

Introduction to Machine Learning

Matthew Nguyen & Michael Schatz
September 25, 2024

Assignment 3 Due Monday September 30

The screenshot shows a web browser window displaying a GitHub README file for Assignment 3. The URL is github.com/schatzlab/appliedgenomics2024/blob/main/assignments/assignment3/README.md. The browser interface includes a sidebar with a 'Files' section containing various assignment files like main, assignments, and README.md.

Assignment 3: BWT and Variant Calling

Assignment Date: Monday, September 23, 2024
Due Date: Monday, September 30, 2023 @ 11:59pm

Assignment Overview

In this assignment you will implement the BWT and explore the requirements for variant calling. The programming exercises can be computed in any programming language, although we recommend python (or C++, Java, or Rust). R is generally inefficient at string processing unless you take great care. See the resources at the bottom of the page for tips for the variant calling exercises.

As a reminder, any questions about the assignment should be posted to [Piazza](#).

Question 1. BWT Encoding [20 pts]

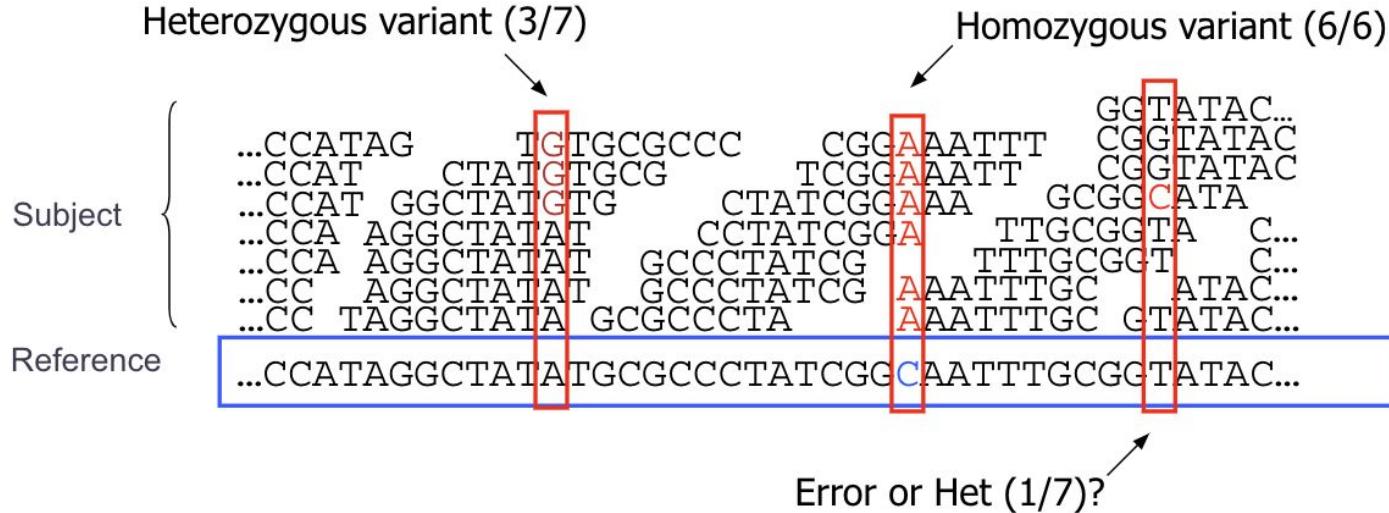
In the language of your choice, implement a BWT encoder and encode the string below. Faster (Linear time) methods exist for computing the BWT, although for this assignment you can use the simple method based on standard sorting techniques. Your solution does *not* need to be an optimal algorithm and can use $O(n^2)$ space and $O(n^2 \lg n)$ time.

Here is the recommended pseudo code (make sure to submit your code as well as the encoded string):

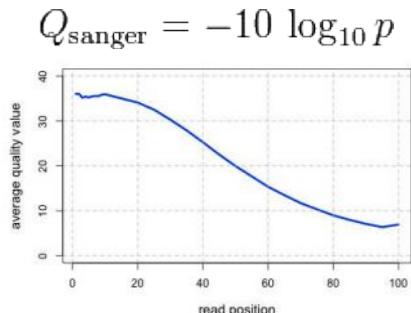
```
computeBWT(string s)
    ## add the magic end-of-string character
    s = s + "$"

    ## build up the BWT from the cyclic permutations
    ## note the ith cyclic permutation is just "s[i..n] + s[0..i]"
    rows = []
    for (i = 0; i < length(s); i++)
        rows.append(cyclic_permutation(s, i))
```

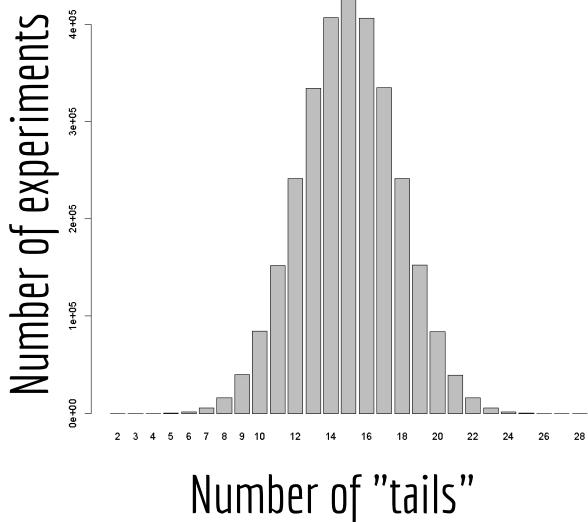
Genotyping Theory



- If there were no sequencing errors, identifying SNPs would be very easy: any time a read disagrees with the reference, it must be a variant!
- Sequencing instruments make mistakes
 - Quality of read decreases over the read length
- A single read differing from the reference is probably just an error, but it becomes more likely to be real as we see it multiple times



So, with 30 tosses (reads), we are much more likely to see an even mix of alternate and reference alleles at a heterozygous locus in a genome



This is why at least a "30X" (30 fold sequence coverage) genome is recommended: it confers sufficient power to distinguish heterozygous alleles and from mere sequencing errors

$$P(3/30 \text{ het}) <?> P(3/30 \text{ err})$$

PolyBayes: The first statistically rigorous variant detection tool.

letter

© 1999 Nature America Inc. • <http://genetics.nature.com>

A general approach to single-nucleotide polymorphism discovery

Gabor T. Marth¹, Ian Korff¹, Mark D. Yandell¹, Raymond T. Yeh¹, Zhijie Gu², Hamideh Zakeri², Nathan O. Stitzel¹, LaDeana Hillier¹, Pui-Yan Kwok² & Warren R. Gish¹

Bayesian posterior probability

$$P(\text{SNP}) = \sum_{\text{all variable } S} \frac{\frac{P(S_1 | R_1)}{P_{\text{Prior}}(S_1)} \cdot \frac{P(S_N | R_N)}{P_{\text{Prior}}(S_N)} \cdot P_{\text{Prior}}(S_1, \dots, S_N)}{\sum_{S_{i_1} \in \{A,C,G,T\}} \dots \sum_{S_{i_N} \in \{A,C,G,T\}} \frac{P(S_{i_1} | R_1)}{P_{\text{Prior}}(S_{i_1})} \cdot \dots \cdot \frac{P(S_{i_N} | R_1)}{P_{\text{Prior}}(S_{i_N})} \cdot P_{\text{Prior}}(S_{i_1}, \dots, S_{i_N})}$$

Probability of observed base composition (should model sequencing error rate)

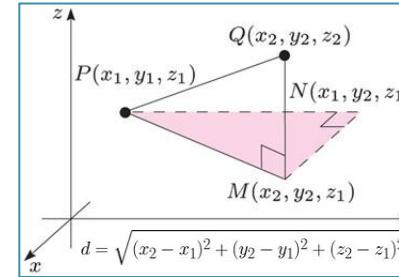
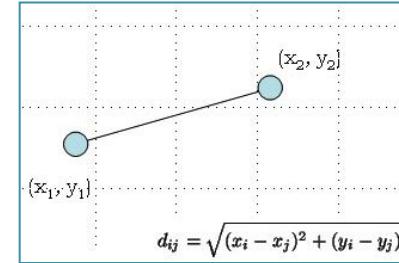
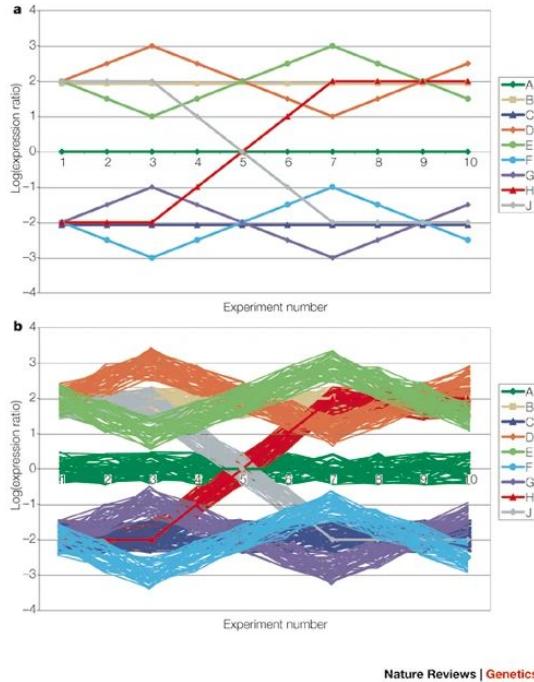
Base call + Base quality

Expected (prior) polymorphism rate

Introduction to Machine Learning: Clustering & Dimensionality Reduction

Clustering Refresher

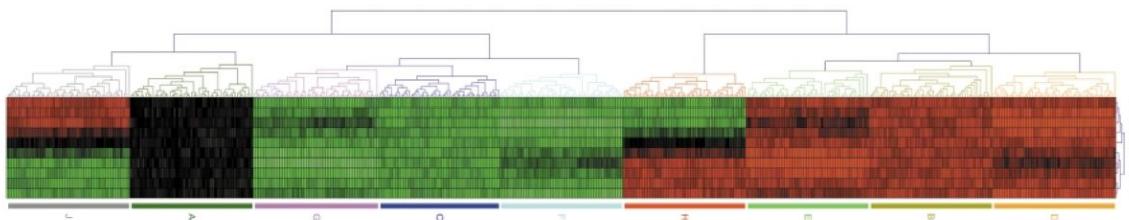
Euclidean Distance



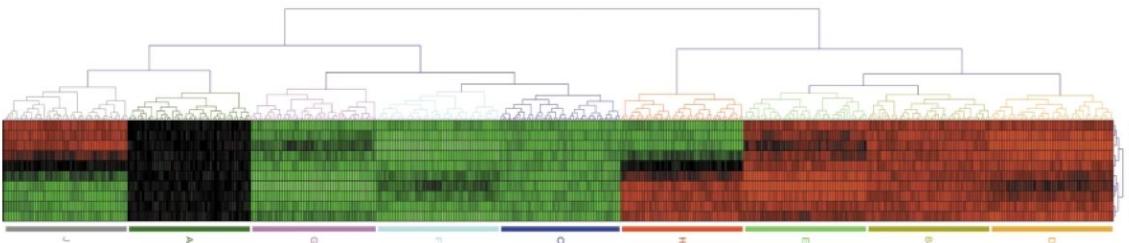
$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \cdots + (q_n - p_n)^2} = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}.$$

Computational genetics: Computational analysis of microarray data
Quackenbush (2001) *Nature Reviews Genetics*. doi:10.1038/35076576

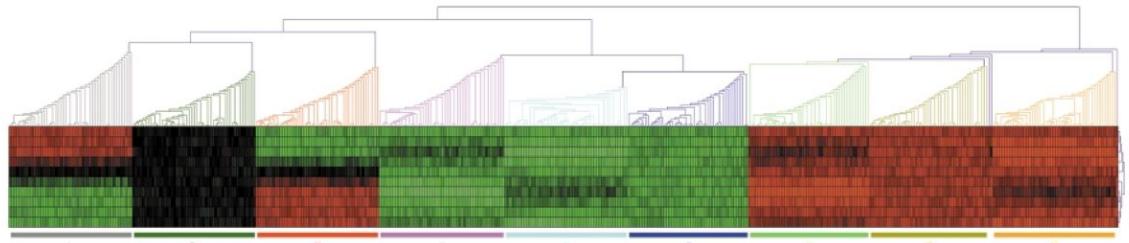
Hierarchical Clustering



average

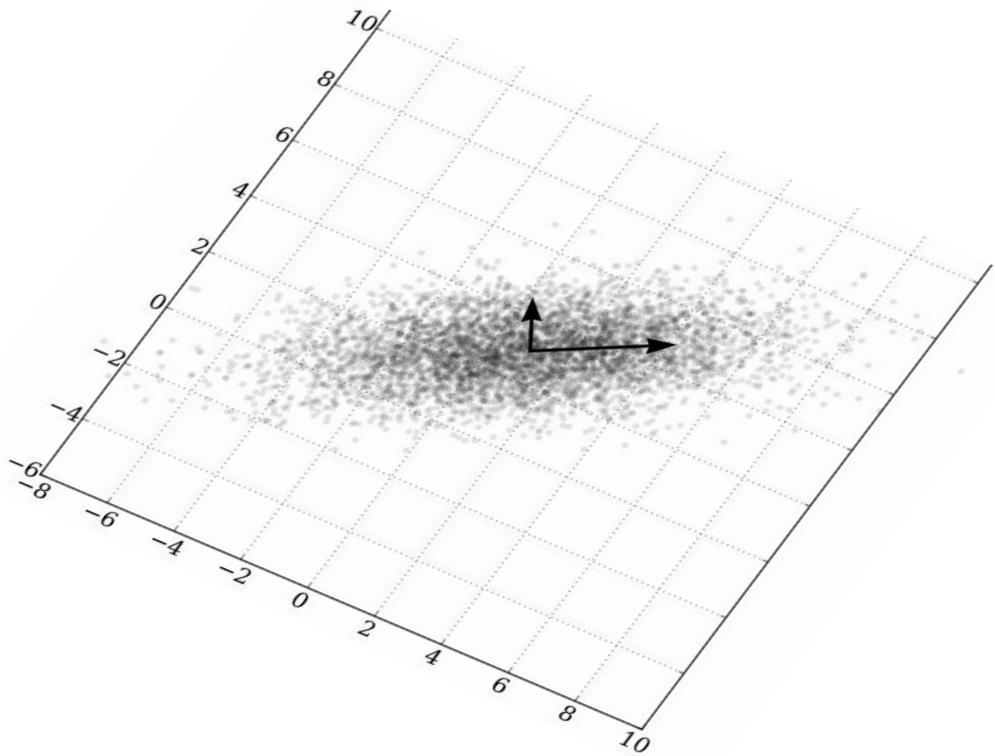
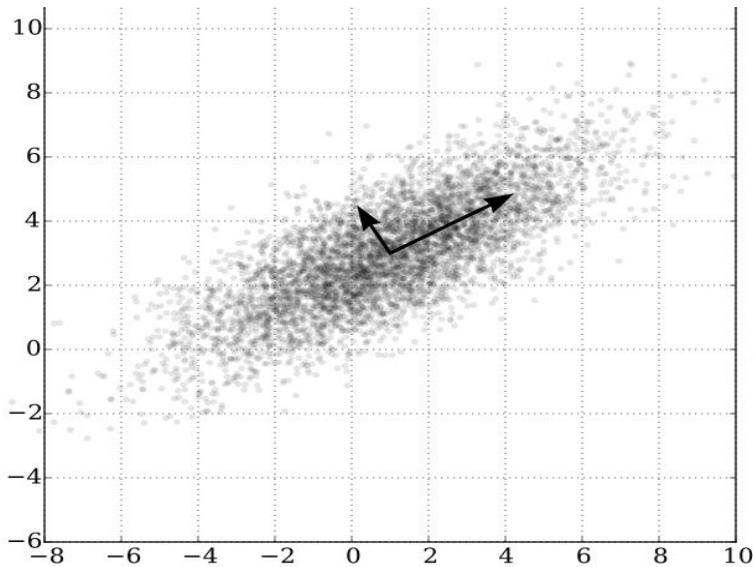


complete



single

Principal Components Analysis (PCA)



