

Biomedical Informatics Research From Statistics to Machine Learning

Horacio Gómez-Acevedo, PhD

Assistant Professor

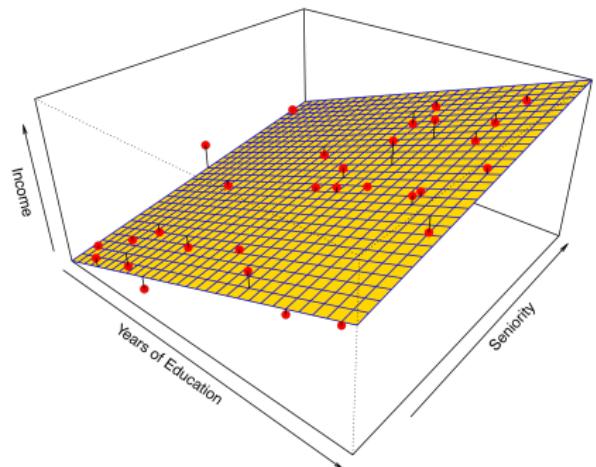
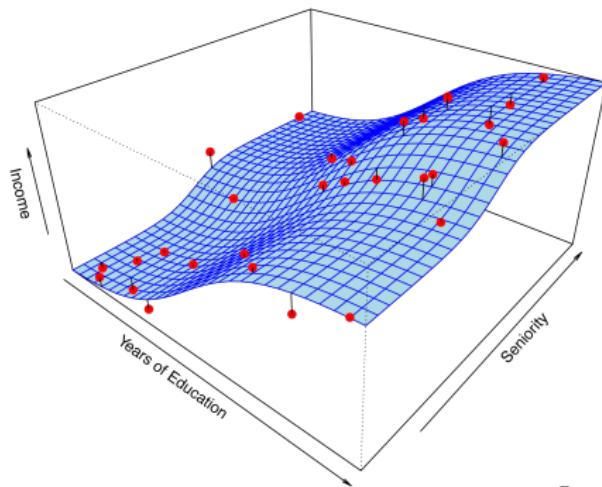
Department of Biomedical Informatics
University of Arkansas for Medical Sciences

March 4, 2022



Statistics and traditional parametric approach

$$\text{Income} = f(\text{Seniority}, \text{Years of Ed}) + \epsilon$$



$$\text{Income} \approx \beta_0 + \beta_1 * \text{Seniority} + \beta_2 * \text{Years of Ed}$$

Multi-linear regression

Statistics and traditional parametric approach

The philosophy of *classical* parametric paradigm is based on three beliefs

- ① The functional dependency from the data can be approximated with a linear function with a small number of parameters
- ② Random errors in real life problems follow a normal distribution.
- ③ Those parameters in the model can be calculated via the maximum likelihood method.

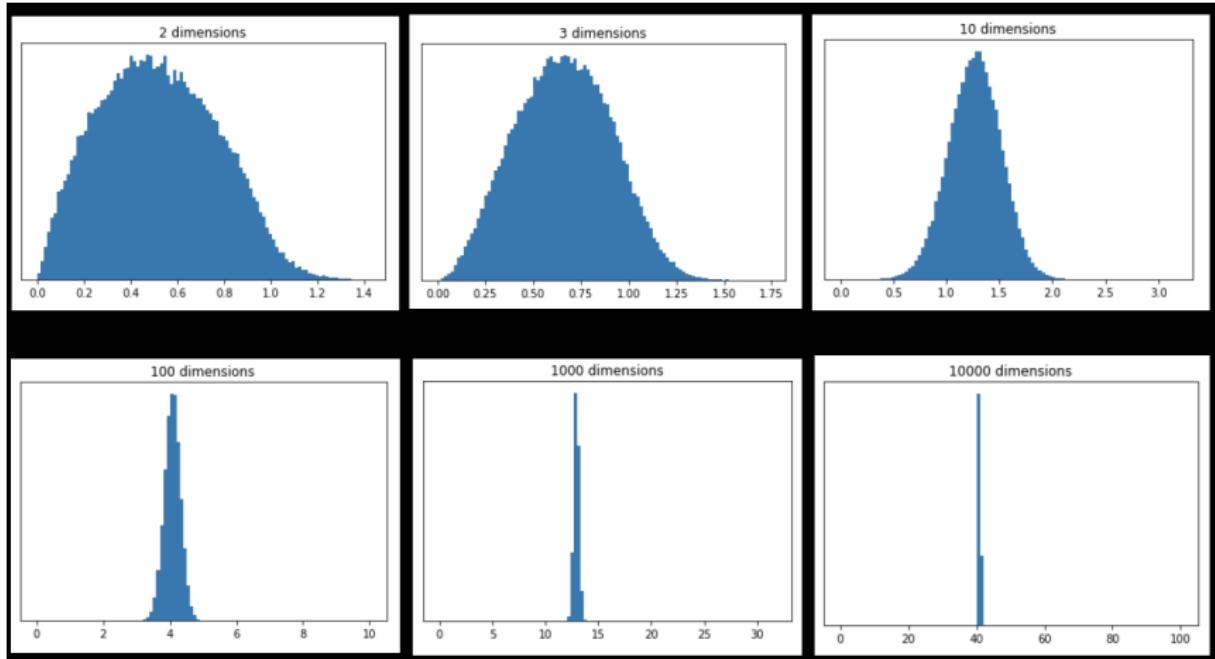
Problems with this approach

- It was until the computers were introduced that more challenge problems were attempted.
- For large multivariate problems it was observed that increasing the number of factors required more and more computational resources.
- R. Bellman called this phenomena the **curse of dimensionality**

New techniques were developed to make informal inferences of data instead of relying in purely statistical techniques.

Curse of Dimensionality

In higher dimensions our intuition is severely impaired



Random vectors in \mathbb{R}^n (Johh Urbanic, Pittsburgh Supercomputing Center)

What is Machine Learning?

