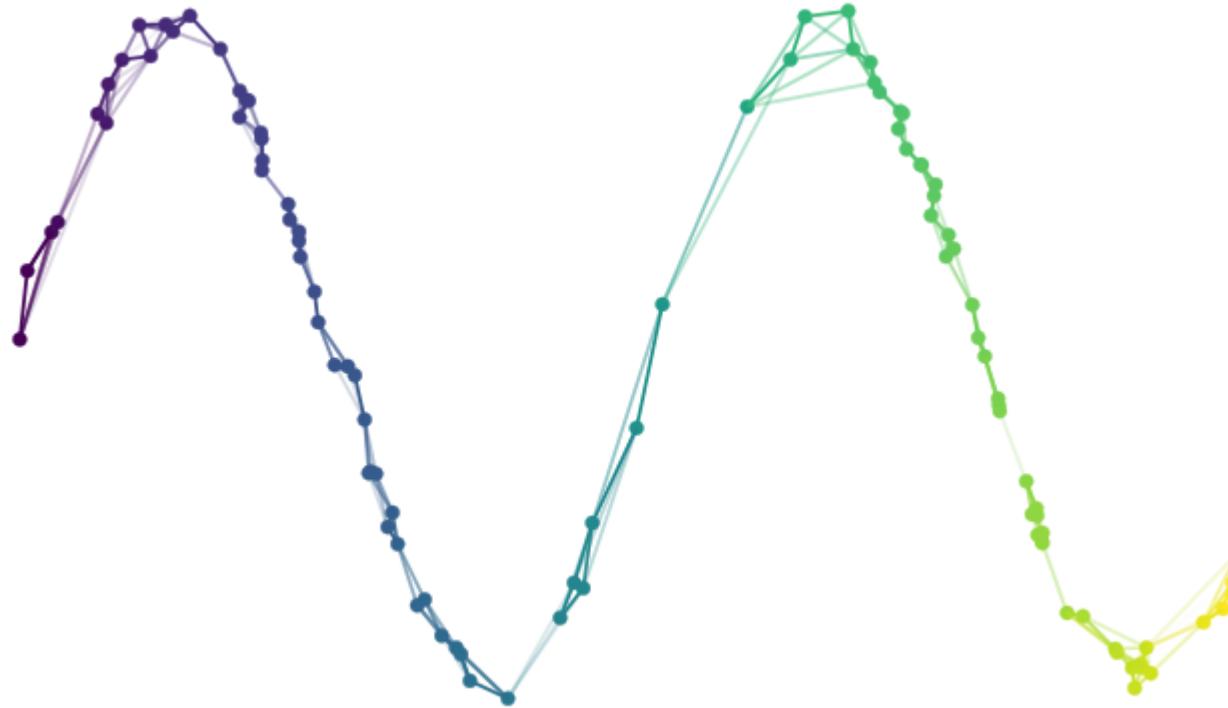


Similar effect as
perplexity in t-SNE

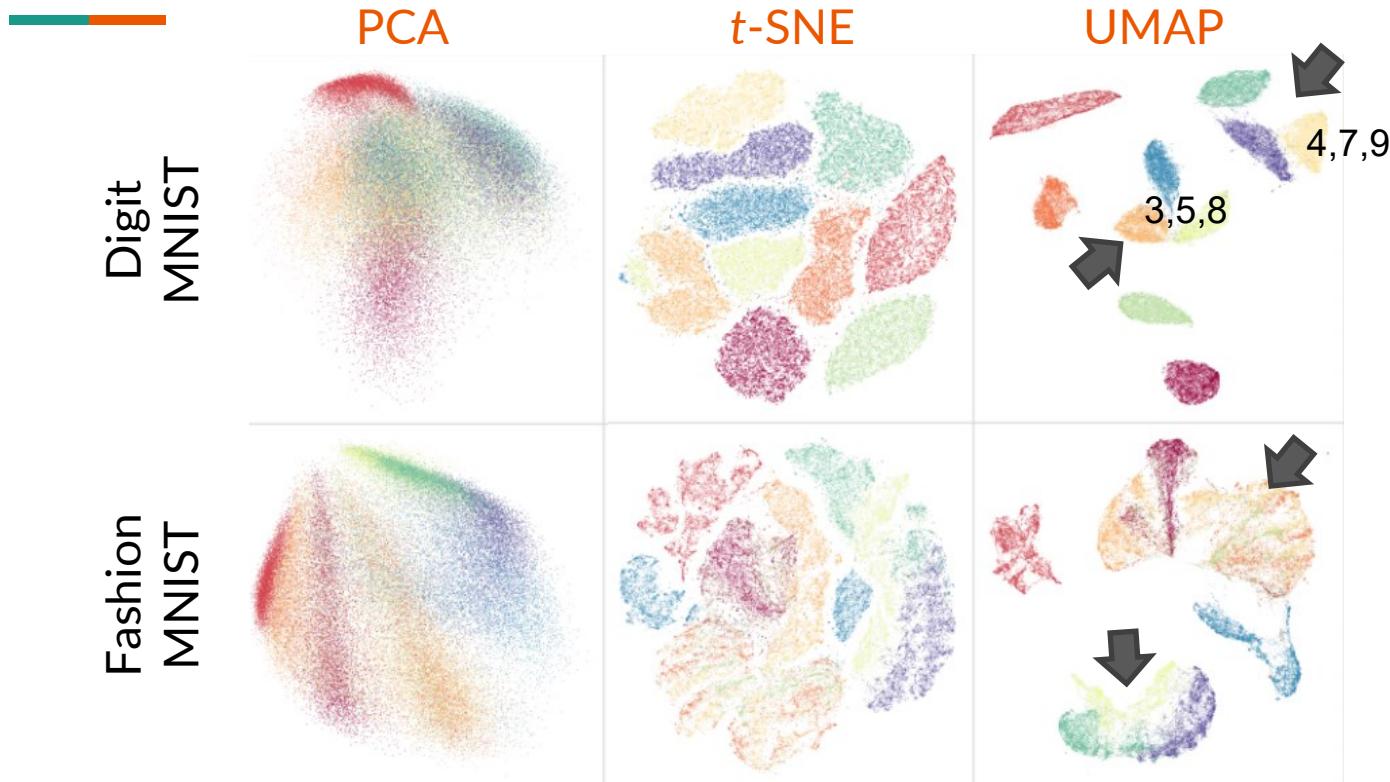
Adding uncertainty between faraway data points