PC1: “New X”- The dimension with the most variability

PC2: “New Y”- The dimension with the second most variability

Principal Components Analysis (PCA)

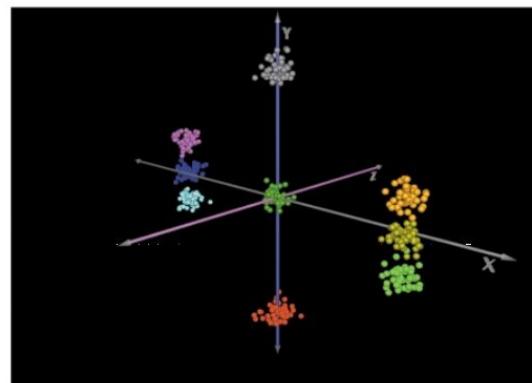
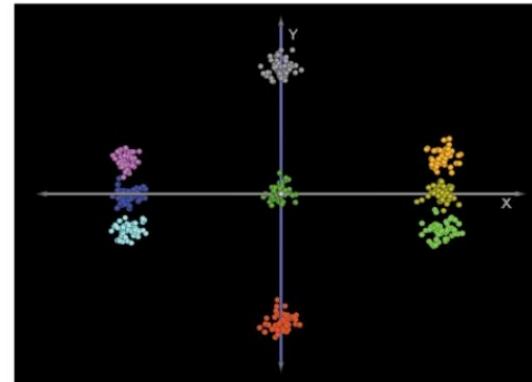
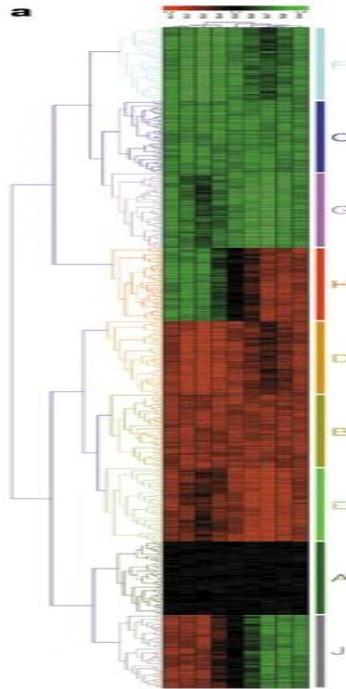
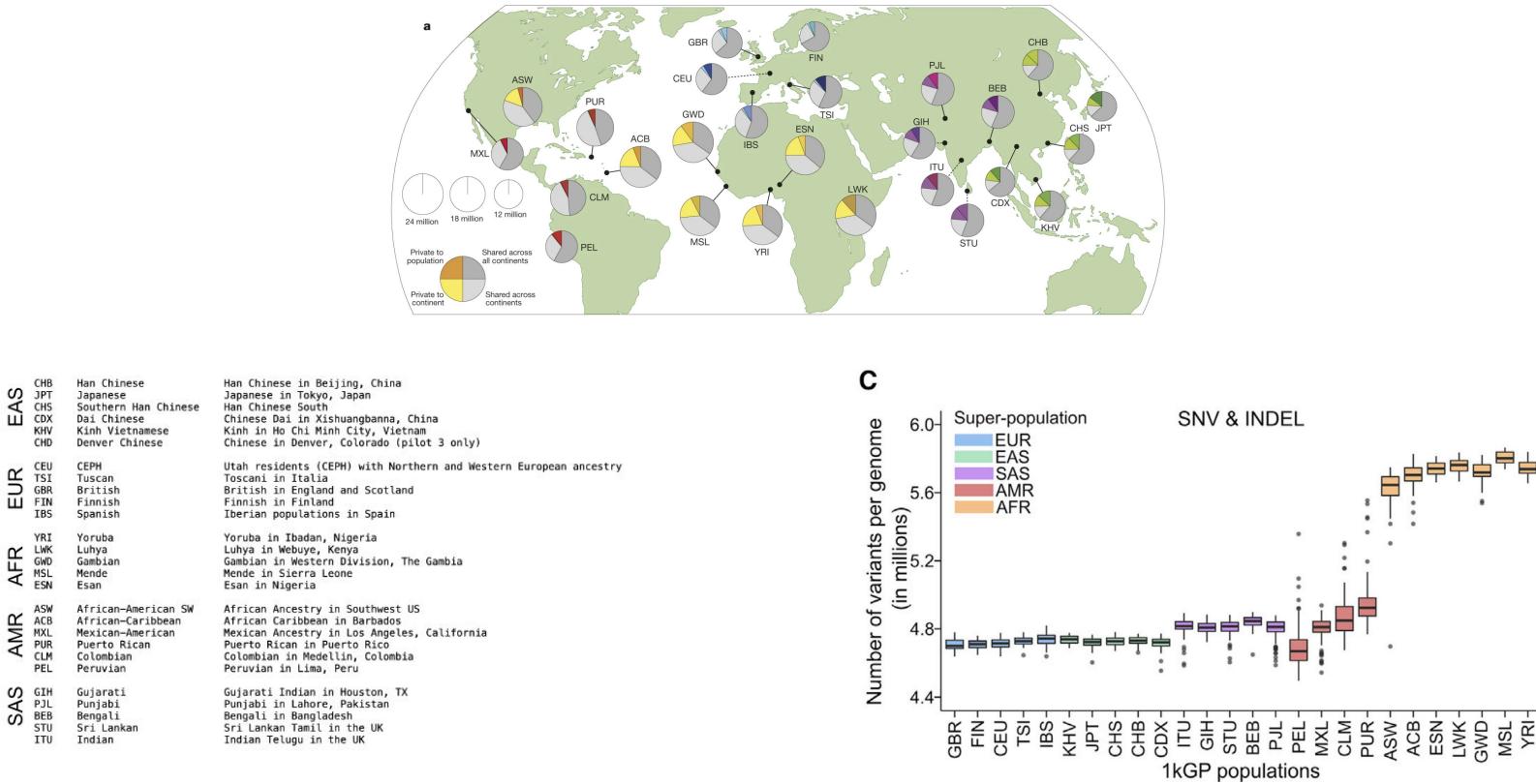


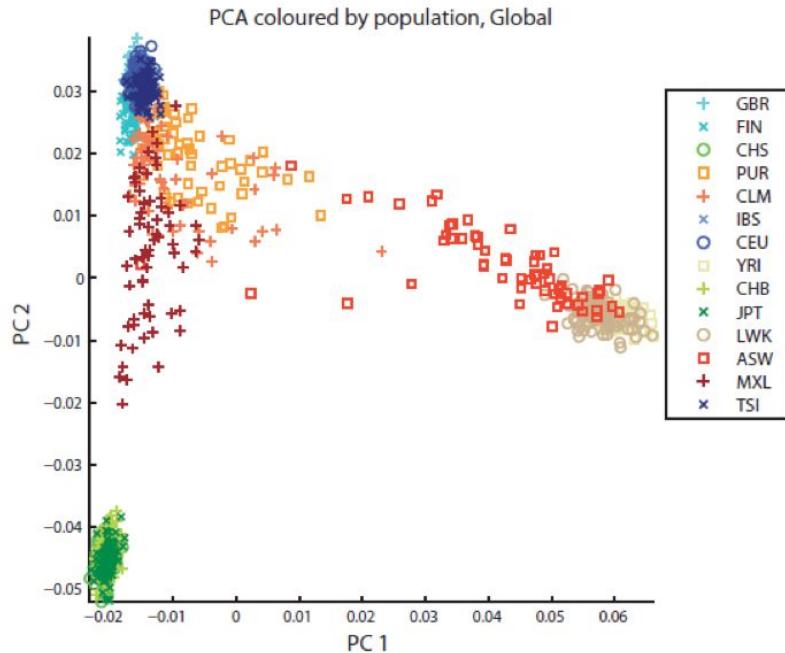
Figure 4 | **Principal component analysis.** The same demonstration data set was analysed using **a** | hierarchical (average-linkage) clustering and **b** | principal component analysis using Euclidean distance, to show how each treats the data, with genes colour coded on the basis of hierarchical clustering results for comparison.

1000 Genomes Project



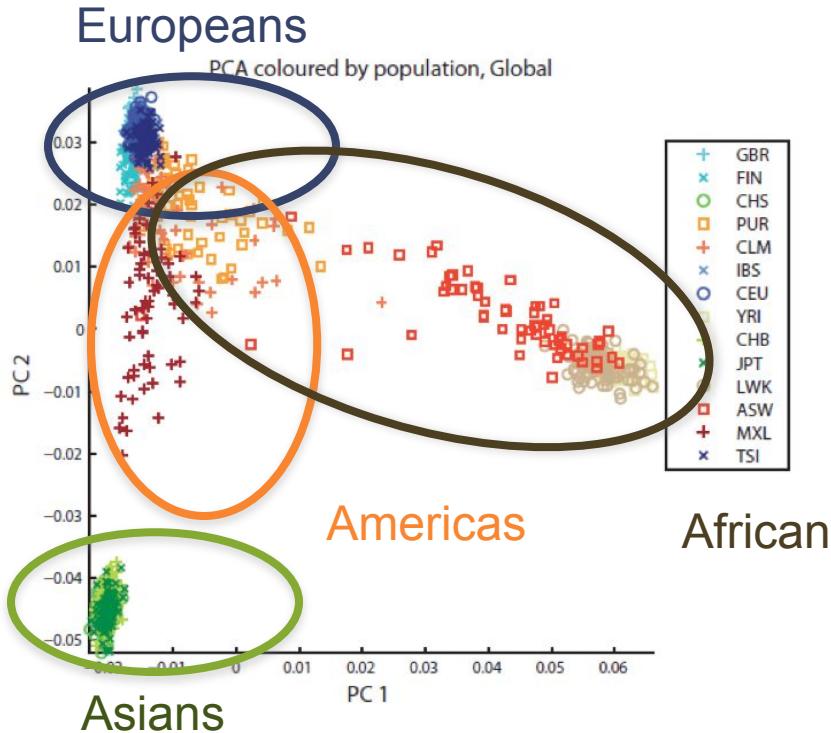
High-coverage whole-genome sequencing of the expanded 1000 Genomes Project cohort including 602 trios
 Byrska-Bishop et al. (2022) Cell. doi: 10.1016/j.cell.2022.08.004

Variation across populations



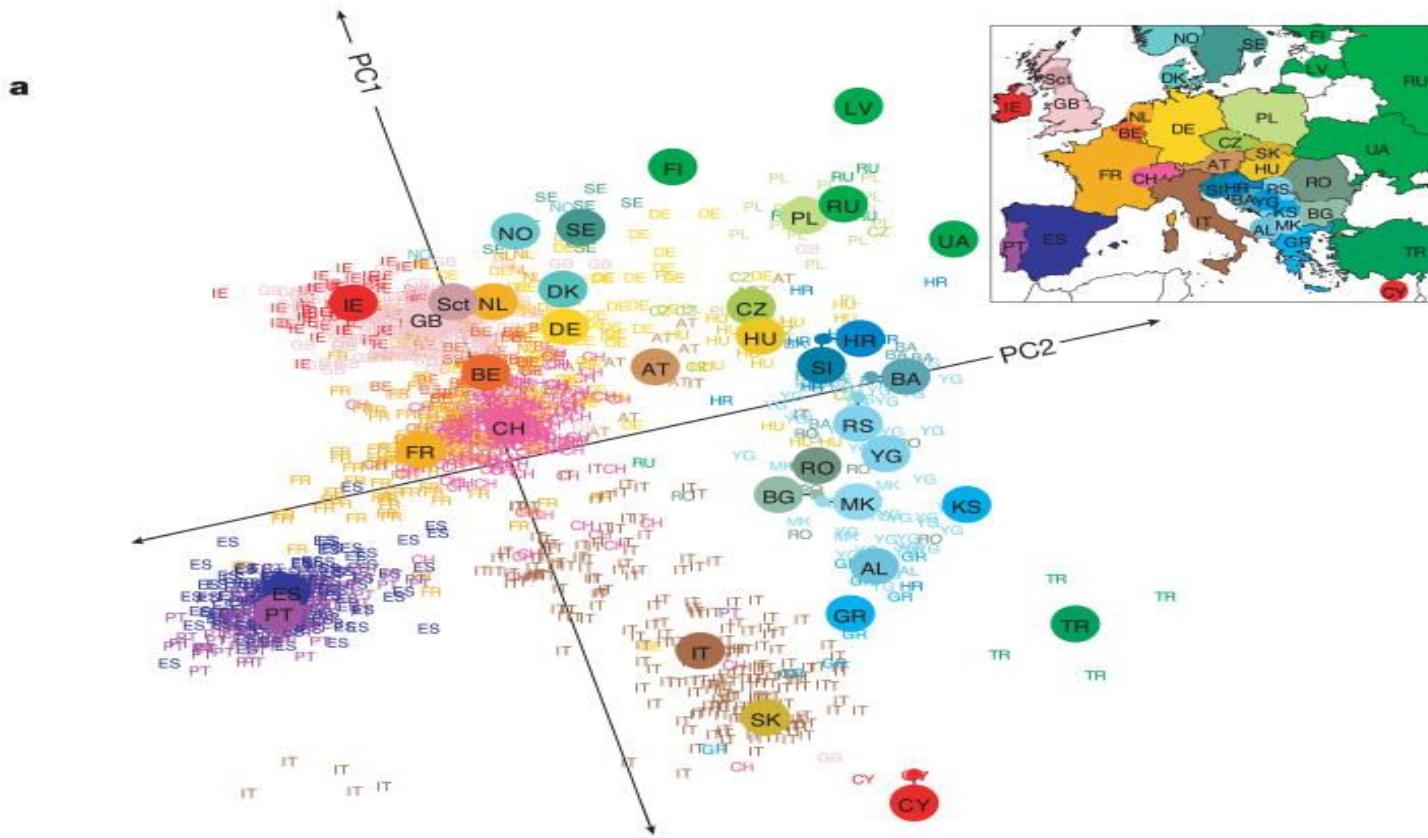
	Person1	Person2	Person3	...
SNP1	0/0	0/1	0/1	
SNP2	0/1	0/1	0/1	
SNP3	0/0	0/0	1/1	
SNP4	1/1	0/1	1/1	
SNP5	0/0	1/1	1/1	
SNP5	0/0	0/0	0/1	
...				