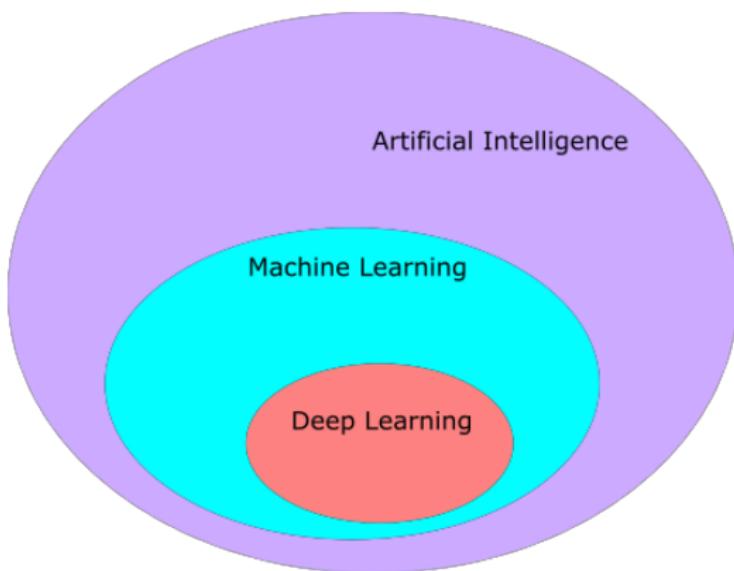
Machine learning (ML) is a vast field and here are some definitions .
(MI is the) field of study that gives computers the ability to learn without being explicitly programmed.

Arthur Samuel 1959

A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T as measured by P, improves with experience E.

Tom Mitchell 1997

Machine Learning ≠ Artificial Intelligence



When do you use ML?

- Problems for which existing solutions require a lot of hand-tuning or long lists of rules.
- Complex problems for which there is no good solution at all using a traditional approaches.
- Fluctuating environments as Machine Learning systems can adapt to not previously seen data.
- **Getting insights about complex problems and large amounts of data.**

Types of ML algorithms

ML algorithms can be classified according to the amount and type of supervision required during training. There are fundamentally four major categories:

- ① supervised
- ② unsupervised
- ③ semisupervised
- ④ reinforced learning

Supervised Learning

In Supervised Learning, the training data you feed the algorithm includes the desired solutions (called labels).

A typical supervised learning task is classification. The algorithm is trained with samples along with their class and the best algorithm that discerns best the labels would be used.

Another typical task is to predict a target numeric value given a set of features called predictors. We normally perform this task with (linear or multilinear) regression.

Supervised Learning (cont)

Other supervised learning algorithms include

- Linear regression
- Logistic regression
- Support Vector Machines (SVMs)
- Decision Trees and Random Forest

Separation by planes

A very common approach for classification of vectors is to find a (hyper)plane that separates the data

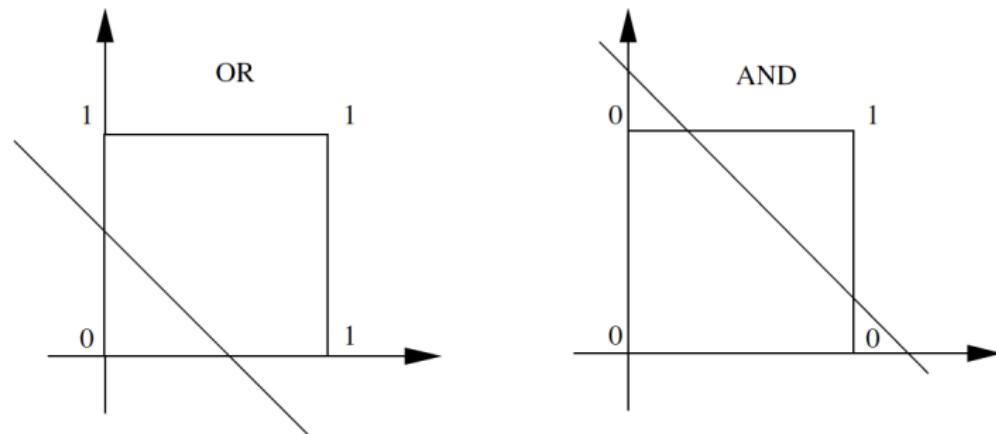
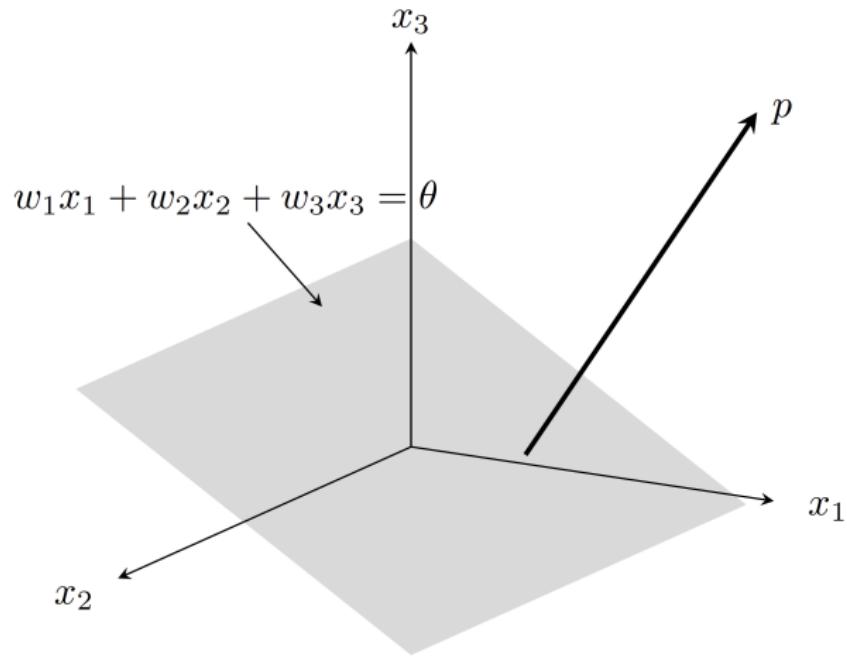


Fig. 3.6. Separations of input space corresponding to OR and AND

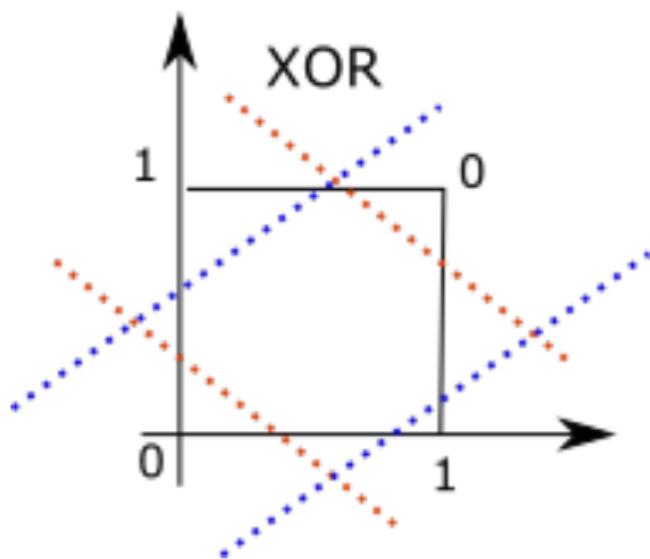
Separation by hyperplanes

More generally, we separate points (vectors) in higher dimensions with hyperplanes