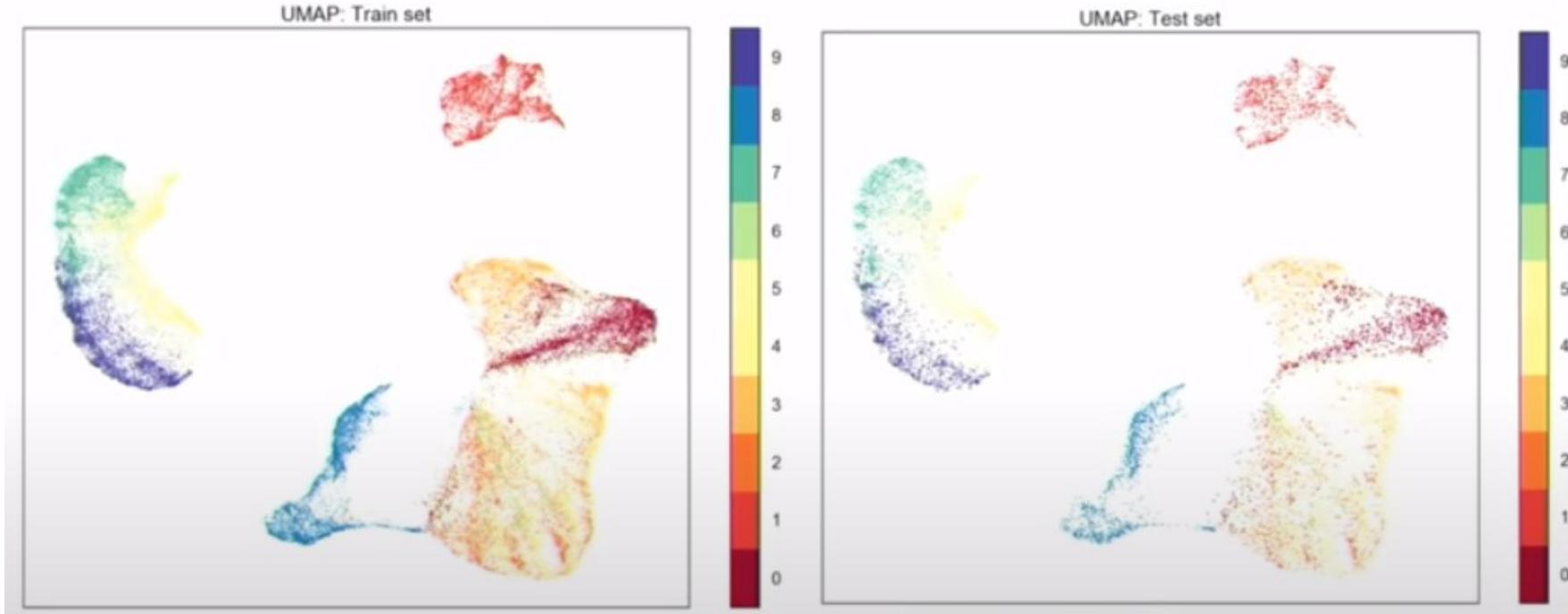
Network representation of neighbor relationship



UMAP can capture long-range relationship



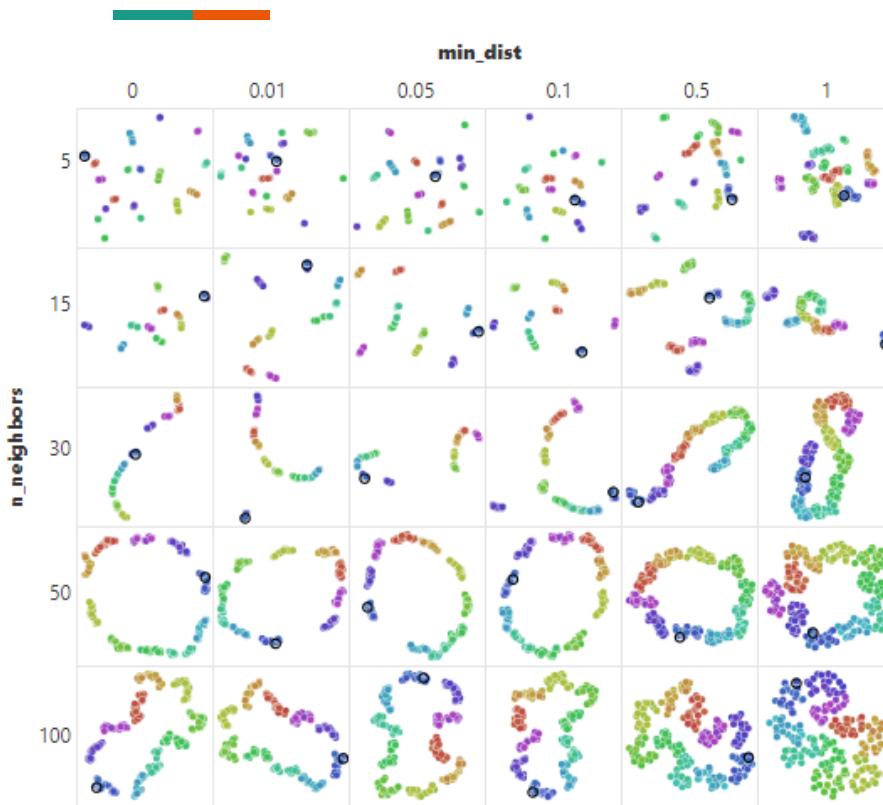
UMAP can transform new data points



Pros and cons of UMAP

- Can capture long-range relationship
- Can be applied to new data points without recomputing
- Require **strong assumptions**