Variation across populations



	Person1	Person2	Person3	...
SNP1	0/0	0/1	0/1	
SNP2	0/1	0/1	0/1	
SNP3	0/0	0/0	1/1	
SNP4	1/1	0/1	1/1	
SNP5	0/0	1/1	1/1	
SNP5	0/0	0/0	0/1	
...				

PC1 is largely correlated with African ancestry:
Human genetic diversity is dominated by African ancestry



Genes mirror geography within Europe

Novembre et al (2008) Nature. doi: 10.1038/nature07331

MNIST



- 70,000 images, each labeled with the correct digit (0-9)
 - 28 x 28 thumbnail images of handwritten digits
 - Grayscale image with 256 possible values (784 bytes per image)
- Widely used dataset for image classification
 - Very simple, the categories are very easy to understand

MNIST: PCA



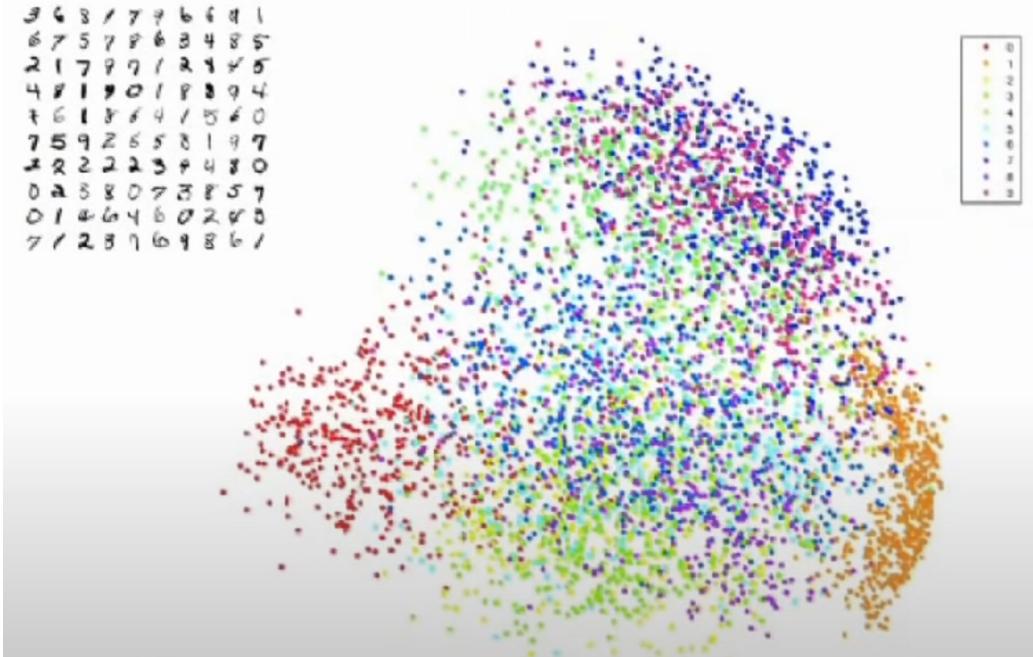
Principal Components Analysis

Each 28x28 image can be considered a point in 784 dimensional space!

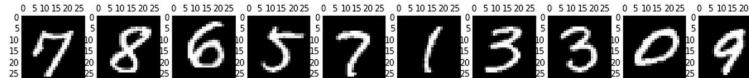
PCA roughly groups digits together

- 0 (red) in bottom left
- 1 (orange) in bottom right

But very incomplete separation, lots of mixing of colors

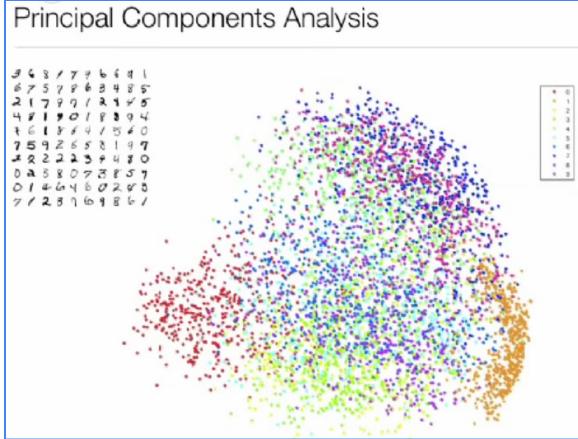


MNIST: PCA and t-SNE



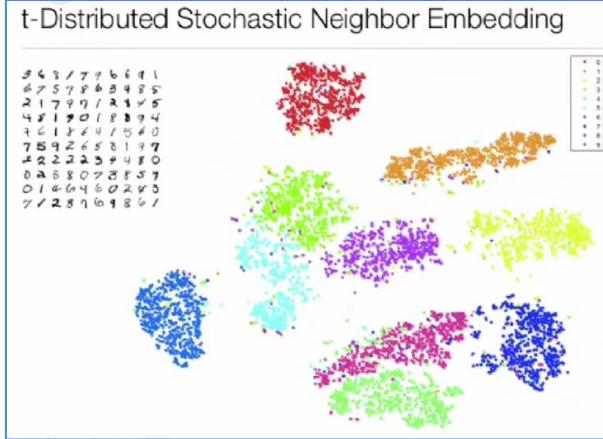
Principal Components Analysis

3 4 8 7 7 9 6 6 4 1
6 7 5 7 9 6 3 4 8 5
2 1 7 9 9 7 3 1 4 5
4 8 1 9 0 1 8 3 9 4
3 6 1 8 4 1 5 2 0
7 5 9 2 6 5 3 1 9 7
4 8 2 2 3 4 1 4 0
0 3 5 8 0 7 2 8 5 7
0 1 4 6 4 6 0 2 9 5
7 2 3 1 6 1 8 1 7



t-Distributed Stochastic Neighbor Embedding

3 4 3 1 7 7 6 6 4 1
6 7 5 7 8 6 5 4 8 5
2 1 7 9 9 7 3 1 4 5
4 8 1 9 0 1 8 3 9 4
3 6 1 8 4 1 5 2 0
7 5 9 2 6 5 3 1 9 7
4 8 2 2 3 4 1 4 0
0 3 5 8 0 7 2 8 5 7
0 1 4 6 4 6 0 2 9 5
7 2 3 1 6 1 8 1 7



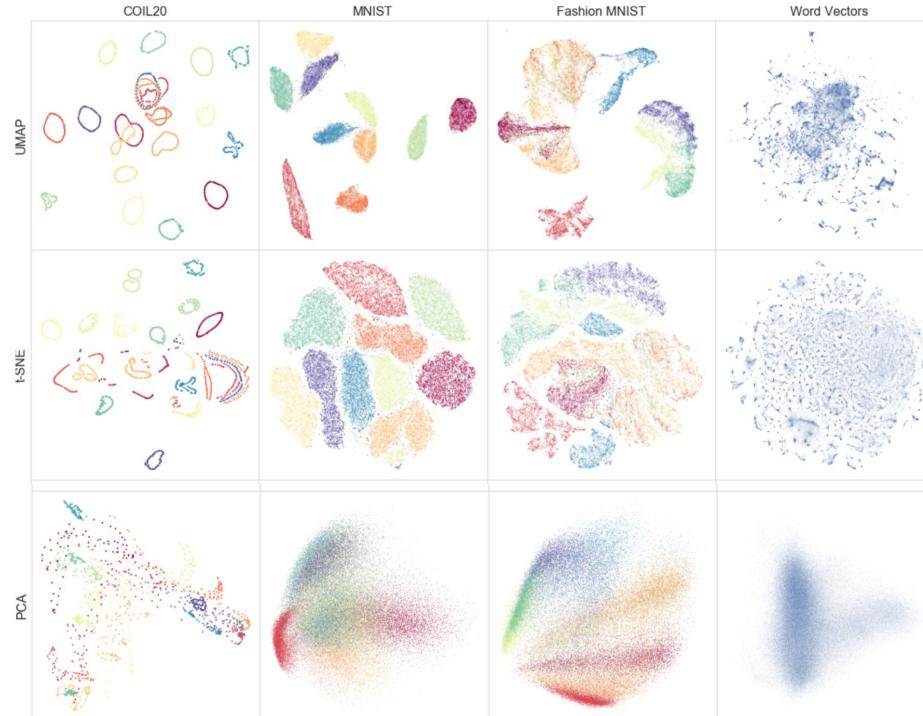
t-distributed Stochastic Neighborhood Embedding

- Non-linear dimensionality reduction technique: distances are only locally meaningful
- Rather than Euclidean distances, for each point fits a Gaussian kernel to fit the nearest N neighbors (perplexity) that define the probabilities that two points should be close together
- Using an iterative spring embedding system to place high probability points nearby

Visualizing Data Using t-SNE

<https://www.youtube.com/watch?v=RJVL80Gg3IA>

UMAP



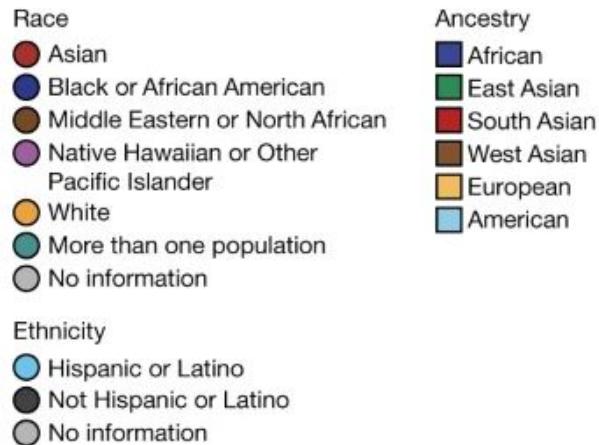
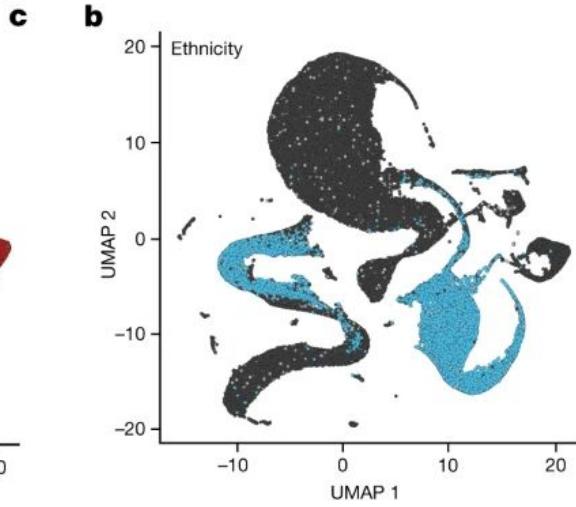
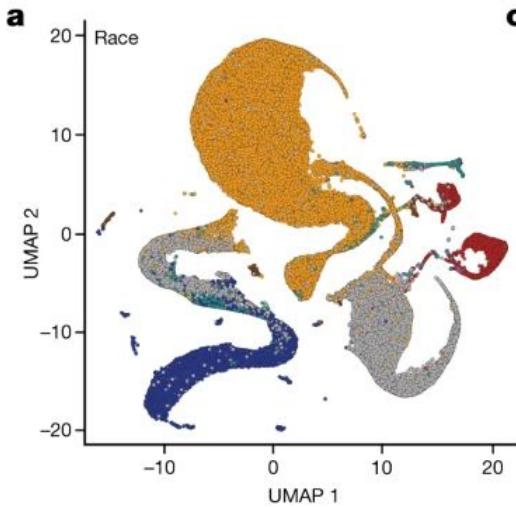
UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction

McInnes et al (2018) arXiv. 1802.03426

<https://www.youtube.com/watch?v=nq6iPZVUxZU>

<https://towardsdatascience.com/how-exactly-umap-works-13e3040e1668>

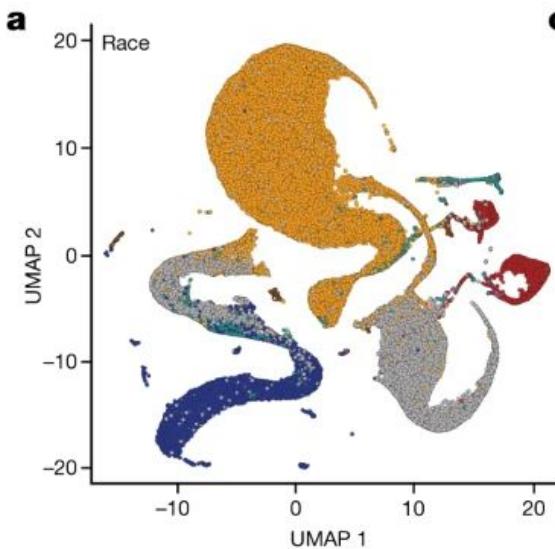
UMAP of Human Variation



Genomic data in the All of Us Research Program

The All of Us Research Program Genomics Investigators (2024) Nature. doi.org/10.1038/s41586-023-06957-x

~~UMAP of Human Variation~~



c “It's a pity that All of Us used UMAP to visualize ancestry variation in their new marker paper, out today in Nature. The UMAP algorithm, by design, exaggerates the distinctiveness of the most frequent ancestries, a message that can be misinterpreted by the public. UMAP pulls unusual genotypes towards the majority clusters; in particular it fails to represent admixture in a sensible way (admixture is fundamentally additive, while UMAP is not). In this setting the messiness of Admixture or PCA plots yield a better reflection of the data”

Jonathan Pritchard
Stanford

Genomic data in the All of Us Research Program

The All of Us Research Program Genomics Investigators (2024) Nature. doi.org/10.1038/s41586-023-06957-x

Introduction to Machine Learning: Classification & Prediction

Classification Problems

Supervised learning problem:

- Given (many) labeled training examples, develop a classifier that will accurately predict an output
- Each example contains **features** used to form a prediction

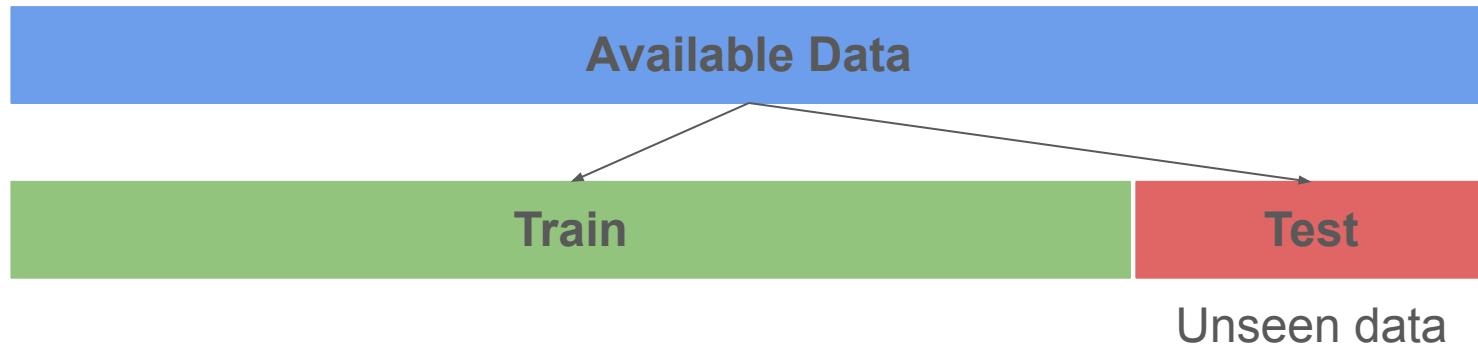
Examples:

- Given a pileup of reads, predict if there really is a variant
- Given SNVs in a person, predict ancestry (or relatedness)
- Given expression matrix, predict sample type
- Given SNVs + cell type, predict changes in expression
- Given SNVs/expression, predict disease risk
- Given SNVs/expression & disease status, predict treatment
- Given an image (photograph, xray, MRI), predict what it is (object, tumor, etc)
- Given English text, classify the sentiment (happy/sad) or predict the next word
- ...