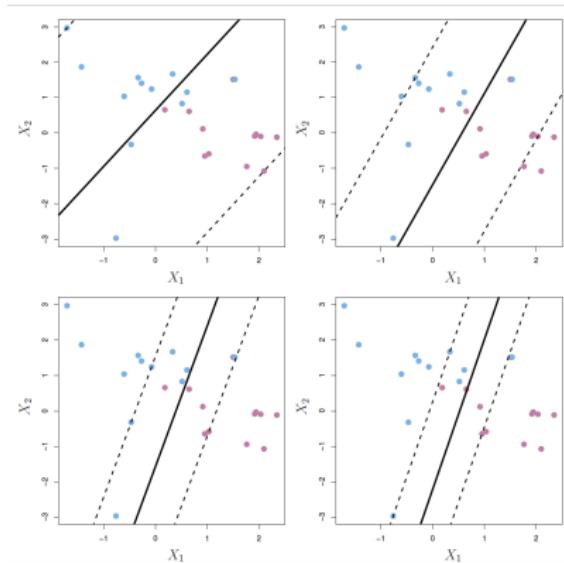
XOR function

However, even when we have a relatively simple function XOR this separation is not possible in \mathbb{R}^2 .



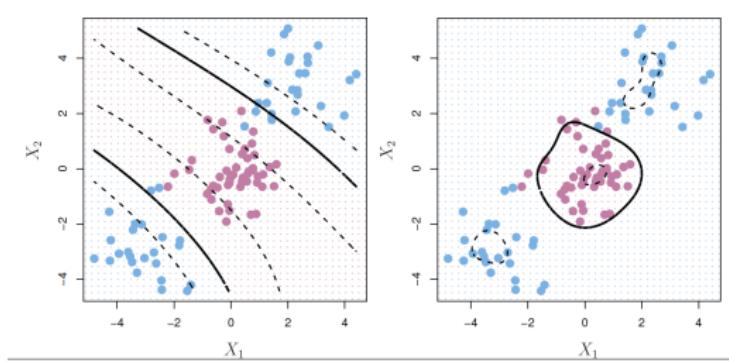
Support Vector Machines

Using an optimization procedure over the number of misclassified elements close to the hyperplane, we develop another separation algorithm based on the support vectors.



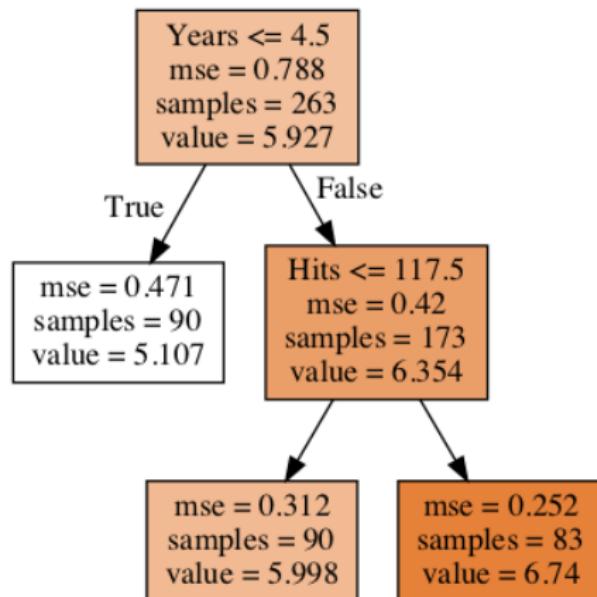
Support Vector Machines

This methodology allows us to use **kernels** to better capture the classification.



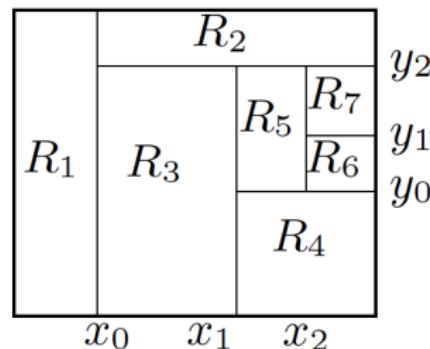
Tree-based methods

One of the most intuitive methods for classification is decision trees.



Partitions

Tree based methods can be thought as a way to partition our data set into smaller samples.



Feasible Partition

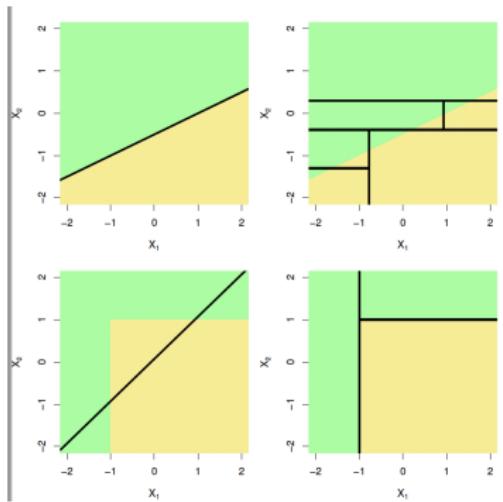
Random Forest

In the setup of classification, we conducted this procedure in two steps

- ① Construct binary decision trees using bootstrapped training sets (without pruning)
- ② Predict the value based on the majority vote. That is the class most commonly occurring will be selected.

Trees vs OLS

There are some setups where tree methods are more appropriate

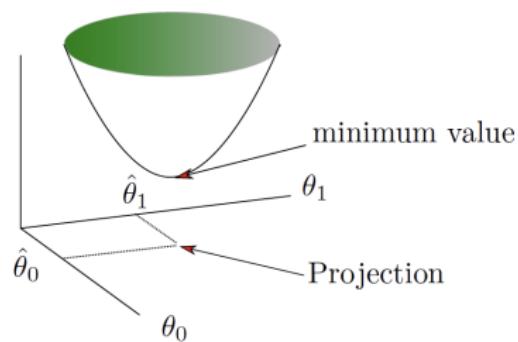


Random forest (boosting/bagging) are very powerful methodologies that compete with regression, their interpretation is complicated.

Model accuracy

Just a quick review of what we called *good fit* in least squares.

$$f(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$



This problem is computationally tracked by using a method called **gradient descent** to find the minima.

Unsupervised Learning

In unsupervised learning the training data is unlabeled and the system tries to learn without programmers intervention.