Customizing UMAP outputs



- Number of neighbors (`n_neighbors`) is perplexity
- Minimum distant for placing similar data point (`min_dist`) is for adjusting the scale of visualization

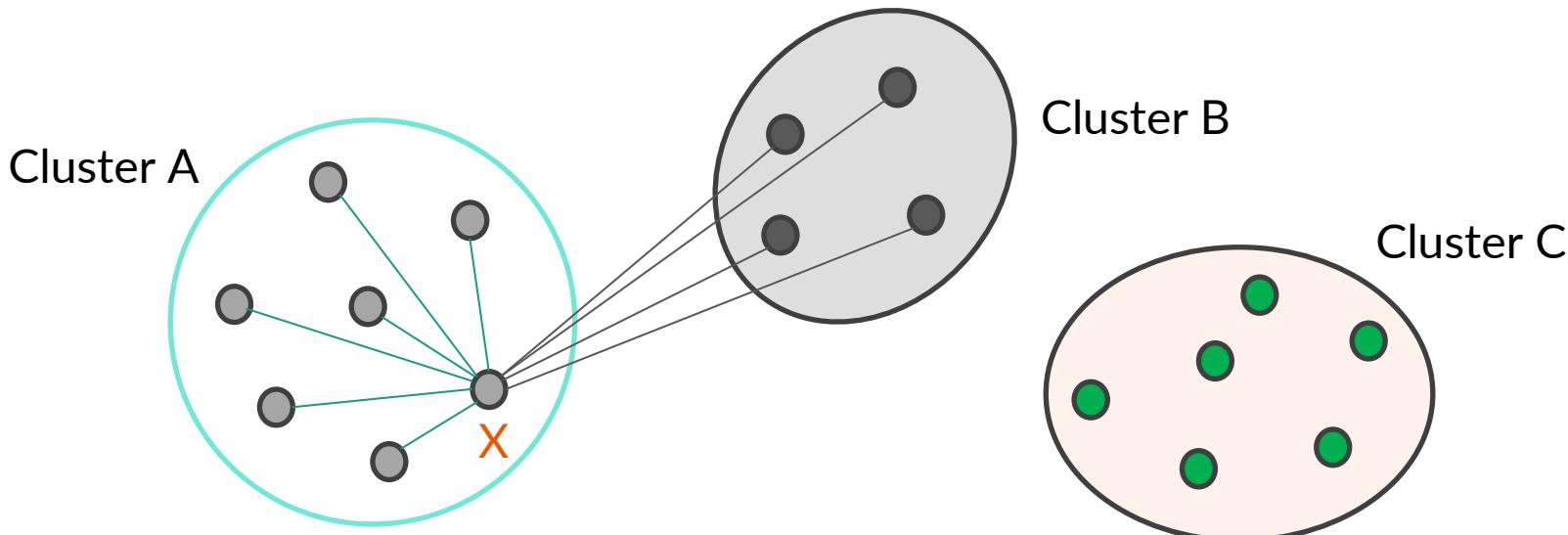


Clustering

The heart of clustering

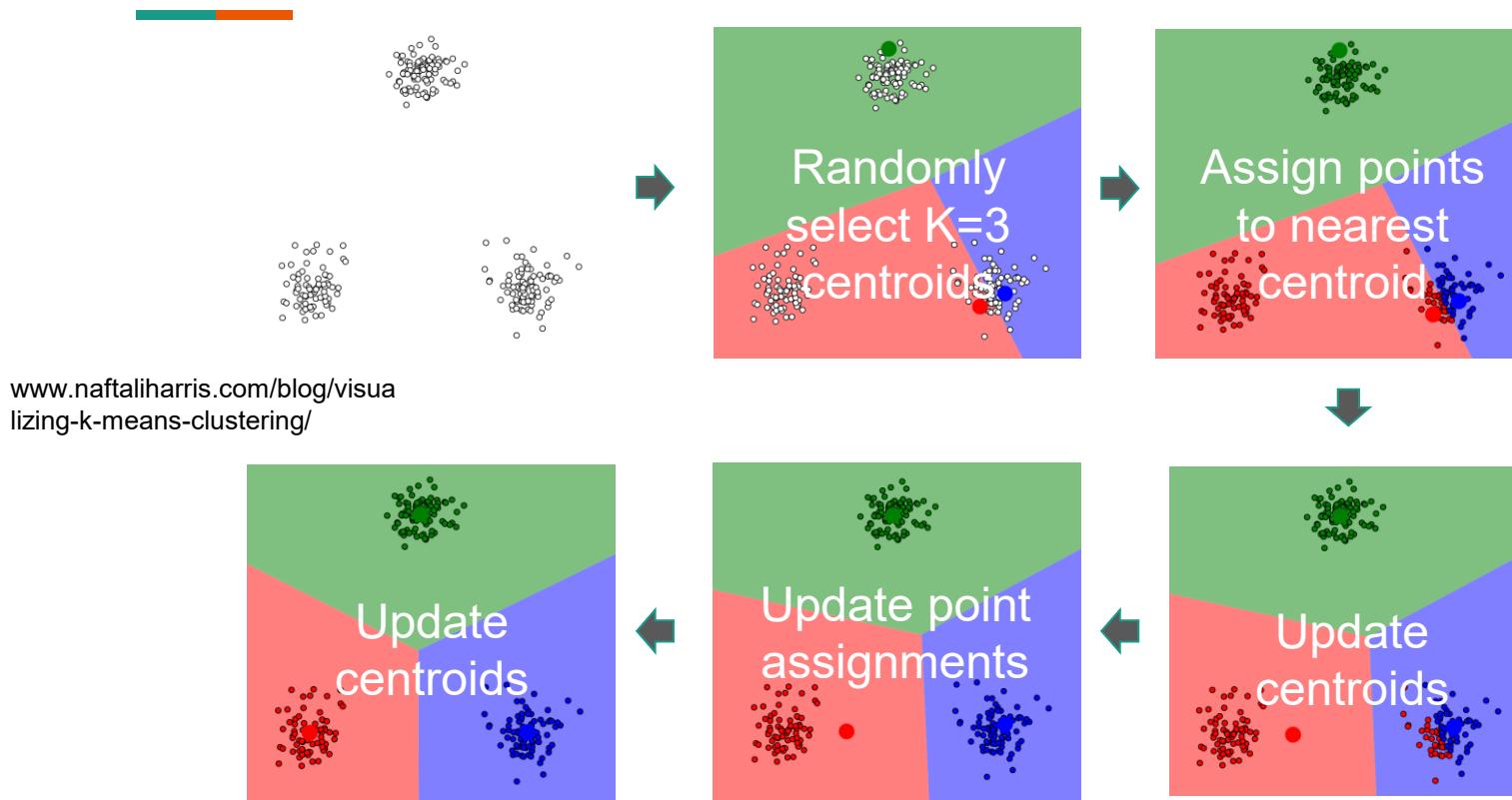
- Goal: Group **similar** data point together
- How to define **similarity**?
 - **Distance**: Between two data points
 - **Linkage**: Between groups of data points
- How many clusters is appropriate?
 - **Within-cluster (small) versus between-cluster (large) distance**