Training datasets and cross validation

We need to do two things with the data:

1. Estimate the parameters for the ML method (**Training**)
2. Evaluate how well the ML method works (**Testing**)

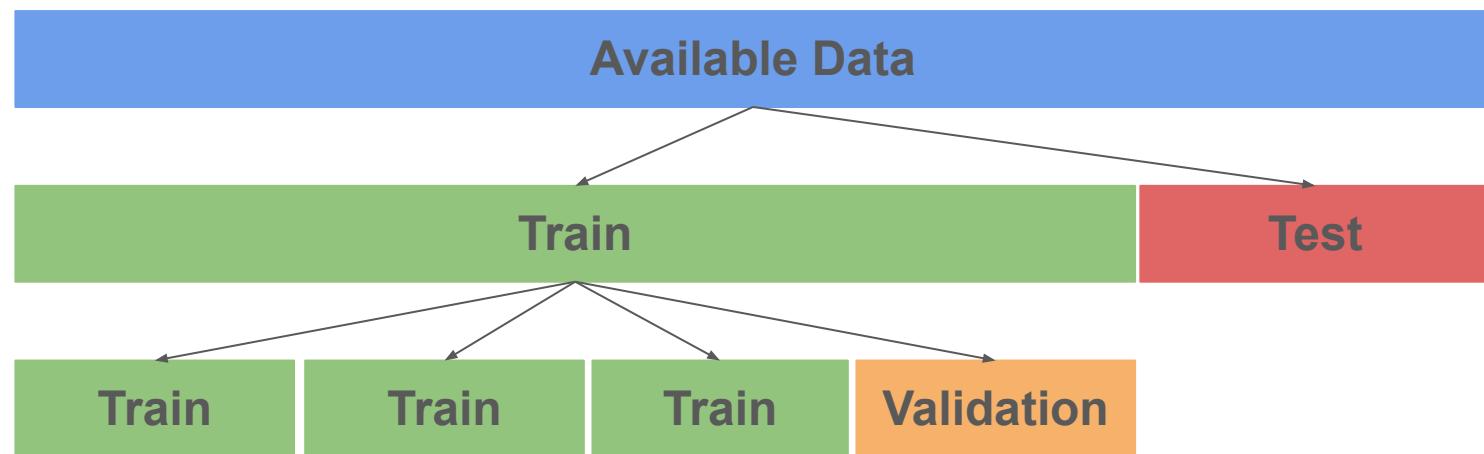


But how do we optimize the model?

Cross validation

We can split our training data into 75% train and 25% validation

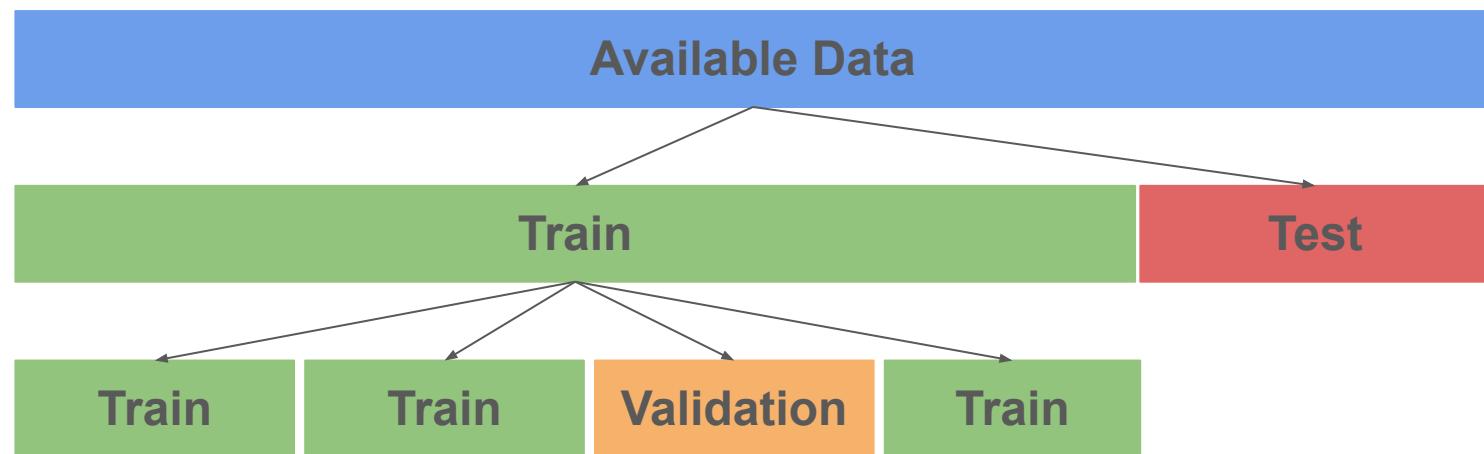
- But which block do we use for validation?



Cross validation

We can split our training data into 75% train and 25% validation

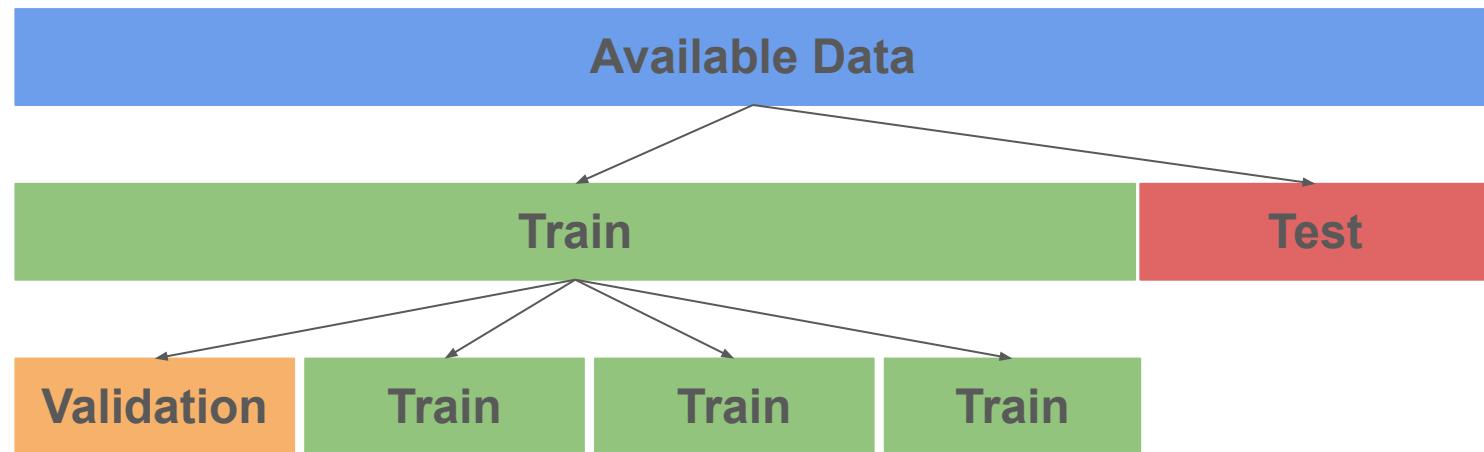
- But which block do we use for validation?



Cross validation

We can split our training data into 75% train and 25% validation

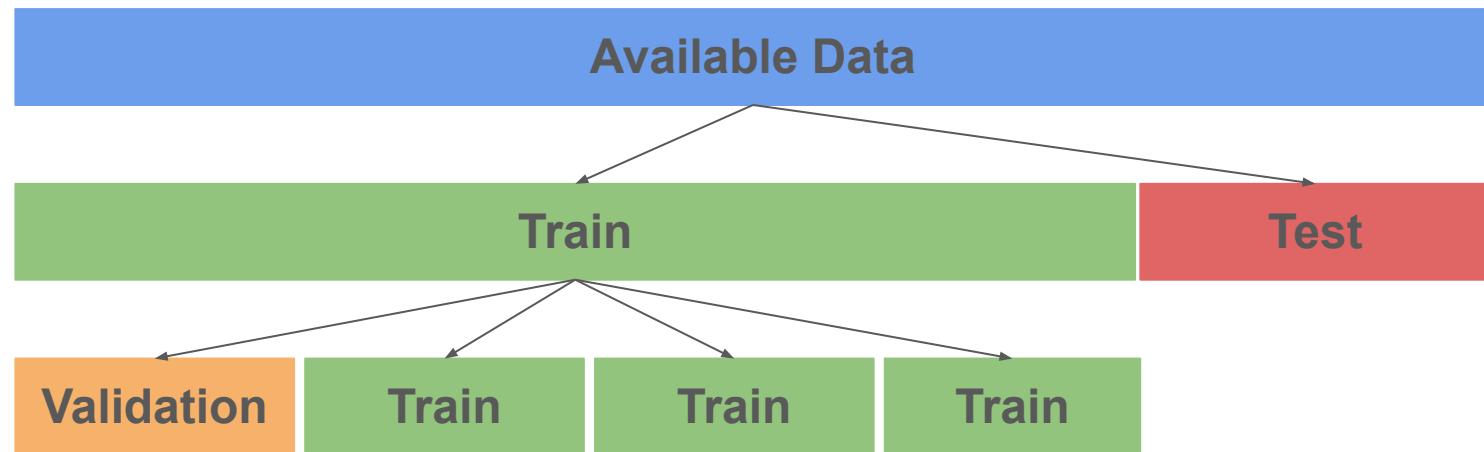
- But which block do we use for validation?



Cross validation

We can split our training data into 75% train and 25% validation

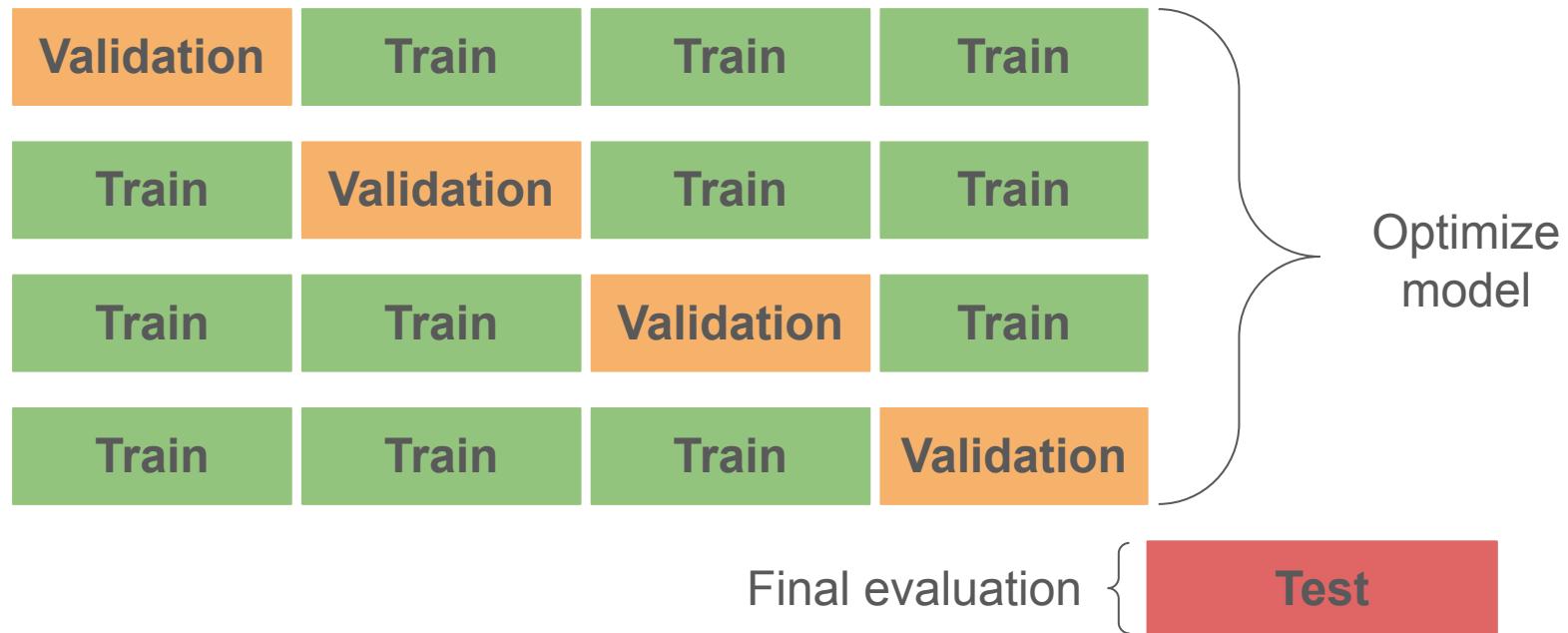
- But which block do we use for validation?
 - Why not all of them?