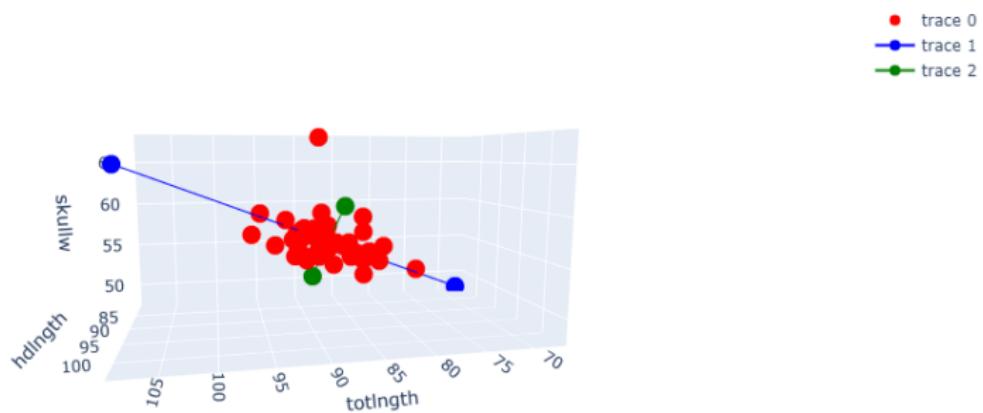
Some of the most important types of unsupervised learning are

- Clustering (K-means)
- Hierarchical Cluster Analysis (HCA)
- Visualization and dimensionality reduction
 - ① Principal Component Analysis (PCA)
 - ② Manifold Learning (t-SNE, MDS)

PCA

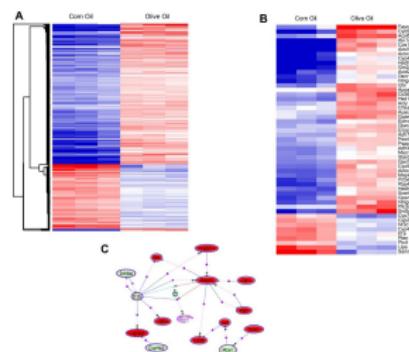
The idea of principal component analysis is to reduce the number of dimensions while preserving the data variability.

Principal Components



Clustering

The goal of clustering is to detect groups with similar characteristics. If you use hierarchical clustering algorithm it divides each group into smaller subgroups based on certain similarities.

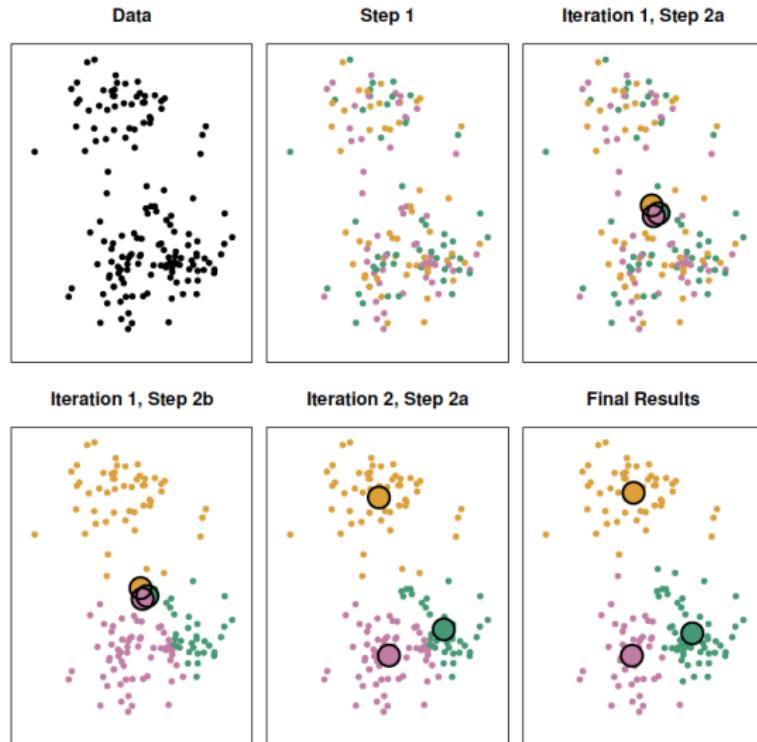


Classification with k-Means Clustering

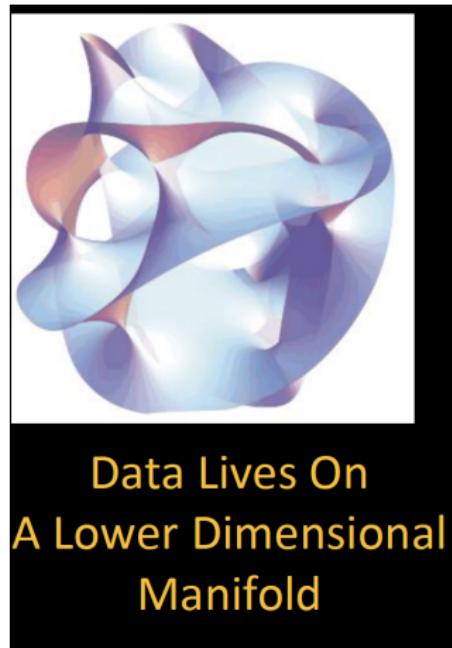
The algorithm goes like this

- ① Assign randomly a value between 1 and K to each data point.
- ② Iterate the following procedure until the clusters assignments stop changing
 - ① Find the centroid for each of the K clusters.
 - ② Each point will be assigned to the cluster K whose distance is the smallest. If two or more are equidistant, select randomly the cluster among the equidistant clusters.

K-Means Clustering (cont.)



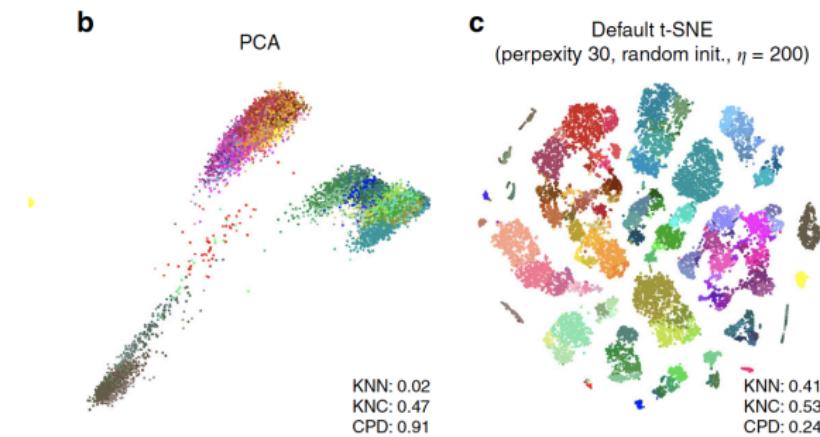
Manifold Hypothesis



John Urbanic, Pittsburgh Supercomputing Center

t-Stochastic Neighbor Embedding

This is another methodology that uses the data and their corresponding embedding (projection) and minimizes the Kullback-Leibler divergence of the normalized distribution of the data and their corresponding embedding.