Silhouette score

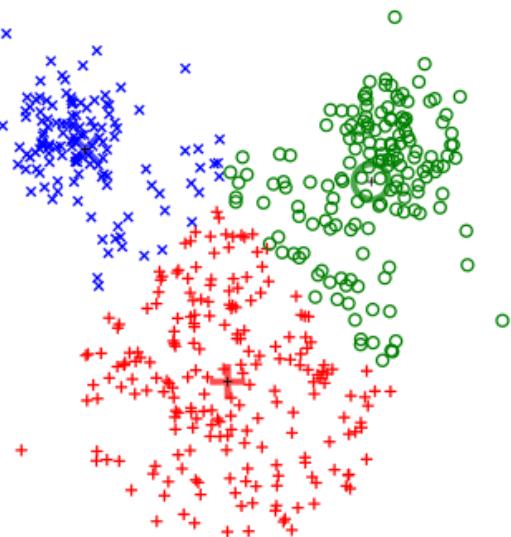


- Compare distances from **X** to other members of cluster A versus distances from **X** to members of cluster B (the closest cluster from A)

k-mean clustering

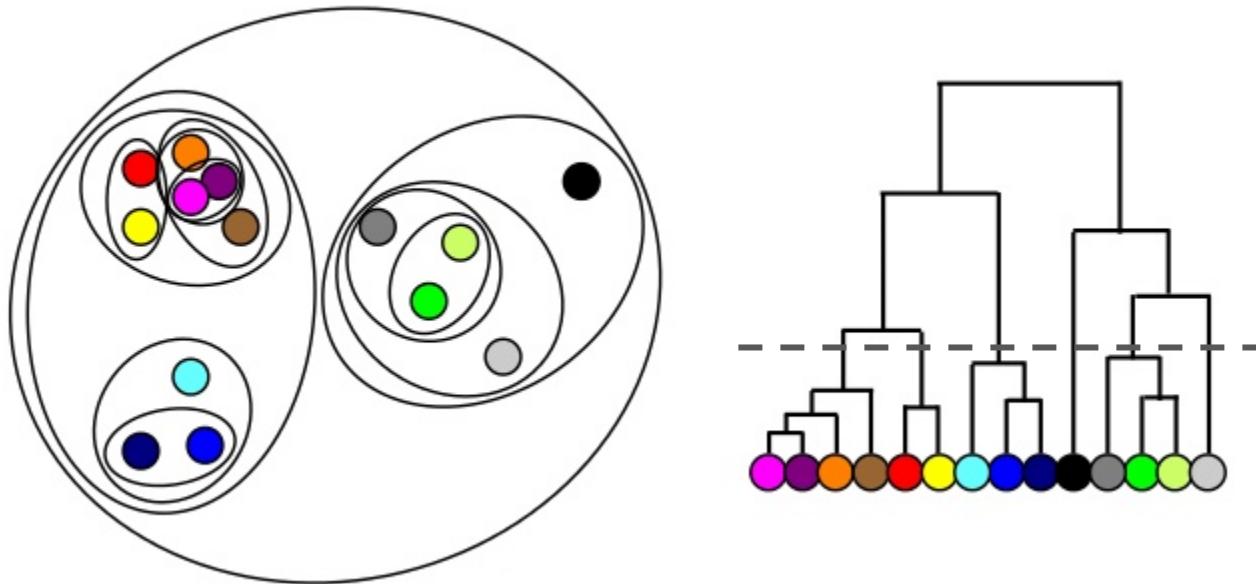


Limitation of k -mean



- Assume Euclidean distance
- Assume that clusters are of equal radius
- The initial guess of the locations of k means can affect the final clusters
 - Repeat multiple times

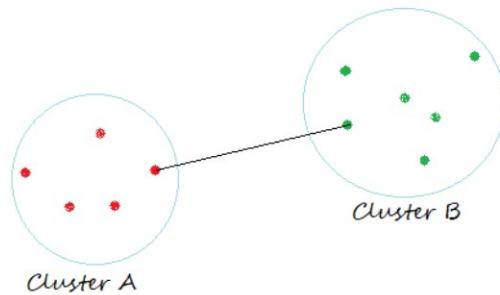
Agglomerative / Hierarchical clustering



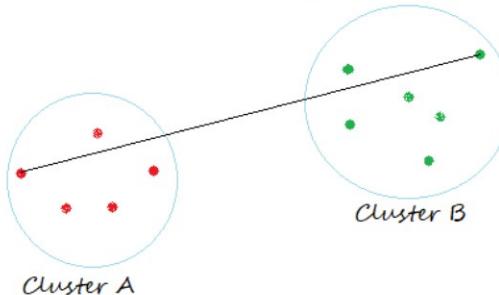
Source: www.slideshare.net/ElenaSgis/data-preprocessing-and-unsupervised-learning-methods-in-bioinformatics

