

BIOS 7747: Machine Learning for Biomedical Applications

Course presentation - Introduction to machine learning

Antonio R. Porras (antonio.porras@cuanschutz.edu)

Department of Biostatistics and Informatics
Colorado School of Public Health
University of Colorado Anschutz Medical Campus

Outline

- Introductions

- Course presentation



Course summary

- ❑ Credits: 2
- ❑ Audience: MS or PHD students in Biostatistics, Bioengineering, Computational Bioscience
- ❑ Prerequisites:
 - Biostatistical methods (e.g., BIOS 6611, BIOS 6612)
 - Linear algebra (e.g., MATH 3191)
 - Python programming (e.g., BIOS 6642)
- ❑ Classes: Tuesdays and Thursdays, 9:00-9:50AM
- ❑ Office hours: Tuesdays, 12:00-1:30pm. Building 500, W4132

Course summary

□ Materials:

- Reading requirements: no book required.
- Programming environment: Python 3
 - Note: non-native Python development environment problems will not be addressed
- Supporting materials:
 - Introduction to Machine Learning. Ethem Alpaydin. Third Edition. 2014. ISBN 0262028182.
 - Deep Learning. Ian Goodfellow and Yoshua Bengio and Aaron Courville. MIT Press. 2016.
<https://www.deeplearningbook.org/>.
 - Introduction to Statistical Learning by Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani. Springer, 2013. Corrected 8th printing, 2017. ISBN 1461471370.
 - Deep Learning with PyTorch: Build, train, and tune neural networks using Python tools 1st Edition. Eli Stevens, Luca Antiga, and Thomas Viehmann. Manning. 1617295264.

□ Evaluation:

- Assignments / homework: 20%
- Paper presentations 15%
- Classroom: 15%
- Final exam: 50%

Course summary

Date	Format	Topic
8/30	Lecture	Course introduction
9/1	Practical	Practical warmup class: Python setup and use of common libraries
9/6	Lecture	Supervised machine learning: regression, regular gradient descent optimization, linear and non-linear regression.
9/8	Practical	Regression and optimization with Python: Statsmodels, Scikit-learn and Scipy.
9/13	Lecture	Feature exploration, visualization, pre-processing and normalization.
9/15	Practical	Feature exploration and pre-processing for a non-linear regression problem.
9/20	Seminar	Invited seminar: Ethics and Bias in Machine Learning: More than the Data. Dr. Matthew DeCamp, MD, PhD
9/22		
9/27	Lecture	Supervised machine learning: classification and logistic regression. Performance evaluation and cross-validation.
9/29	Practical	Cross-validation of logistic regression-based classification methods.
10/4	Flipped classroom	Supervised machine learning: Naïve Bayes, K-nearest neighbors, decision trees and random forests (boosting, bootstrap and bagging)
10/6	Practical	K-nearest neighbors, decision trees and random forests in Python.
10/11	Flipped classroom	Supervised machine learning: Lagrange multipliers and support vector machines. The kernel trick. Platt's algorithm. Recursive feature elimination for feature selection. Support vector regression.
10/13	Practical	Support vector machines, class imbalance, Platt's algorithm, understanding and visualizing overfitting in Python.
10/18	Flipped classroom	Unsupervised learning: clustering, mixture models and other alternatives. Selecting the appropriate data for clustering. Performance evaluation.
10/20	Practical	Clustering and visualization in Python.

Course summary

Date	Format	Topic
10/25	Lecture	The curse of dimensionality and dimensionality reduction. Unsupervised dimensionality reduction using principal component analysis. Generalized principal component analysis and principal component analysis-based modeling. Supervised dimensionality reduction using linear discriminant analysis.
10/27	Practical	Implementation and visualization of principal component analysis and linear discriminant analysis in Python.
11/1	Student presentations	Presentations of feature-based machine learning research papers.
11/3	Student presentations	Presentations of feature-based machine learning research papers.
11/8	Lecture	Introduction to neural networks. Feed-forward networks, activation and backpropagation. Examples of biomedical applications.
11/10	Practical	Introduction to Neural Networks with Pytorch and Tensorboard in Python.
11/15	Lecture	Working with time series and images: convolutional neural networks. Examples and application to biomedical data.
11/17	Practical	Convolutional Neural Networks with Pytorch.
11/22	Lecture	Design alternatives in neural networks, examples in biomedical applications. Incorporating diverse data.
11/24		Break
11/29	Practical	Practical considerations when training a neural network. Understanding the effect of design modifications
12/1	Practical	Practical considerations when training a neural network. Understanding the effect of design modifications
12/6	Student presentations	Presentations of deep learning research papers.
12/8	Student presentations	Presentations of deep learning research papers.
12/13		Exams week
12/15		Exams week