Semisupervised Learning

Algorithms that handle a mixture of labeled and unlabeled data are called semisupervised. As the name suggests, this type of algorithms use a combination of supervised and unsupervised algorithms. For instance deep belief networks (DBNs) are based on unsupervised components called restricted Boltzmann machines that are trained sequentially in an unsupervised manner, and then the whole system is fine-tuned using supervised learning techniques.

Reinforced Learning

The learning system is called an agent in this context, and can observe the environment, select and perform actions and get rewards or penalties. It must then learn by itself what is the best strategy, called a policy, to the most reward over time. A policy defines what action the agent should choose when it is in a given situation. DeepMind's Alpha Go program used reinforced learning to beat world champion Ke Jie at the game of Go. It learned its winning policy by analyzing millions of games, and then playing many games against itself.

Challenges in ML

There are either bad algorithms or bad data that can derail any serious ML effort.

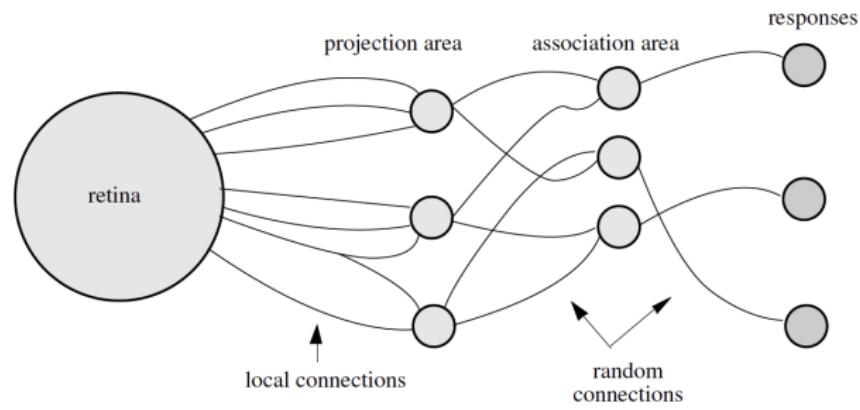
- Insufficient quality of training data. For very simple problems you typically need thousands of samples, and for complicated such as image or speech recognition you may need millions of samples.
- Nonrepresentative training data To generalize a ML algorithm you need that your training data be representative of the new cases you want to generalize. If the sample is too small, you will have sampling noise, but even very large samples can be nonrepresentative if the sampling method is flawed. This is called sampling bias

Challenges in ML (cont)

- Poor quality data. If your training data is full of errors, outliers and noise, it will make the ML methodology to under perform. Data quality is a field on its own right. Missing features can happen at random or being systematic. Spend enough time with your data to decide quality or features that are missing loads of information.
- Overfitting This occurs when the model is too complex relative to the amount and noisiness of the training data. Regularization is the process of reducing overfitting by making the model simpler by introducing a hyperparameter (this is a parameter not for the model but for the algorithm).

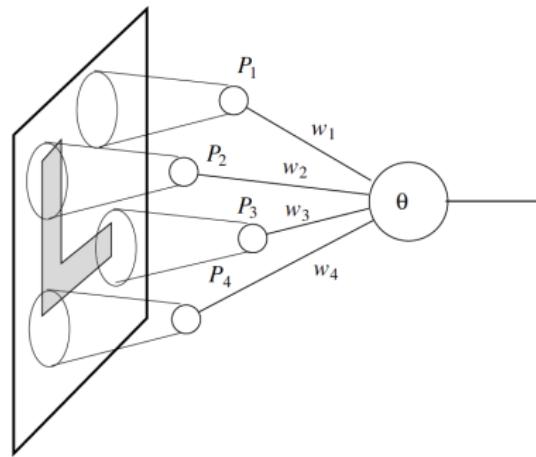
The perceptron

In 1958 Frank Rosenblatt proposed the perceptron which is a more general computational model. The essential innovation was the introduction of numerical weights and a special interconnection pattern. In the original Rosenblatt model the computing units are threshold elements and the connectivity is determined stochastically. Learning takes place by adapting the weights of the network with a numerical algorithm. The so-called classical perceptron is depicted below



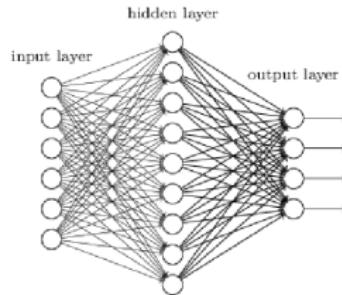
Threshold logic

The computational equivalent is a network with weights w_i and a number of input units P_i . If the combined value of the inputs is larger than θ , then it fires a 1, otherwise a 0.

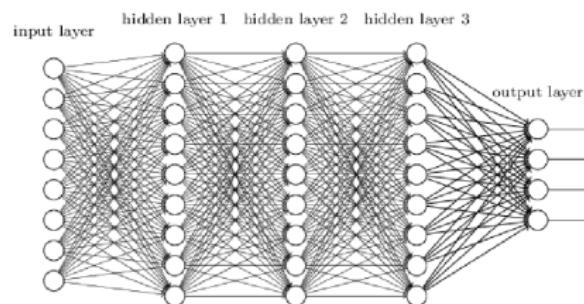


Deep Neural Networks

"Non-deep" feedforward neural network



Deep neural network



Deep Learning

"Deep learning achieves great power and flexibility by representing the world as a nested hierarchy of concepts, with each concept defined in relation to simpler concepts and more abstract representations computed in terms of less abstract ones."

Goodfellow et al.

Deep learning is only possible because of the advances in computer power (GPUs in particular). The basic concepts date back to the end of 1990's with Y. LeCun and colleagues.

Deep Learning in real time

One impressive representation of deep learning applied to a classification problem

<https://www.cs.ryerson.ca/~aharley/vis/conv/flat.html>

Deep Learning pitfalls

Applications of deep learning and architectures is abundant, and results really impressive. However, there are still some fundamental theoretical problems that have not been solved. For instance the **Black box problem** in which is not clear how the deep neural network makes decisions and how can it be modified.