Linkage = distance metric for groups of data points

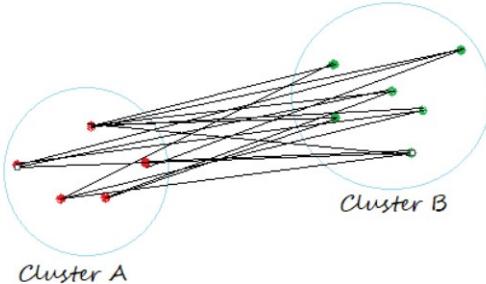
Single Linkage



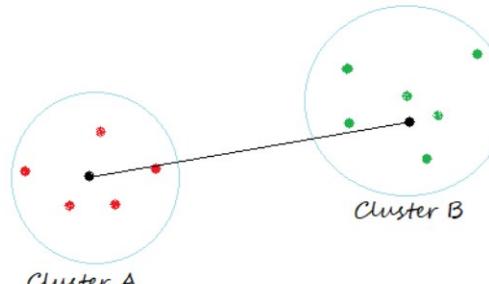
Complete Linkage



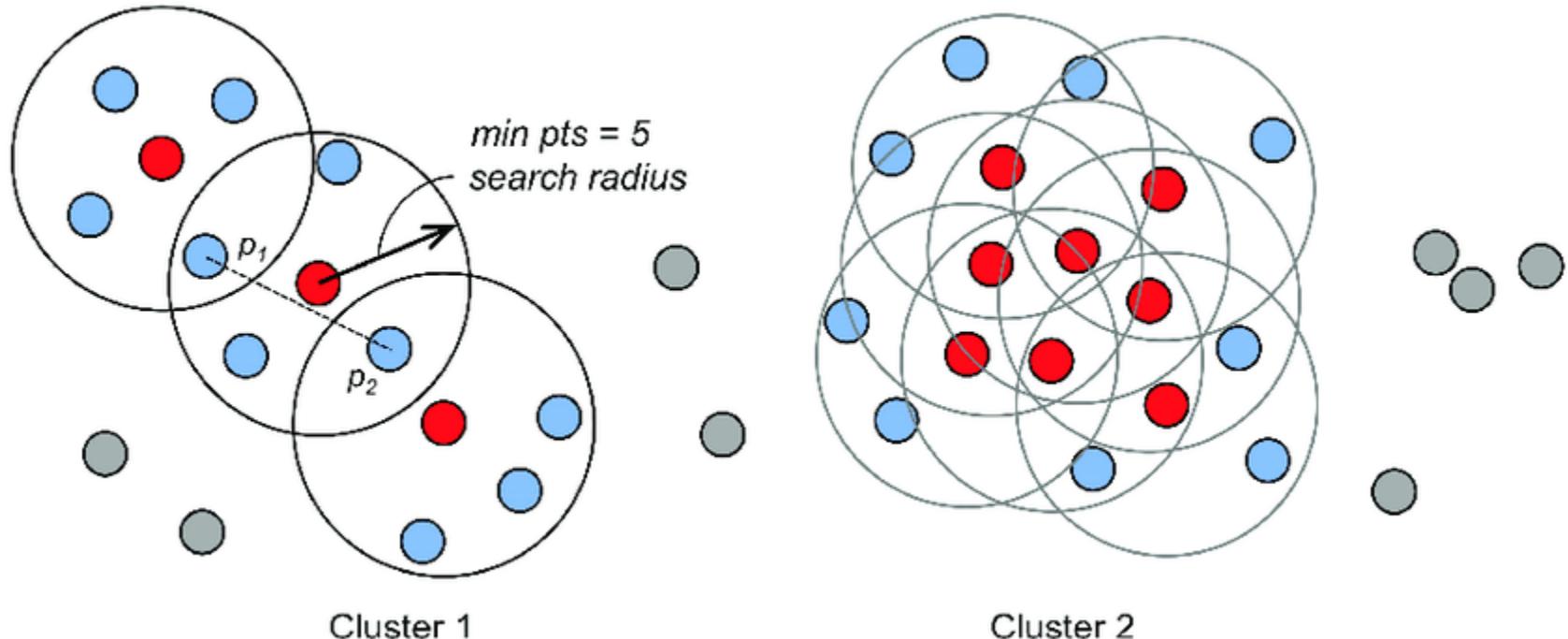
Average Linkage



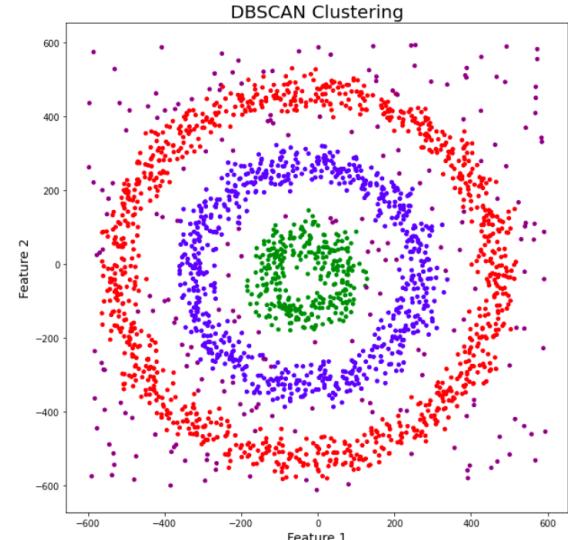
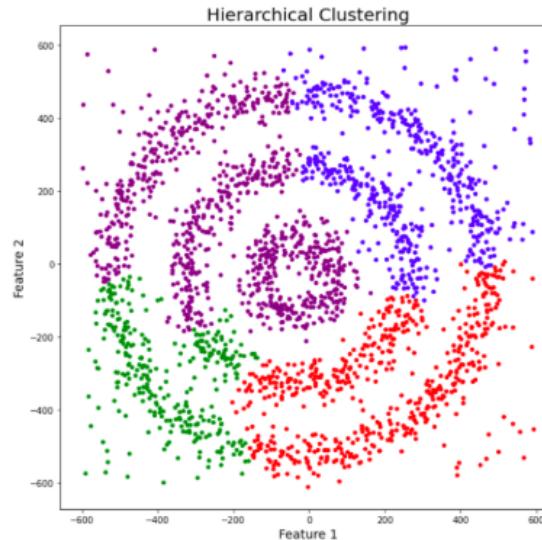
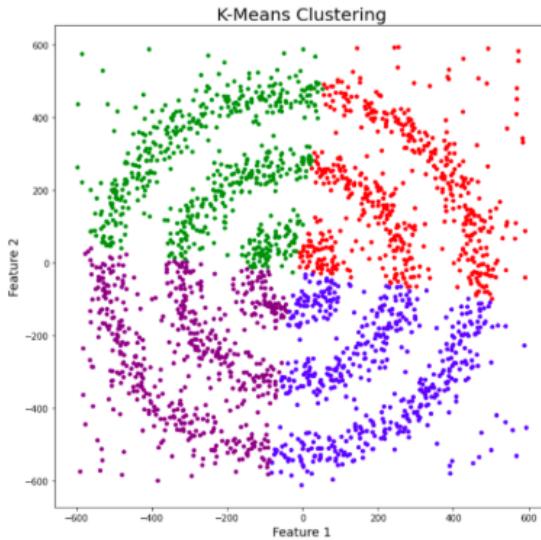
Centroid Linkage