Cross validation

k-Fold cross validation

- Split training data into k splits



Evaluation metrics

Confusion matrix:

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Evaluation metrics

Confusion matrix:

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Accuracy:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Evaluation metrics

Confusion matrix:

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Accuracy:

- But what if the classes are unbalanced?

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Evaluation metrics

Precision:

- How many correctly predicted cases turned out to be positive
- False positive is of higher concern

$$Precision = \frac{TruePositive}{TruePositive + FalsePositive}$$

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Evaluation metrics

True Class

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Precision:

- How many correctly predicted cases turned out to be positive
- False positive is of higher concern

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

Recall (sensitivity):

- How many of the actual positive cases we were able to predict correctly
- False negative is of higher concern

$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

Evaluation metrics

Precision:

- How many correctly predicted cases turned out to be positive
- False positive is of higher concern

$$Precision = \frac{TruePositive}{TruePositive + FalsePositive}$$

Recall (sensitivity):

- How many of the actual positive cases we were able to predict correctly
- False negative is of higher concern

$$Recall = \frac{TruePositive}{TruePositive + FalseNegative}$$

F1 Score:

- Combines precision/recall (harmonic mean)

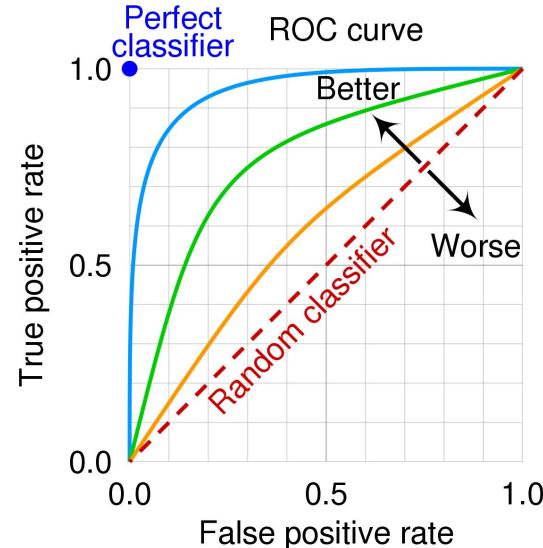
$$F1 = 2. \frac{Precision \times Recall}{Precision + Recall}$$

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Evaluation metrics

AUC-ROC

- Receiver Operator Characteristic (ROC) plots the true positive rate against the false positive rate
- Area Under the Curve (AUC) measure how the model distinguishes classes

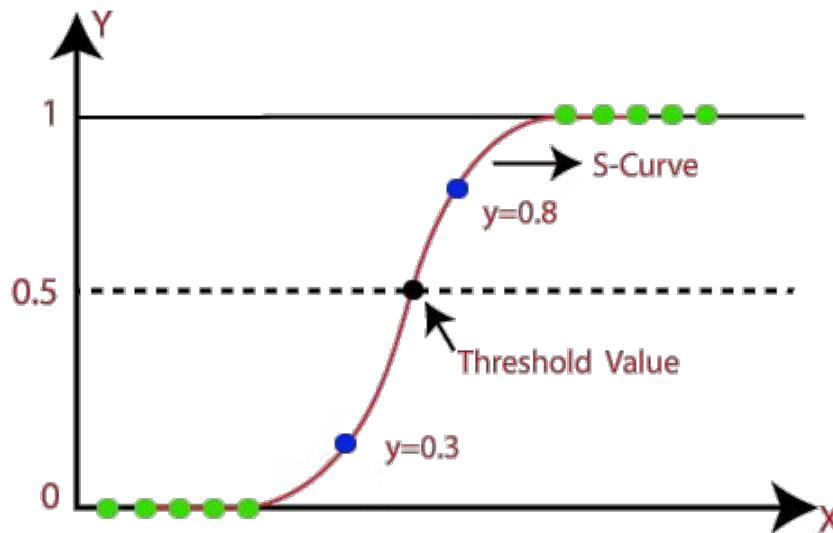


		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Logistic Regression

Logistic regression assigns a weight to each feature to form a prediction

$$\hat{y}(w, x) = w_0 + w_1x_1 + \dots + w_px_p$$

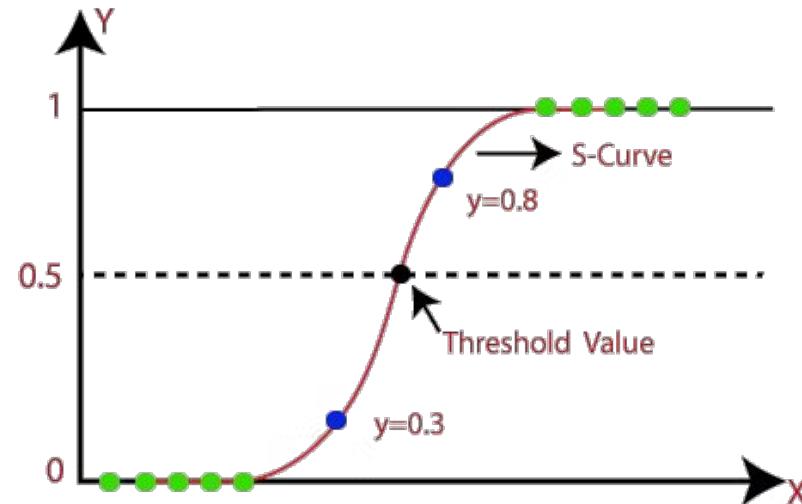


Logistic Regression

Logistic regression assigns a weight to each feature to form a prediction

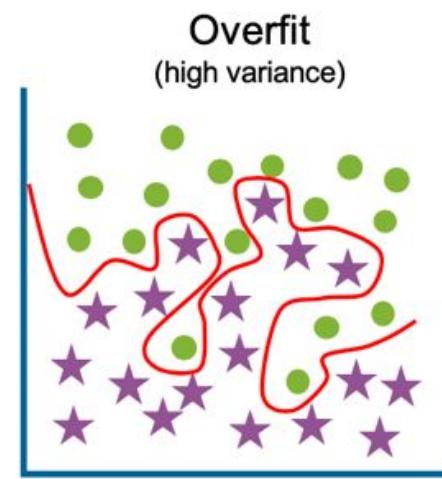
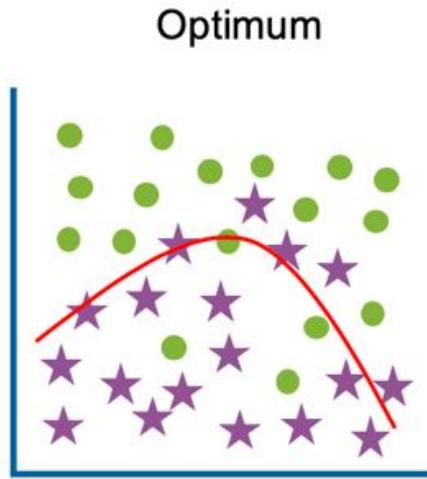
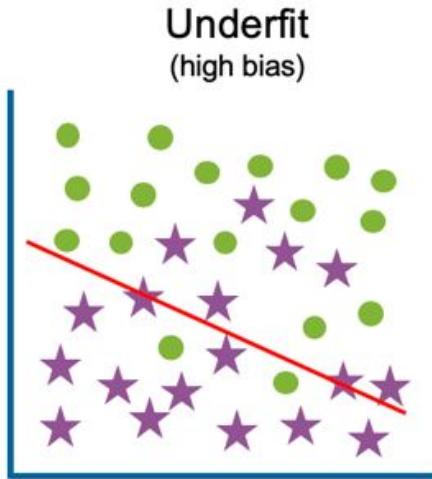
Pros	Cons
<ul style="list-style-type: none">• Easy to understand• Good accuracy when data is linearly separable• Less inclined to overfit	<ul style="list-style-type: none">• Assumes linearity between output and features• No multicollinearity between features

$$\hat{y}(w, x) = w_0 + w_1x_1 + \dots + w_px_p$$



Overfitting

Model is very accurate on training data but can't generalize to new data



High training error
High test error

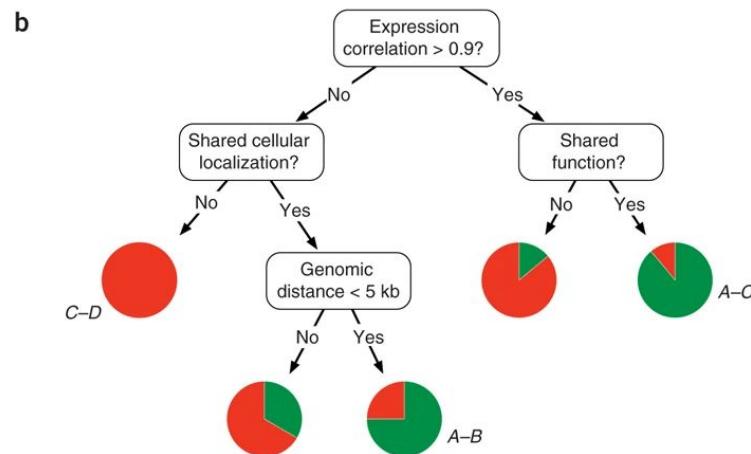
Low training error
Low test error

Low training error
High test error

Decision Trees

Decision tree learns simple decision rules inferred from features

a	Gene Pair	Interact?	Expression correlation	Shared localization?	Shared function?	Genomic distance
	A-B	Yes	0.77	Yes	No	1 kb
	A-C	Yes	0.91	Yes	Yes	10 kb
	C-D	No	0.1	No	No	1 Mb
	⋮					



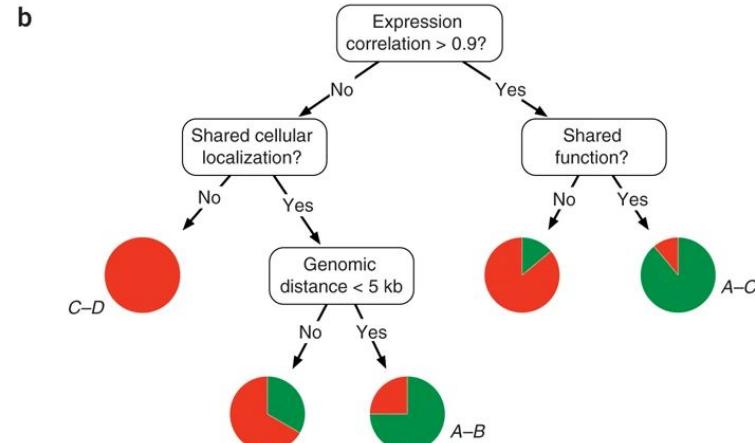
Decision Trees

Decision tree learns simple decision rules inferred from features

Pros	Cons
<ul style="list-style-type: none">• Easy to understand• Can be visualized• Can handle multiple outputs	<ul style="list-style-type: none">• Can create overcomplicated trees that overfit• Small variations in data can lead to completely different tree

a	Gene Pair	Interact?	Expression correlation	Shared localization?	Shared function?	Genomic distance
	A-B	Yes	0.77	Yes	No	1 kb
	A-C	Yes	0.91	Yes	Yes	10 kb
	C-D	No	0.1	No	No	1 Mb

:

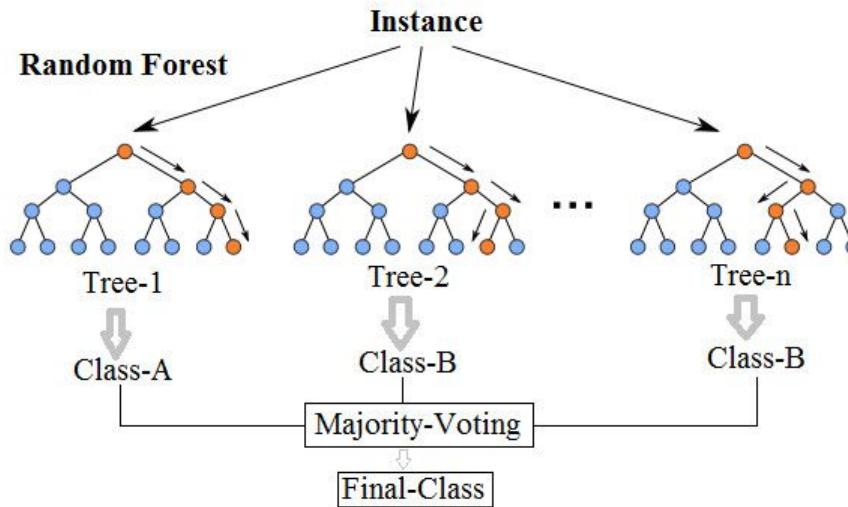


Random Forests

An ensemble of decision trees built from sample of training data

- Prediction is the averaged output of all decision trees

Random Forest Simplified

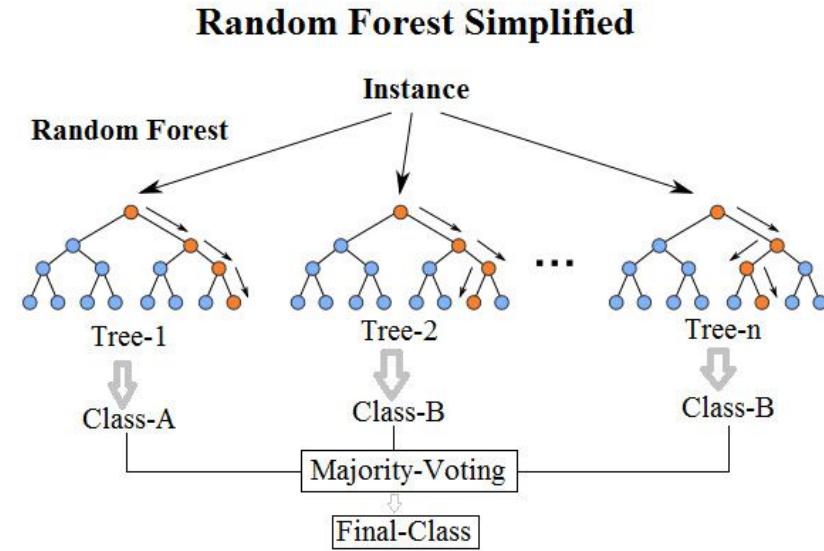


Random Forests

An ensemble of decision trees built from sample of training data

- Prediction is the averaged output of all decision trees

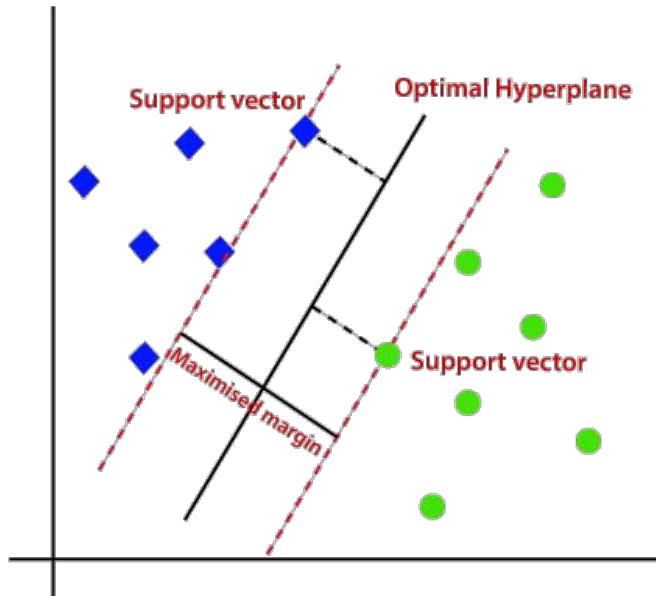
Pros	Cons
<ul style="list-style-type: none">Highly accurateIntrinsically interpretableCan learn non-linear decision boundariesParallel processing	<ul style="list-style-type: none">High computational complexityCan overfit because of noise