What is Biomedical Informatics (BMI)

The American Medical Informatics Association (AMIA) gives a broad definition of the subject.

Biomedical and health informatics applies principles of computer and information science to the advancement of life sciences research, health professions education, public health, and patient care.

- BMI develops, studies and applies theories, methods and processes for the generation, storage, retrieval, use, and sharing of biomedical data, information, and knowledge.
- BMI builds on computing, communication and information sciences and technologies and their application in biomedicine.
- BMI investigates and supports reasoning, modeling, simulation, experimentation and translation across the spectrum from molecules to populations, dealing with a variety of biological systems, bridging basic and clinical research and practice, and the healthcare enterprise

Areas of Research at DBMI

The Department of Biomedical Informatics (DBMI) has four speciality areas

- Translational Bioinformatics
- Clinical Research Informatics
- Imaging Informatics
- Clinical Informatics (mostly for MDs)

Omics in medicine

Since the molecular function is not restricted to the genome, other modalities are also informative

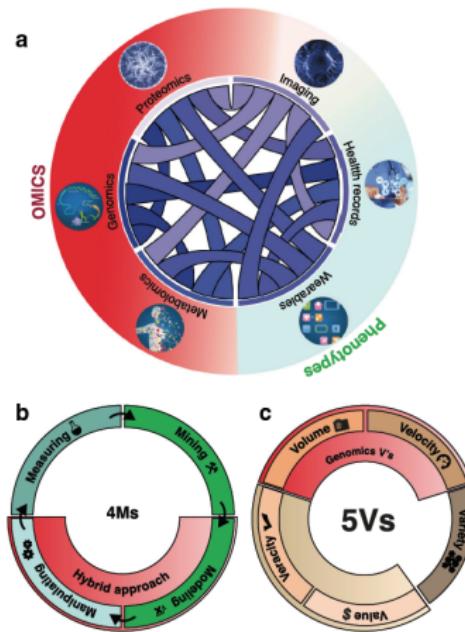


Fig. 8: Other Omics in healthcare

Where does Computer Sciences fit?

The short answer is EVERYWHERE!

- Algorithm development
- Pipeline development
- Cloud computing
- Data mining
- Data quality
- Data integration

Where does Statistics/Biostatistics fit?

Statistical machine learning is one of the driven forces for Biomedical Informatics Research.

- Random Forest
- Lasso and Ridge Regression
- Support Vector Machines
- Clustering Algorithms
- Bayesian modeling

Algorithm Development

In Genomics, one of the first tasks is alignment of reads (small pieces of DNA from samples) to the reference genome (transcriptome).

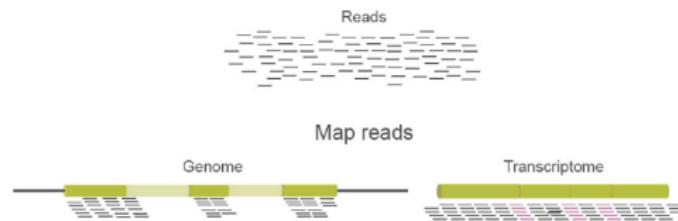


Fig 9: Alignment of reads to Genome/Transcriptome

Alignment Problem

There are algorithms to find what is the best “matched” configurations and also determine what are the possible variants.

Pos	1	2	3	4	5	6	7	8	9
Ref	T	A	G	C	C	G	A	T	C
r1	T	A	G	C	C	G	A		
r2	T	A	G	C	C	G	A		
r3	T	A	—	C	C	A	G	A	
r4	T	A	G	C	C	H	H		
r5	T	A	G	C	C	G	A	T	C
r6	S	S	G	C	C	G	A		
r7			G	C	C	G	A		

Pipeline Development

Since software implementations differ (inputs/outputs) and newly developed software keeps appearing, there is a need to keep up to date solutions.

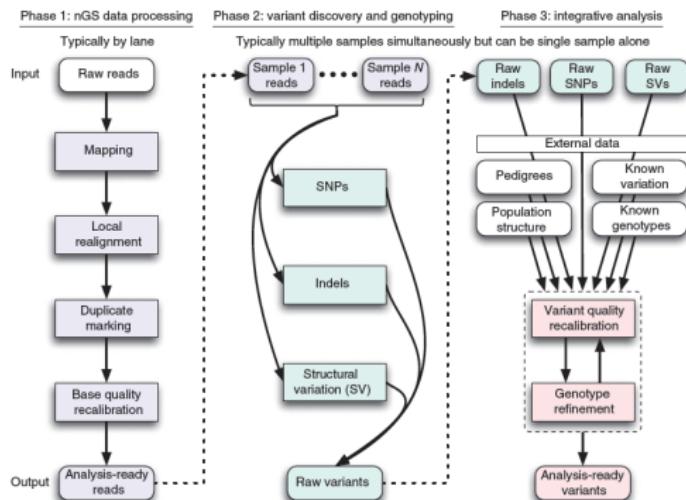


Figure 1 Framework for variation discovery and genotyping from next-generation DNA sequencing.

Fig 12: Broad Institute Best Practices (old)

Cloud computing solutions for Genomics

Google, Amazon and other cloud providers have already proposed platforms for genomics data that depends heavily on containers (Docker or Kubernetes) and workflow management systems (e.g., Cromwell)

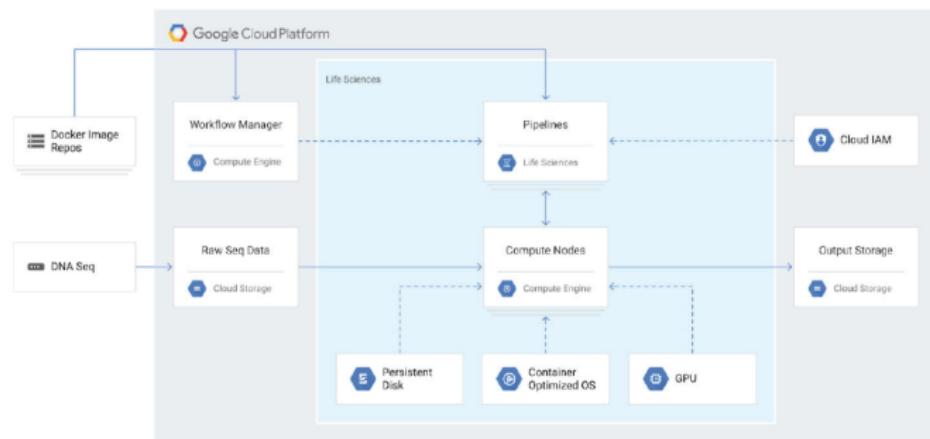
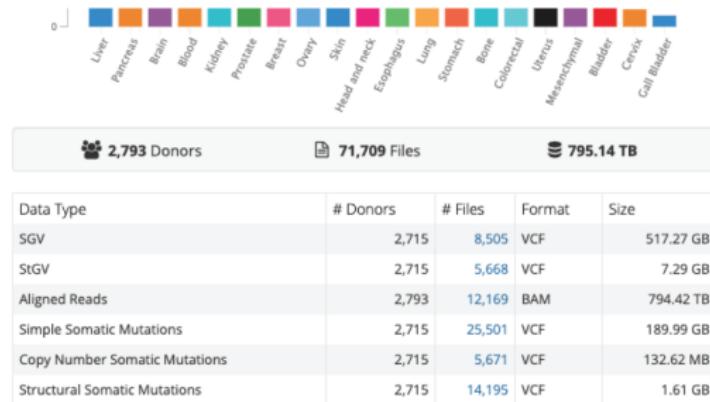