DBSCAN: A density-based technique



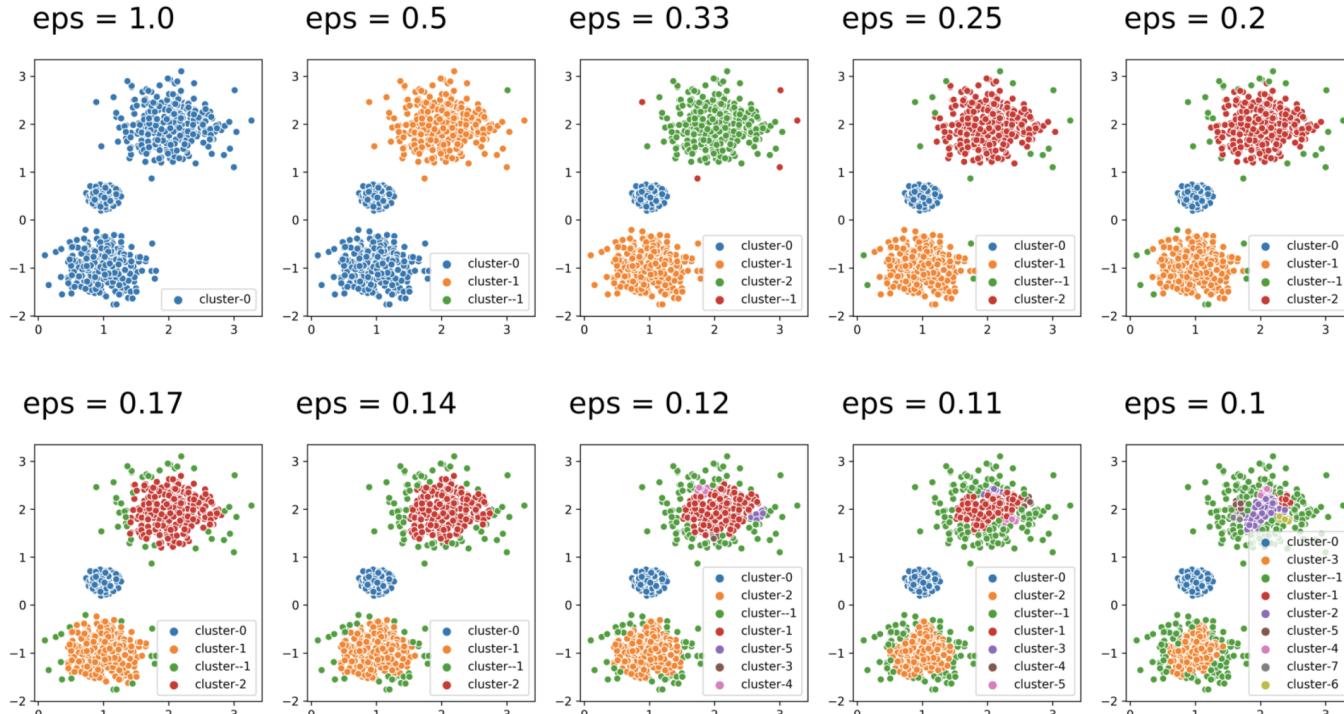
DBSCAN can handle complex cluster shape



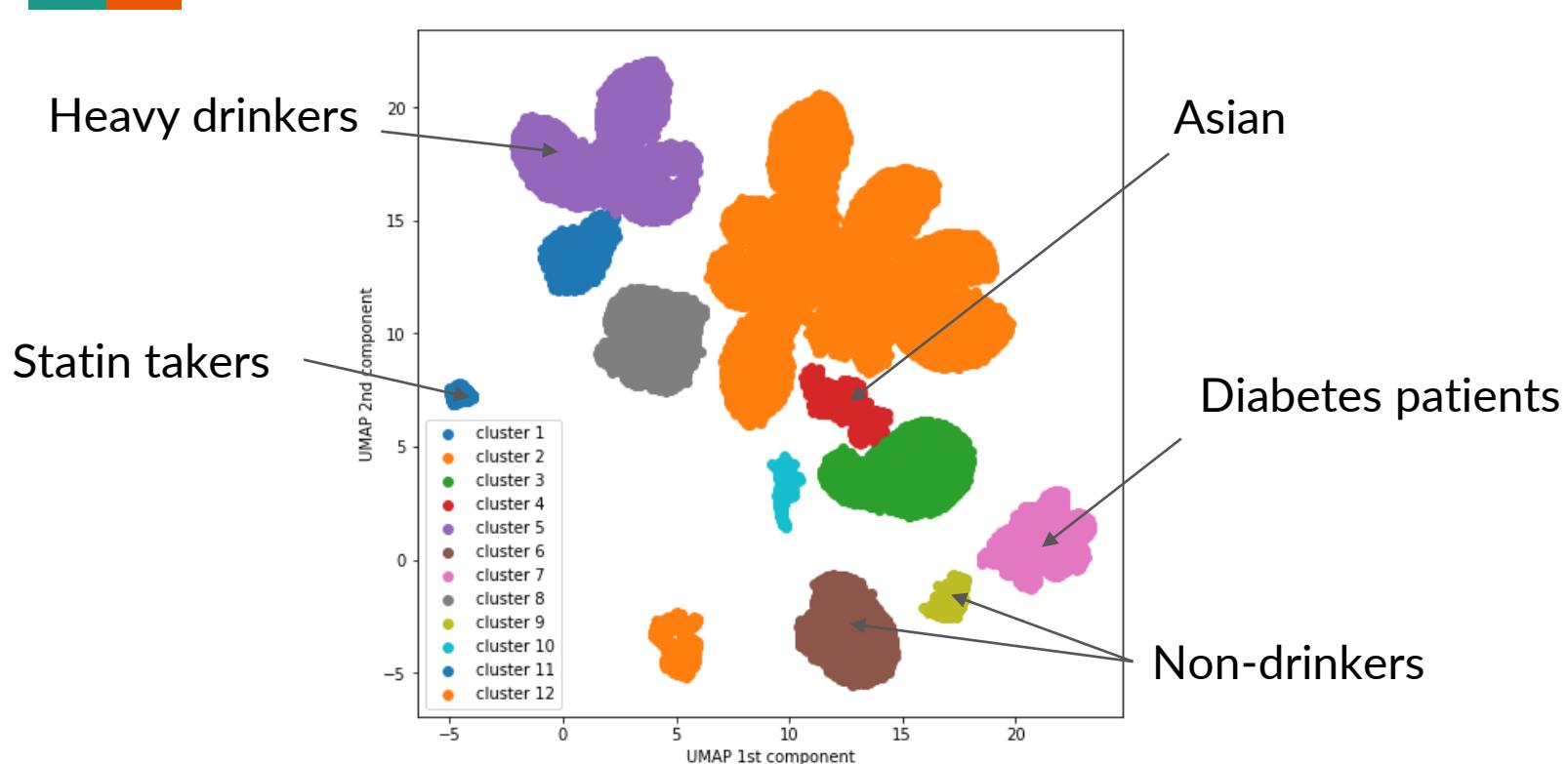
<https://www.analyticsvidhya.com/blog/2020/09/how-dbscan-clustering-works/>