Support Vector Machine (SVM)

SVMs find a hyperplane that best separates the data into classes

- Maximize margin between different classes

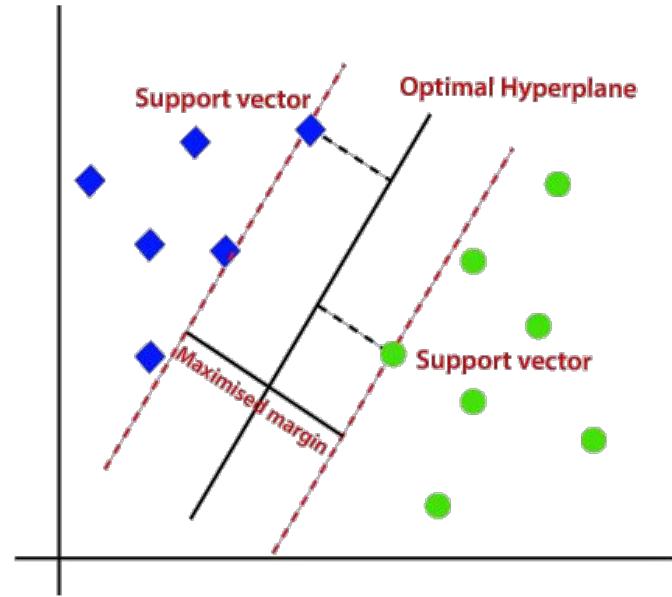


SVM

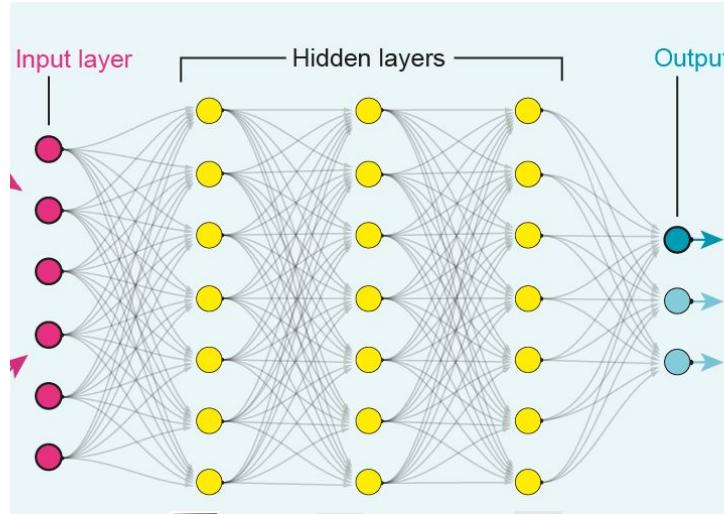
Support Vector Machines find a hyperplane that best separates the data into classes

- Maximize margin between different classes

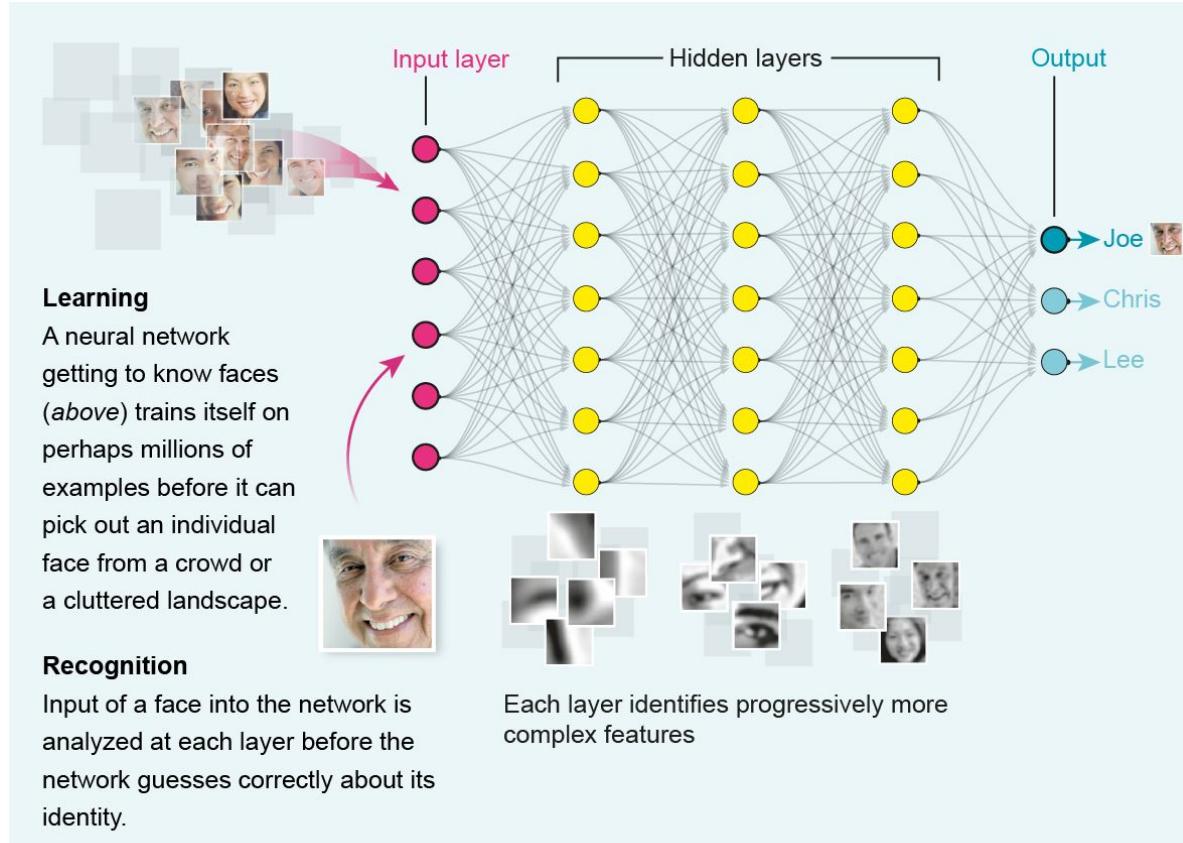
Pros	Cons
<ul style="list-style-type: none">• Effectively handles high dimensional data• Can work well with small datasets• Can model non-linear decision boundaries	<ul style="list-style-type: none">• Does not work for large data sets• Need to choose kernel• Limited to two-class problems• No probabilistic interpretation



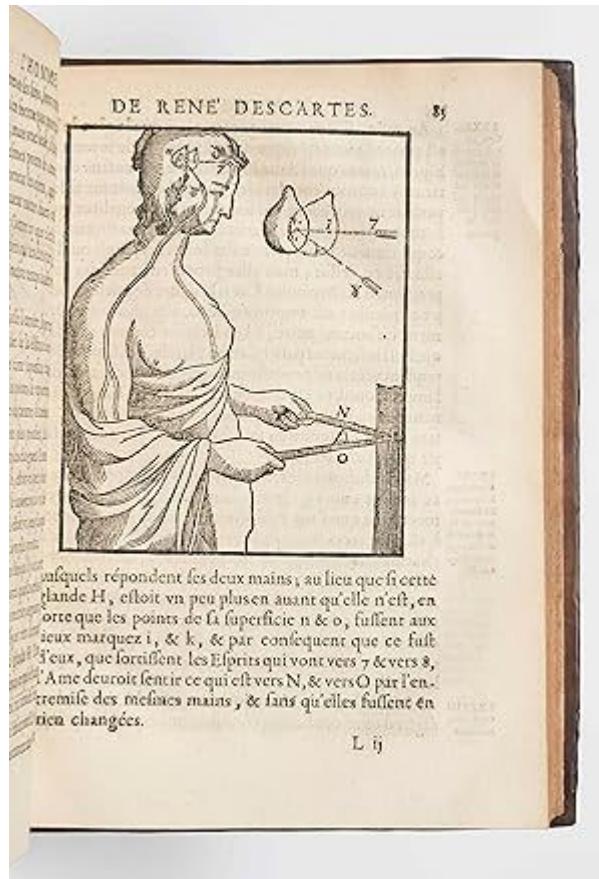
Artificial Neural Networks



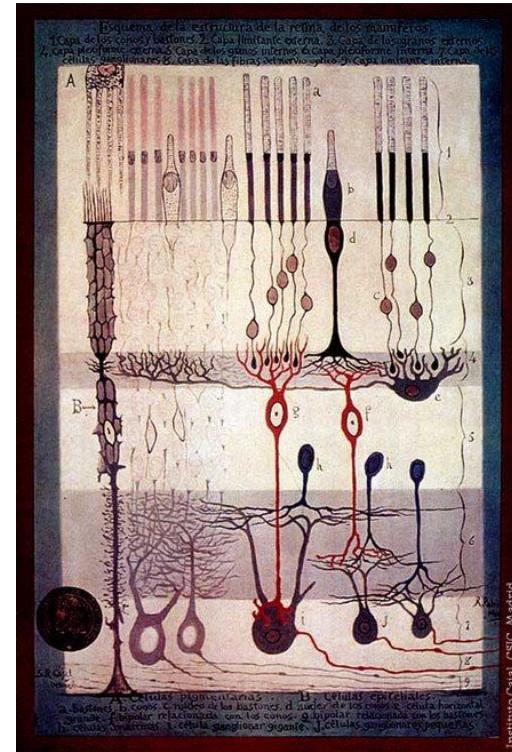
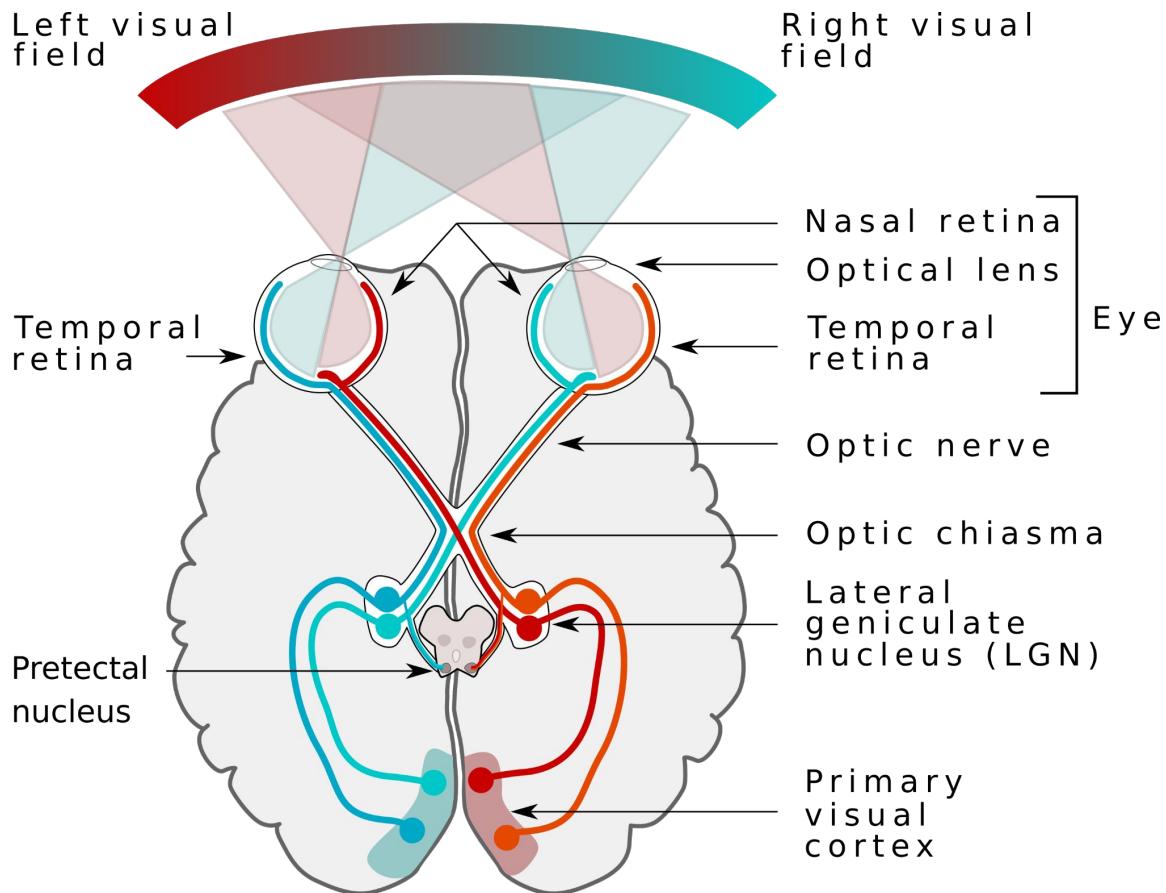
Artificial Neural Networks