Fig 13: Google Cloud Arquitecture

Data Mining

If you are interested in bioinformatics for Cancer Research



Available data as of Jan 23, 2020

Fig 14: Pan-Cancer Consortium

Machine Learning

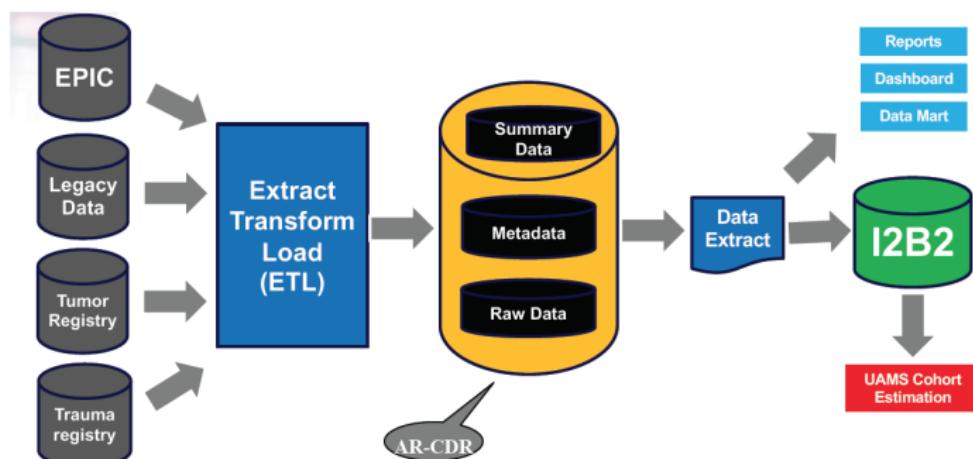
There is current research in methods for current events

The screenshot shows the homepage of the COVID-19 Open Research Dataset Challenge (CORD-19). At the top, there is a navigation bar with links for "Dataset", "Tasks", "Kernels", "Discussion", "Activity", and "Metadata". Below the navigation bar, the title "COVID-19 Open Research Dataset Challenge (CORD-19)" is displayed, along with the subtitle "An AI challenge with AI2, CZI, MSR, Georgetown, NIH & The White House". To the right of the title is a circular icon with a yellow dot and several red arrows pointing towards it, with the number "4910" next to it. Below the title, there is a banner for "AI2 Allen Institute For AI and 8 collaborators • updated 2 days ago (Version 6)". Further down, there are tabs for "Data", "Tasks (10)", "Kernels (606)", "Discussion (221)", "Activity", and "Metadata". A "Download (6 GB)" button is also present. On the far right, there is a "New Notebook" button and a "Back to task list" link. The main content area features a section titled "What is known about transmission, incubation, and environmental stability?" with a "900" count next to it. Below this section, there is a link to "COVID-19 Open Research Dataset Challenge (CORD-19)" and a mention of "Paul Mooney · 68 Submissions".

Fig 15: CORD-19 challenge

Research in Clinical Informatics

As an example of research in clinical informatics, Dr. Ahmad Baghal has created the UAMS Cohort Estimation Tool which is a web-based front end to discover patient cohorts during the early stages of research to obtain counts of patients based on inclusion and/or exclusion criteria.



I2B2: *Informatics for Integrating Biology & the Bedside*



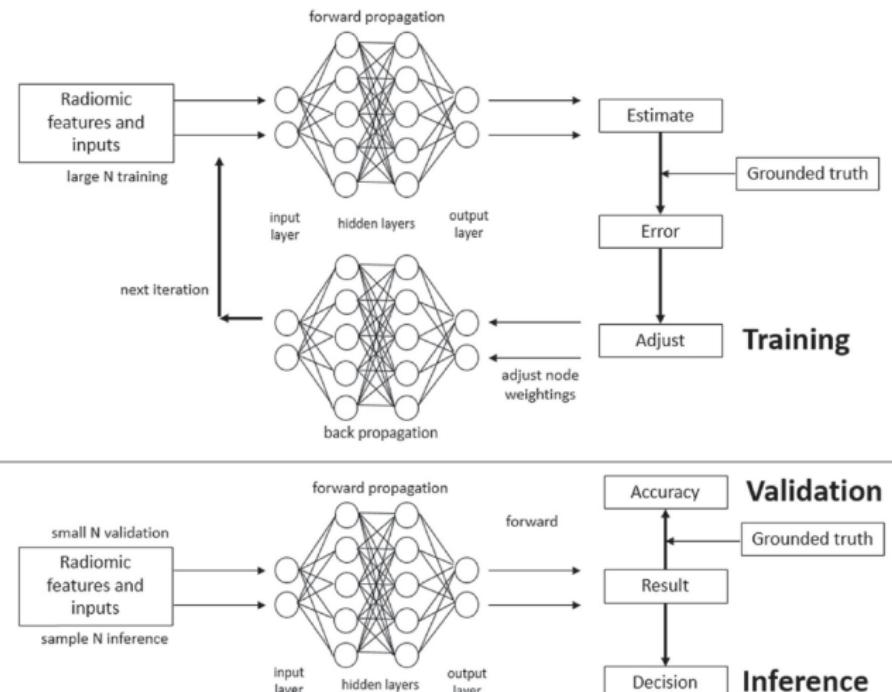
Cancer Imaging Archive

UAMS hosts the Cancer Imaging Archive



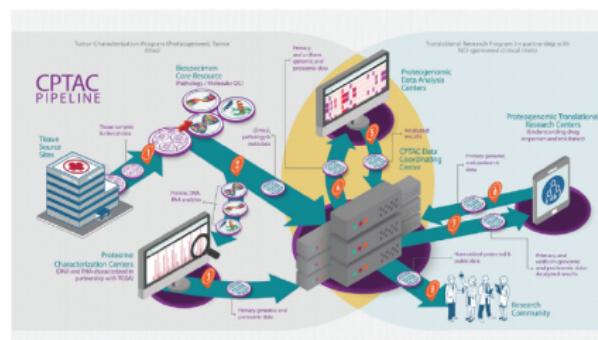
Imaging Informatics

In medical imaging, AI aims to enhance outcomes, quality and efficiency while maintaining ethical and regulatory requirements