- Distance-based techniques assume that data are spread in all directions

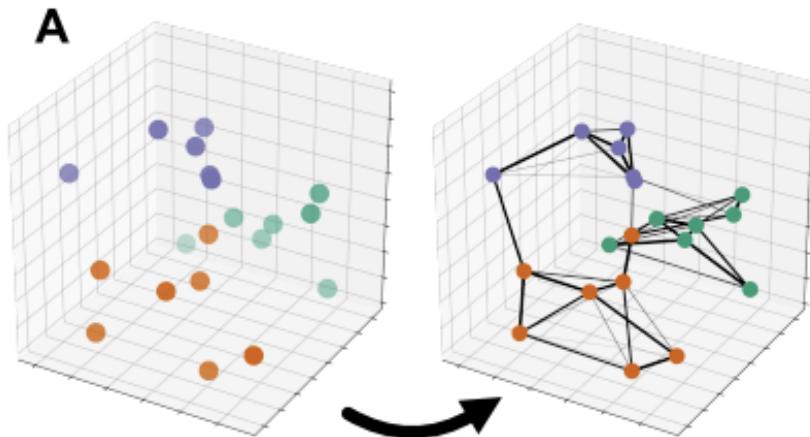
But may be difficult to tune



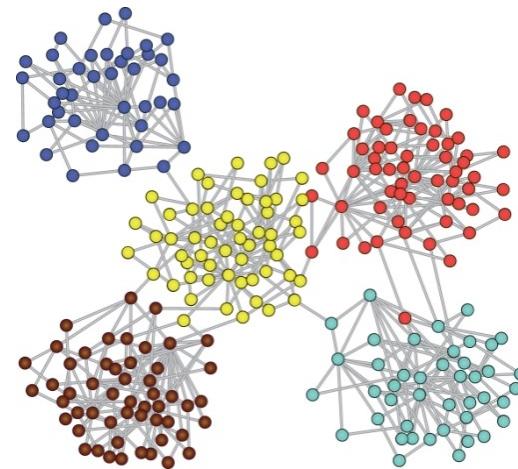
DBSCAN on patient data



Spectral clustering and network clustering



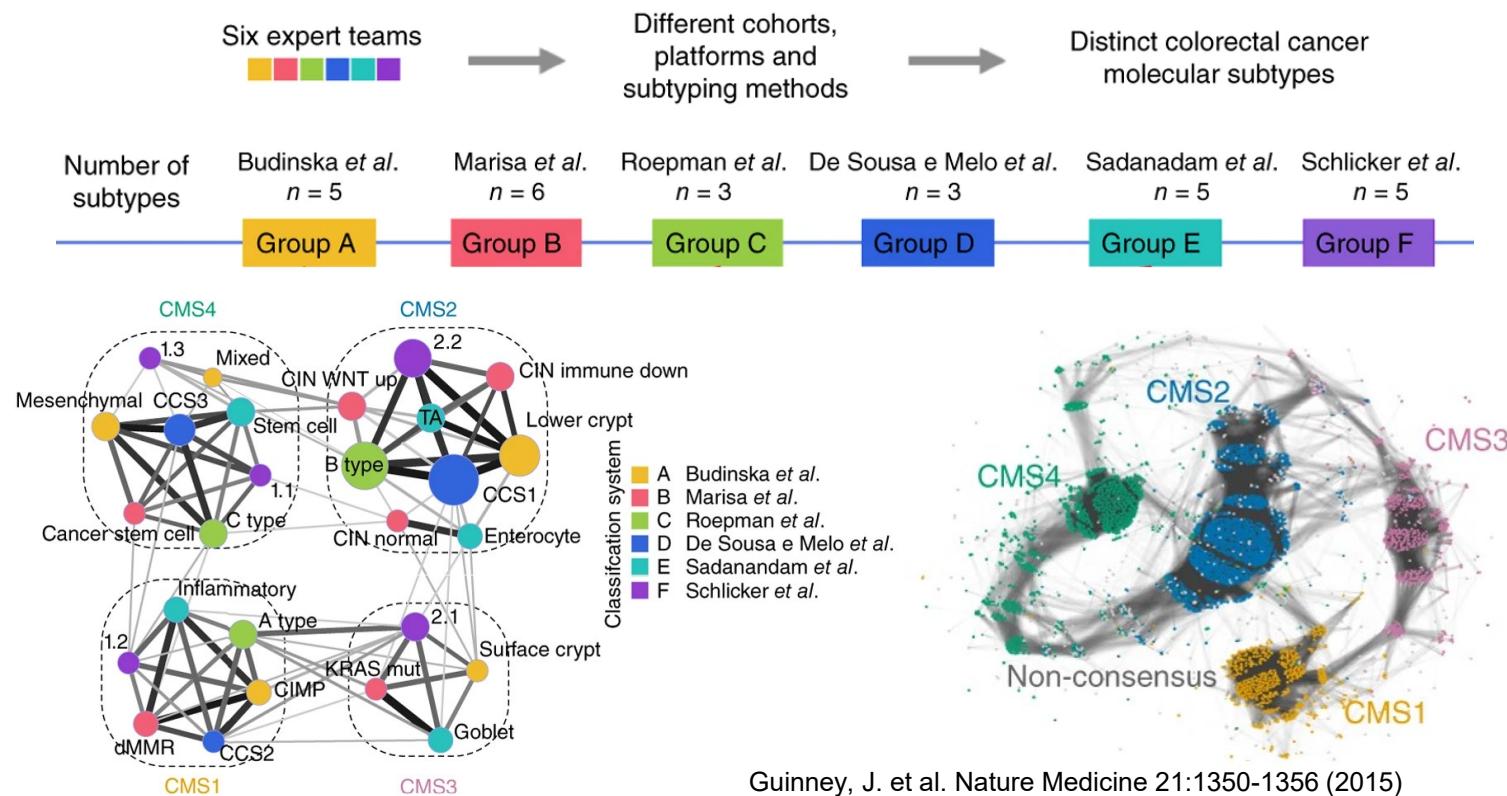
Sainburg, T. et al., Neural Comput 33(11):2881-2907 (2021)



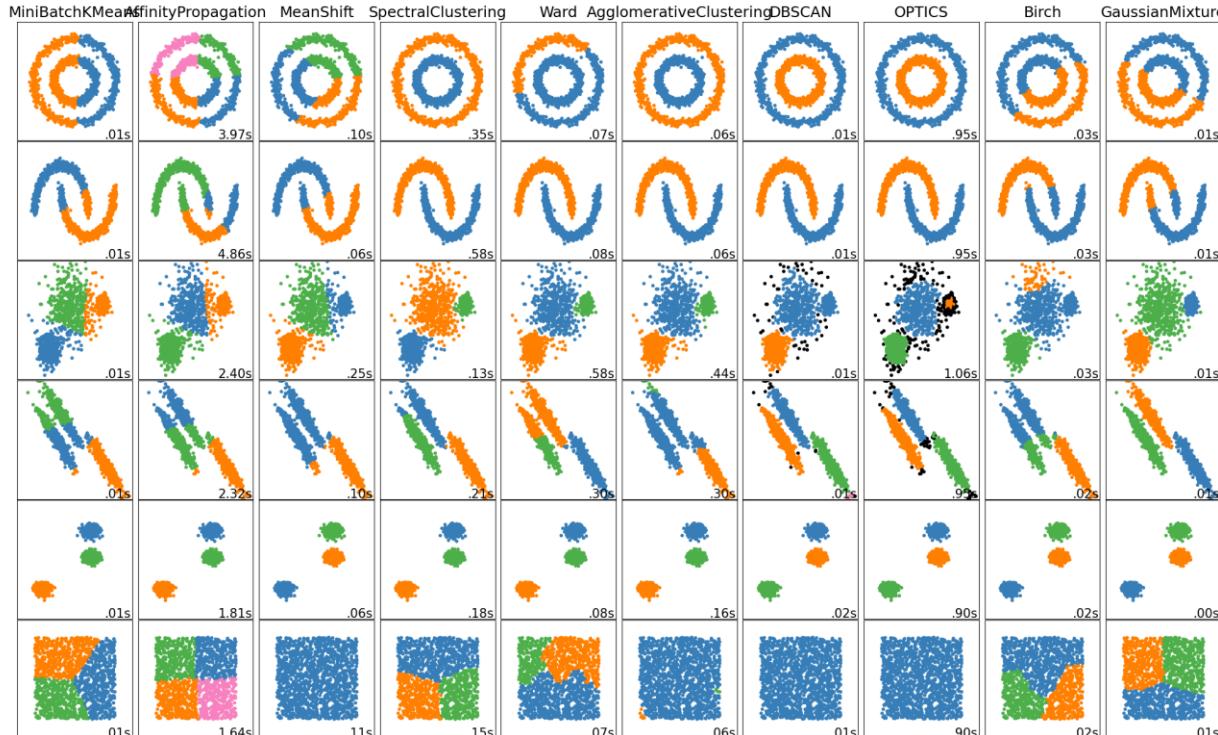
<https://github.com/topics/graph-clustering>

- View distance matrix as network
- Apply some threshold on the distance to create sparse network
- Split network into **modules with dense edges**

Colorectal cancer subtyping



No universal solution



Key points

- Explore distribution of feature values on the 2D/3D scatter plot
- Distance calculation influences everything!
 - Using *all* features vs *informative* features
 - Euclidean(x, y) = $\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$
 - Pick the appropriate distance metric for the data type
 - Correlation vs Euclidean for gene expression
- Unsupervised learning needs domain knowledge

Any question?

- See you on November 23