Neural Circuitry

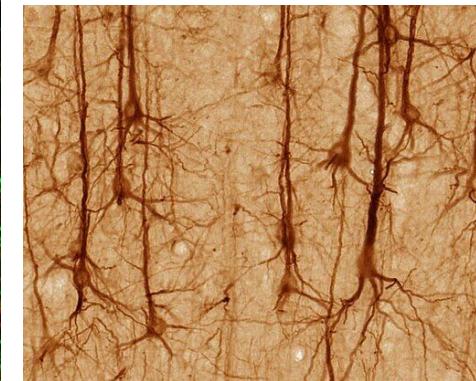
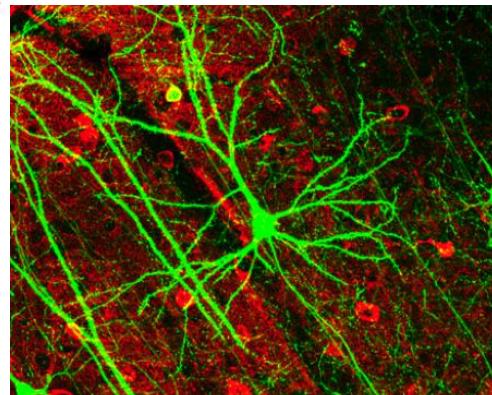
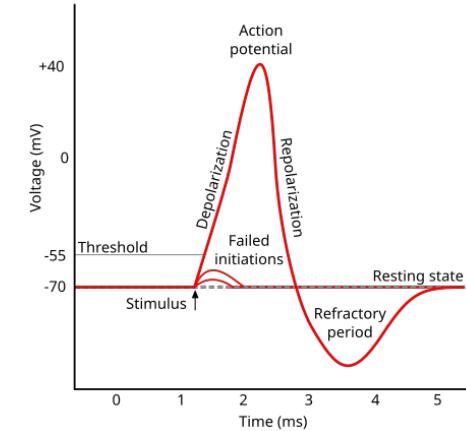
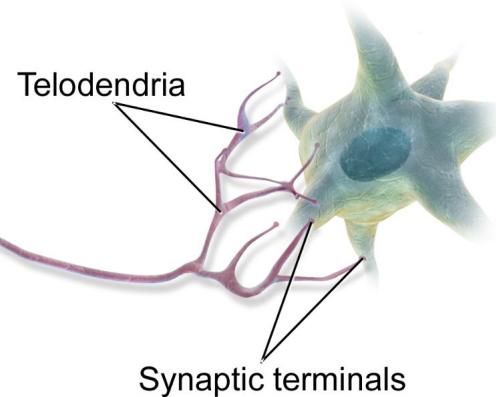
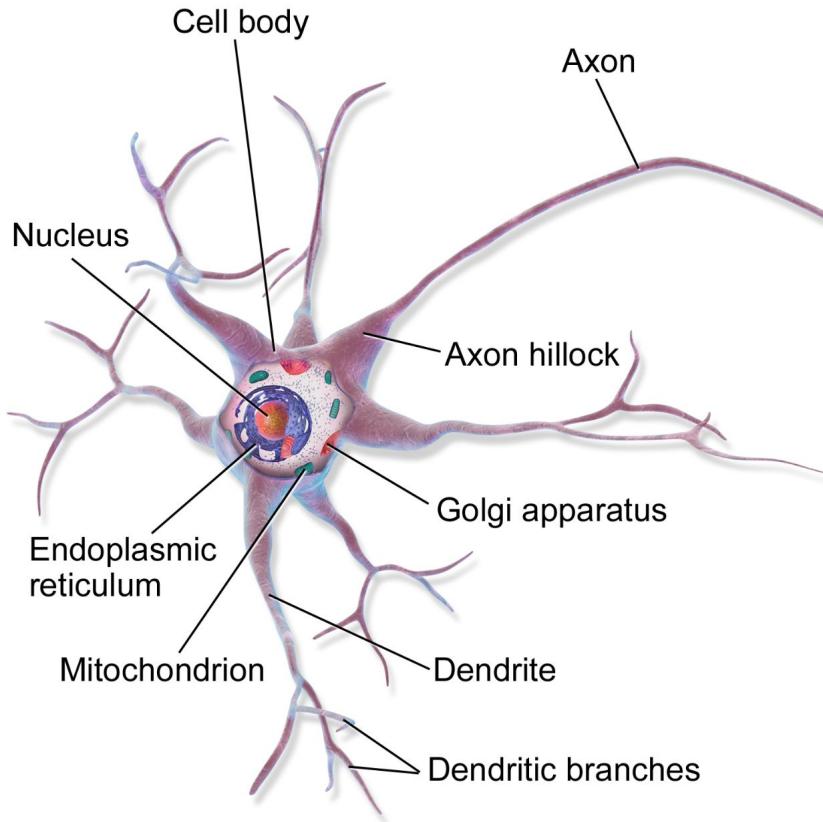


Human Visual Cortex



Structure of the Mammalian Retina S. Ramón y Cajal (1900)

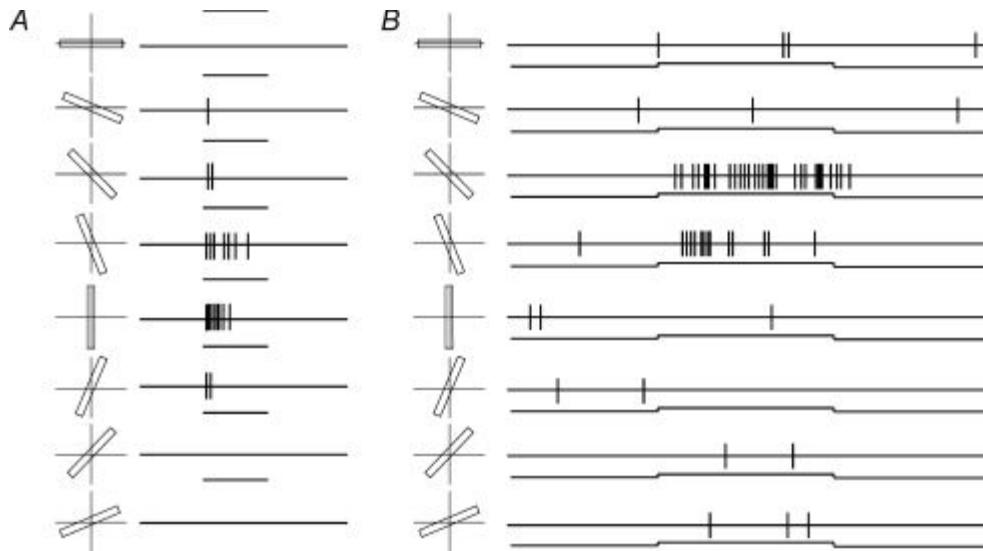
Neurons



Vision & Pattern Recognition

Receptive Fields, Binocular Interaction and Functional Architecture in the Cat's Visual Cortex

D. H. HUBEL AND T. N. WIESEL • Neurophysiology Laboratory, Department of Pharmacology Harvard Medical School, Boston, Massachusetts

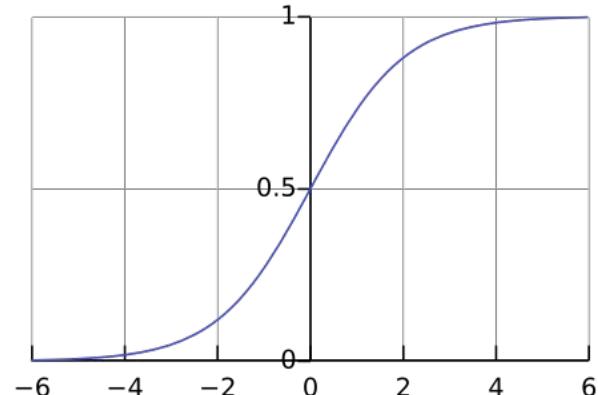
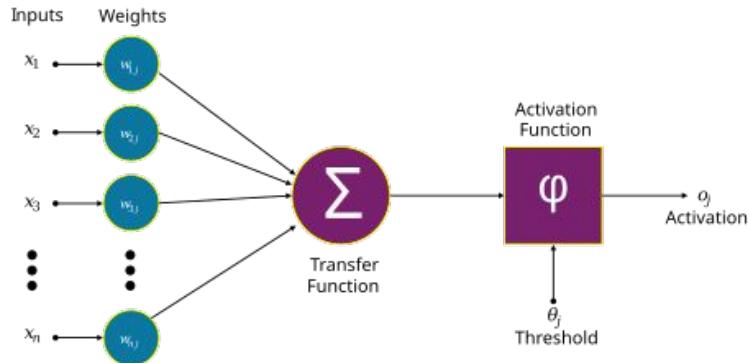
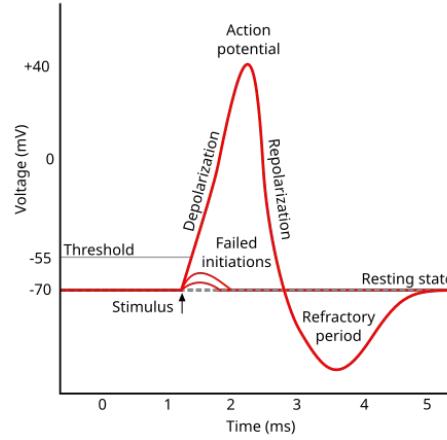
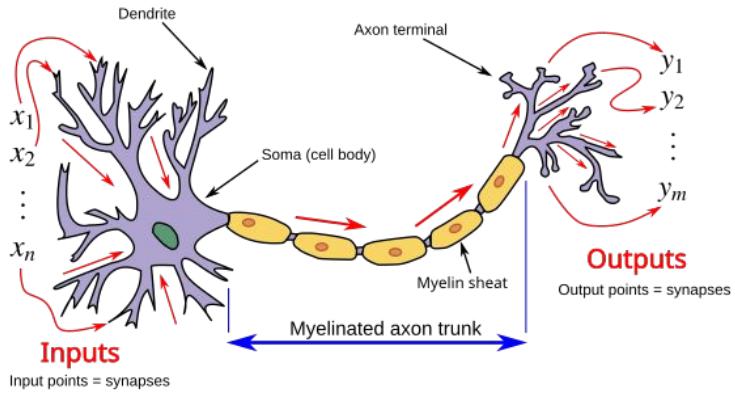


Example simple cells in the striate cortex of anaesthetized paralysed cat (A) and the awake fixating monkey (B). The comparison illustrates the similar response to oriented slits of light in the anaesthetized, paralysed cats and awake behaving monkeys. The example from the awake monkey shows the same qualitative orientation tuning but a slightly higher background rate than from the anaesthetized cat. Traced from Fig. 3A of Hubel & Wiesel (1959) and from Fig. 5 of (Wurtz, 1969c).

Table 1. Simple Cortical Fields

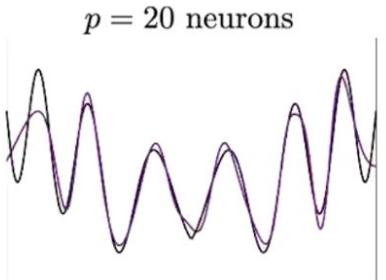
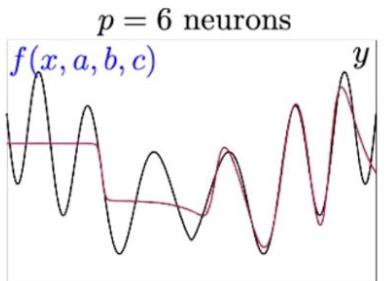
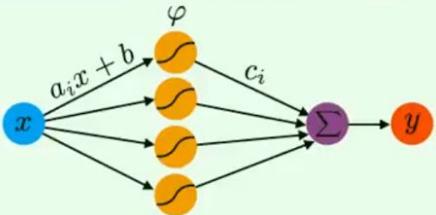
	Text-fig.	No. of cells
(a) Narrow concentrated centres		
(i) Symmetrical flanks		
Excitatory centres	2C	23
Inhibitory centres	2D	17
(ii) Asymmetrical flanks		
Excitatory centres	—	28
Inhibitory centres	2E	10
(b) Large centres; concentrated flanks	2F	21
(c) One excitatory region and one inhibitory	2G	17
(d) Uncategorized	—	117
Total number of simple fields		233

Biological vs Artificial Neurons



The non-linear activation function (sigmoid) is needed to model non-linear relationships

ANNs are “Universal Approximators”



1 hidden layer perceptron:

$$y \approx f(x, a, b, c) \stackrel{\text{def.}}{=} \sum_{i=1}^p c_i \varphi(a_i x + b_i)$$

Approximation by Superpositions of a Sigmoidal Function*

G. Cybenko†

Abstract. In this paper we demonstrate that finite linear combinations of compositions of a fixed, univariate function and a set of affine functionals can uniformly approximate any continuous function of n real variables with support in the unit hypercube; only mild conditions are imposed on the univariate function. Our results settle an open question about representability in the class of single hidden layer neural networks. In particular, we show that arbitrary decision regions can be arbitrarily well approximated by continuous feedforward neural networks with only a single internal, hidden layer and any continuous sigmoidal nonlinearity. The paper discusses approximation properties of other possible types of nonlinearities that might be implemented by artificial neural networks.

Key words. Neural networks, Approximation, Completeness.



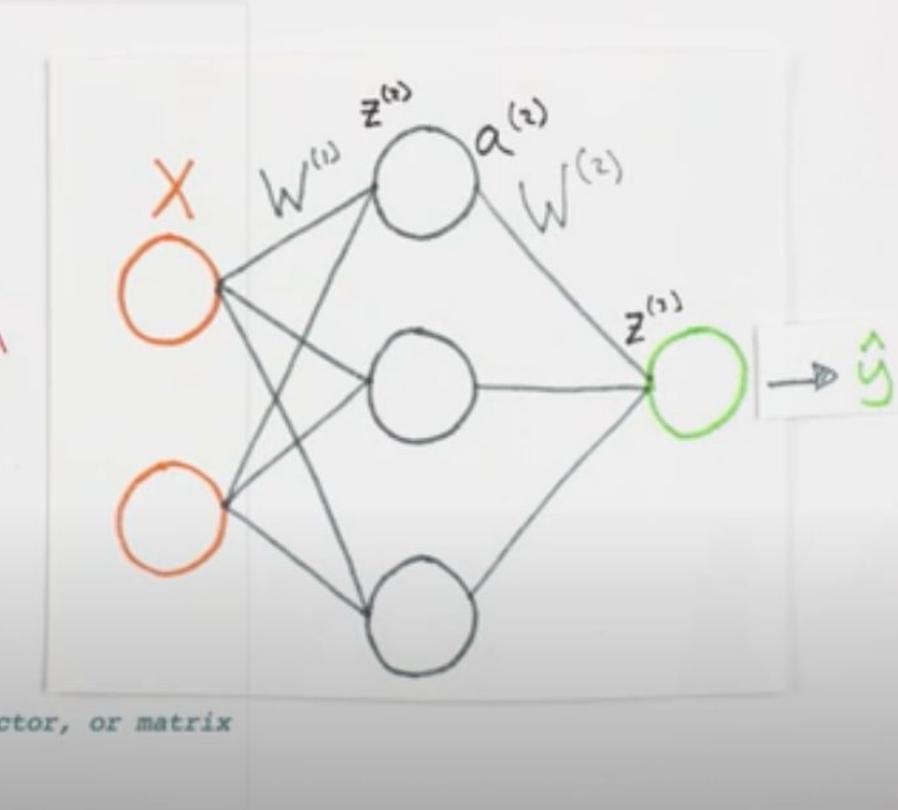
Training with Backpropagation: Random initialization

```
class Neural_Network(object):
    def __init__(self):
        #Define Hyperparameters
        self.inputLayerSize = 2
        self.outputLayerSize = 1
        self.hiddenLayerSize = 3

        #Weights (Parameters)
        self.W1 = np.random.randn(self.inputLayerSize, \
                                self.hiddenLayerSize)
        self.W2 = np.random.randn(self.hiddenLayerSize, \
                                self.outputLayerSize)

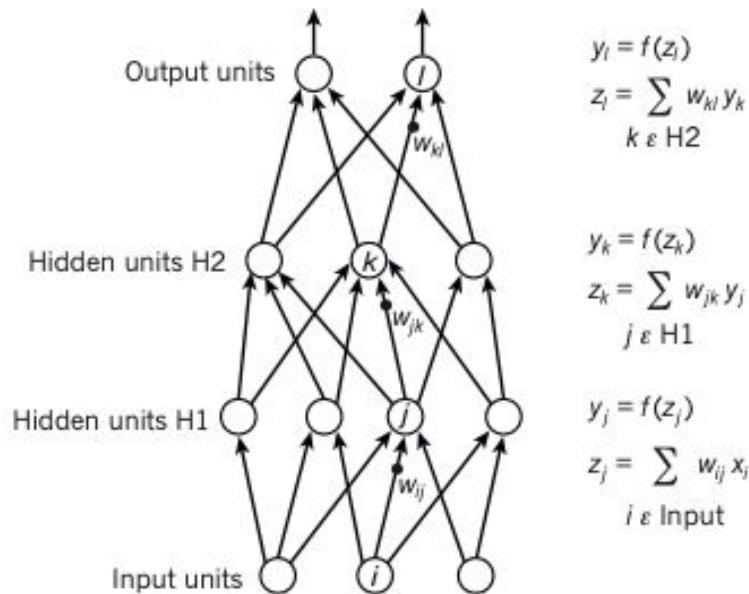
    def forward(self, X):
        #Propagate inputs though network
        self.z2 = np.dot(X, self.W1)
        self.a2 = self.sigmoid(self.z2)
        self.z3 = np.dot(self.a2, self.W2)
        yHat = self.sigmoid(self.z3)
        return yHat

    def sigmoid(self, z):
        #Apply sigmoid activation function to scalar, vector, or matrix
        return 1/(1+np.exp(-z))
```