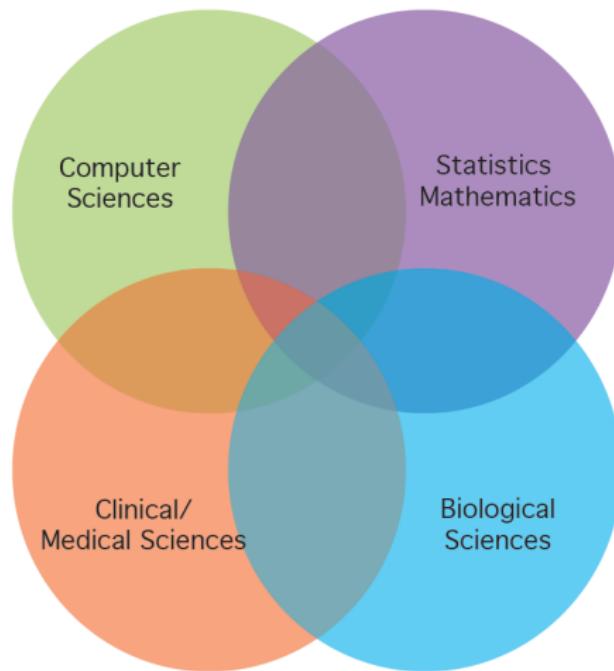
Multiomics for cancer research

National Cancer Institute's Clinical Proteomic Tumor Analysis Consortium (CPTAC) has collected data (genomics and proteomics) and recently linked to the imaging data



What skills do I need?

BMI is highly interdisciplinary



What computer language do I need?

It really depends on the area that you want to specialize and your background. At DBMI, all our students must acquire certain competency in Python



Computer/Scripting Languages

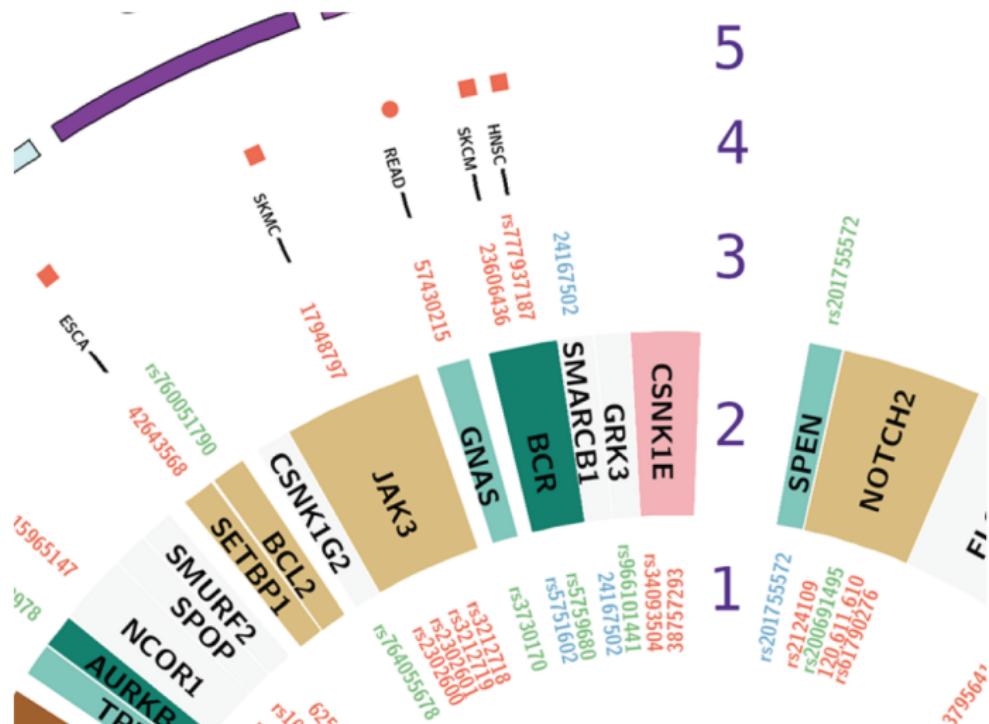
What do I do in DBMI?

I have research projects in:

- Imaging analysis
- Deep Learning in imaging
- Cancer genomics/imaging data integration
- Limb Development and Bone homeostasis
- Genomics of rare diseases
- Statistical Machine Learning

Research Collaborations

A Bioinformatician is a key research partner in molecular biology



DBMI graduate program

- Graduate Certificate -15 credits
- Master of Science - 36 credits
- Doctor of Philosophy - 55 credits (minimum)

**DEADLINE for Fall 2022 is April 1st
Flexible & Distance Learning Options Available for Working Professionals!**

Contact:

Crystin Mullins cmullins@uams.edu

Horacio Gomez-Acevedo gomezacevedohoracio@uams.edu



Figure credits

Some figures come from James et al. An introduction to statistical learning, and G. Rojas, Neural Networks, both from Springer Verlag.

- Figure 1. Science Photo Library
- Figure 2. Nature Education
- Figure 3. National Human Genome Research Institute
- Figure 4. Integrating Genomics into Healthcare: A Global Responsibility
- Figure 6 Big Data-The power of petabytes
- Figure 7 Exponential scaling of single-cell RNA-seq in the past decade
- Figure 8. Genomics and data science: an application within an umbrella
- Figure 11 Mapping single molecule sequencing reads using basic local alignment with successive refinement BLASR: application and theory
- Figure 12. A Framework for Variation Discovery and Genotyping Using Next-Generation DNA Sequencing Data

Figure credits

- Figure 13. Google Cloud Architecture
- Figure 14. ICGC Data Portal
- Figure 19. Machine Learning and Deep Learning in Medical Imaging: Intelligent Imaging
- Figure 20. A longitudinal analysis of data quality in a large pediatric data research network
- Figure 21. CPTAC