Training with Backpropagation: Forward evaluation

c



$$y_l = f(z_l)$$
$$z_l = \sum_{k \in H2} w_{kl} y_k$$

$$y_k = f(z_k)$$
$$z_k = \sum_{j \in H1} w_{jk} y_j$$

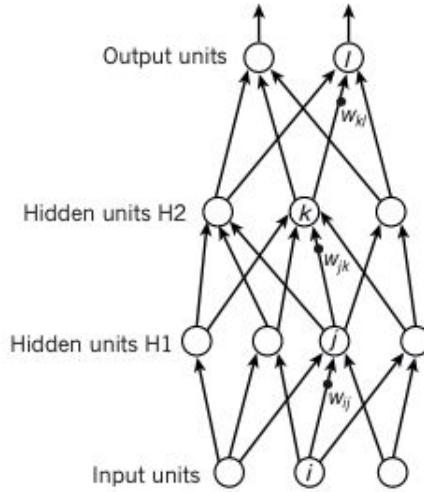
$$y_j = f(z_j)$$
$$z_j = \sum_{i \in \text{Input}} w_{ij} x_i$$

Deep Learning

LeCun, Bengio, Hinton (2020). Nature. doi:10.1038/nature14539

Training with Backpropagation: Backpropagation

c



$$y_i = f(z_i)$$
$$z_i = \sum_{k \in H2} w_{ki} y_k$$

$$y_k = f(z_k)$$
$$z_k = \sum_{j \in H1} w_{jk} y_j$$

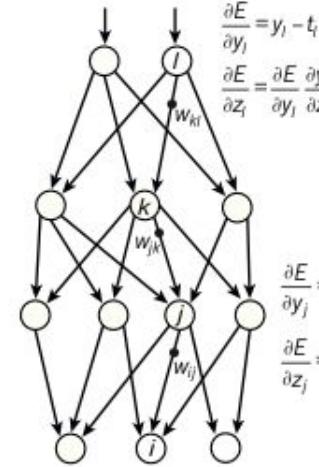
$$y_j = f(z_j)$$
$$z_j = \sum_{i \in \text{Input}} w_{ij} x_i$$

d

Compare outputs with correct answer to get error derivatives

$$\frac{\partial E}{\partial y_k} = \sum_{l \in \text{out}} w_{kl} \frac{\partial E}{\partial z_l}$$

$$\frac{\partial E}{\partial z_k} = \frac{\partial E}{\partial y_k} \frac{\partial y_k}{\partial z_k}$$



$$\frac{\partial E}{\partial y_j} = y_j - t_j$$

$$\frac{\partial E}{\partial z_j} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial z_j}$$

$$\frac{\partial E}{\partial y_j} = \sum_{k \in H2} w_{jk} \frac{\partial E}{\partial z_k}$$

$$\frac{\partial E}{\partial z_j} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial z_j}$$

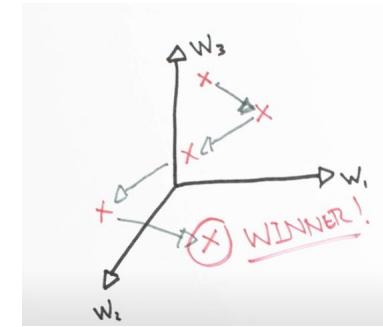
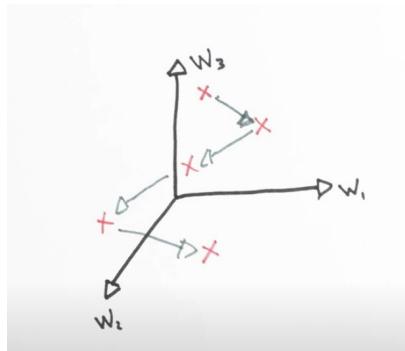
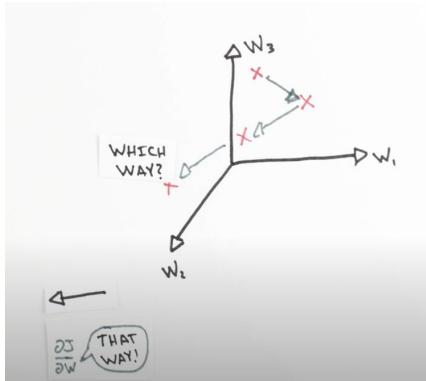
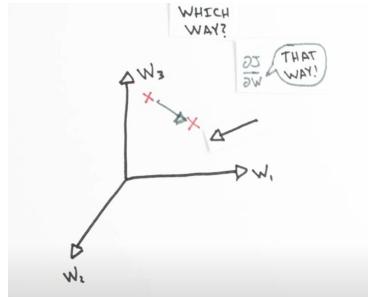
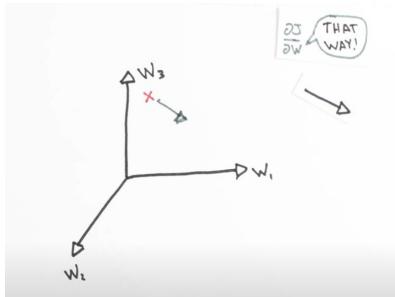
- Calculations along every edge, but the math is generally very basic
 - Addition, multiplication, exponentiation, etc
- Also easy to parallelize (especially with GPUs)

Deep Learning

LeCun, Bengio, Hinton (2020). Nature. doi:10.1038/nature14539

Training with Backpropagation: Gradient Descent

Minimize errors between predictions and true labels

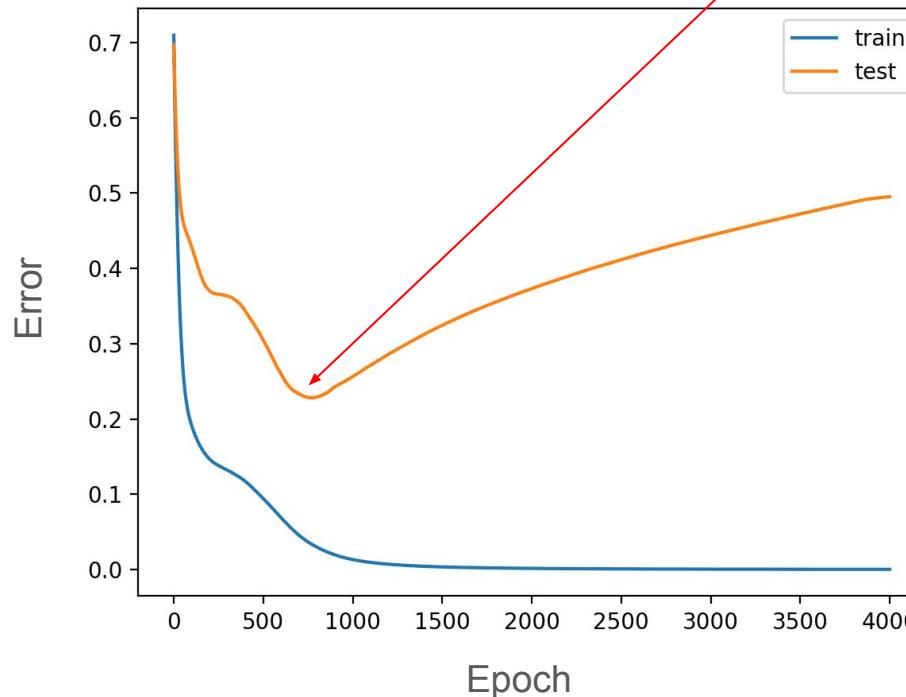


Ball rolling down a hill: start with a random initialization, then stochastic gradient descent by iteratively computing the partial derivatives of the functions

Training ANN through multiple epochs

Error goes down with more training epochs

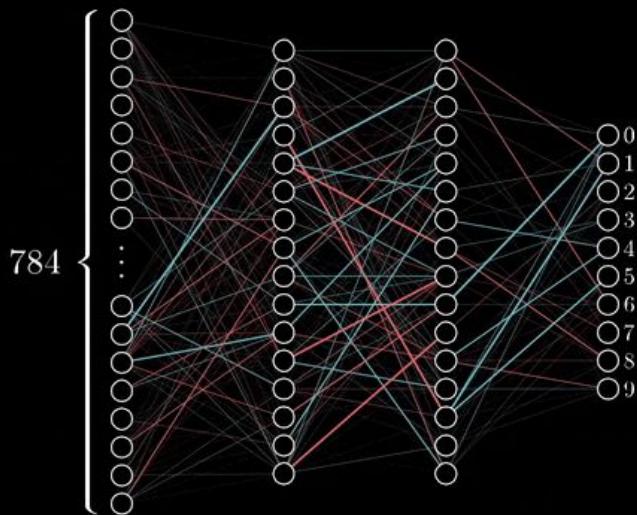
But what's happening here?



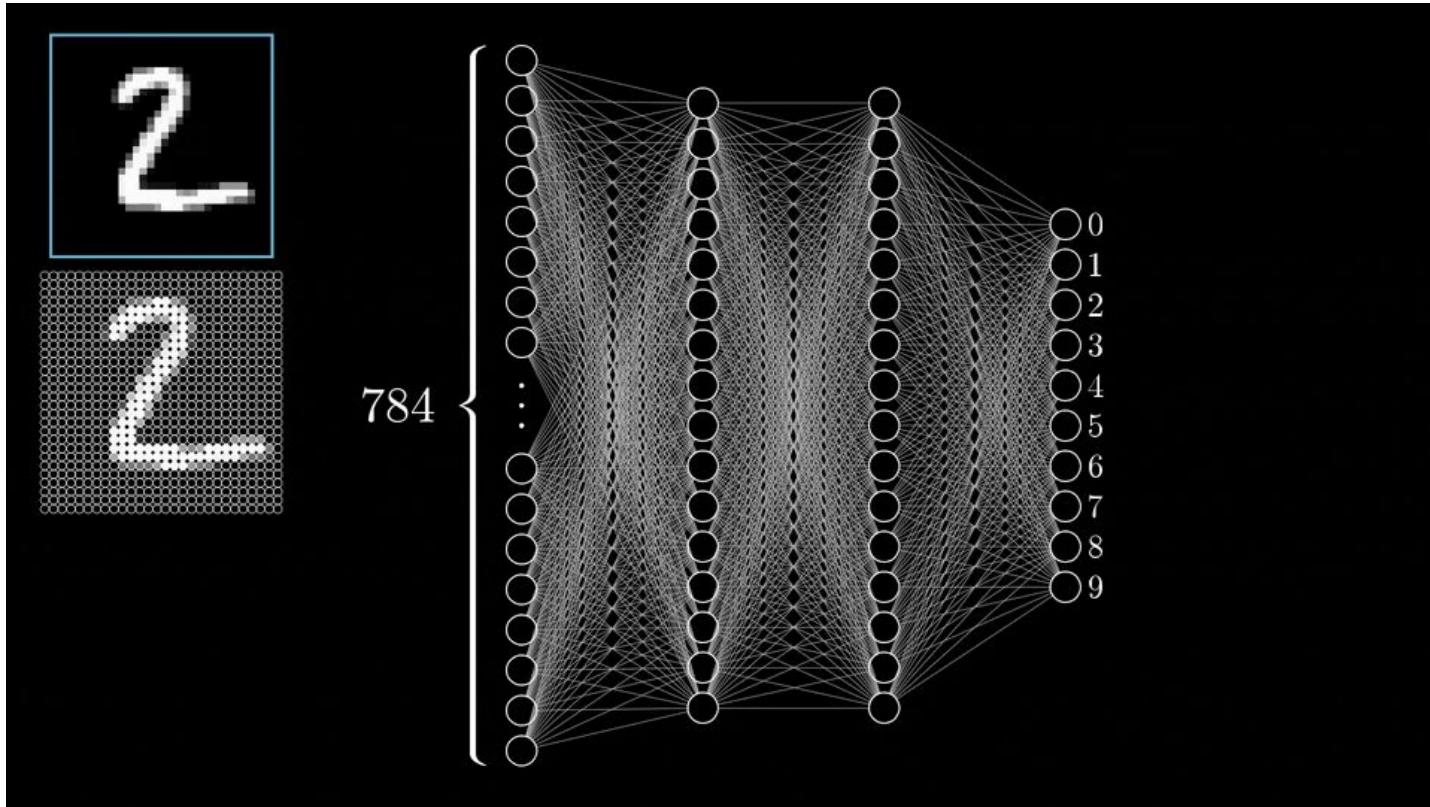
Training for too long
can lead to
overfitting!

Recap - Training

Training in
progress. . .



Recap - Evaluation



ANNs in Bioinformatics

GenNet: predicting phenotypes from genetic data

- Connections between layers defined through prior biological knowledge
- Basis for predicting many phenotypes
 - E.g. Schizophrenia, hair and eye color