Introduction to machine learning

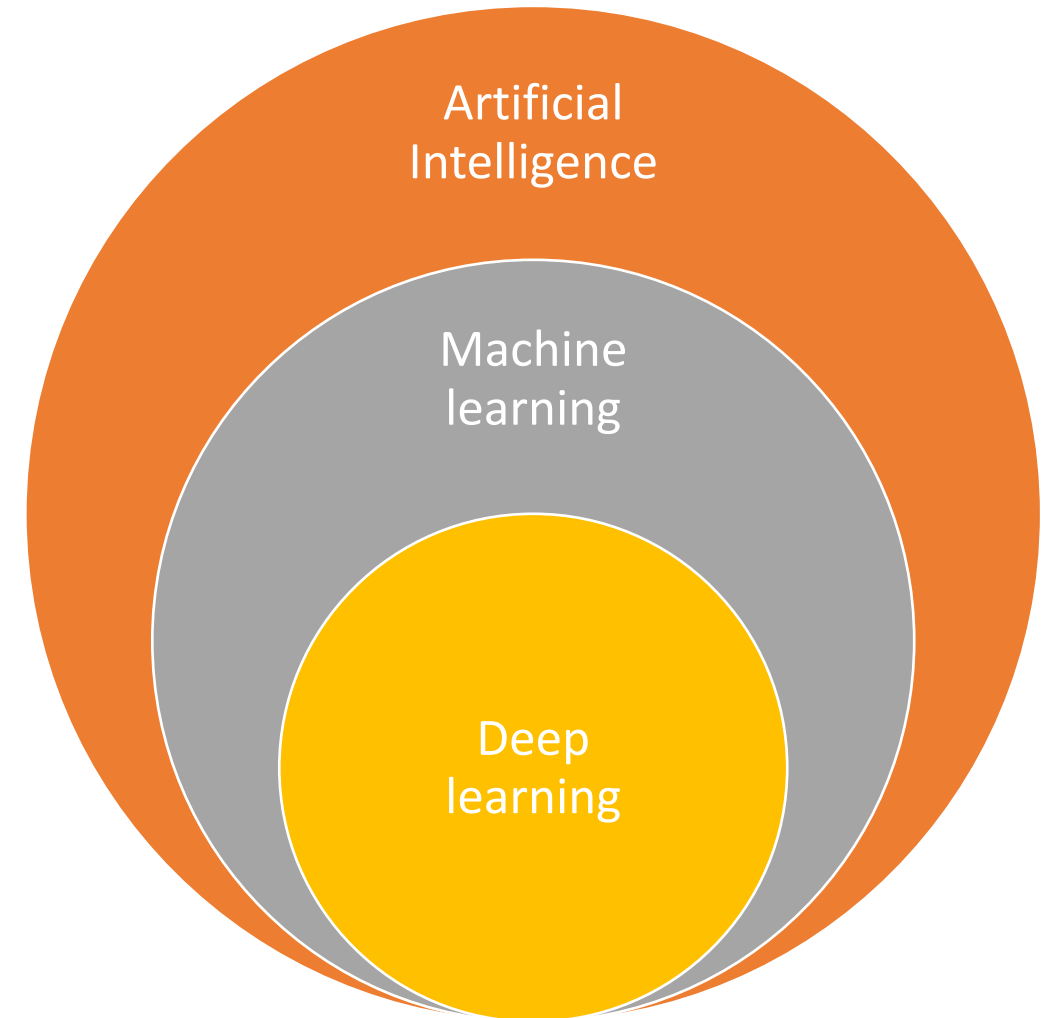
Intelligence: capability of inferring new information, retaining it as knowledge that can be applied within a context or environment

Human intelligence: capability of humans to reach correct conclusions about what is true and false, and to solve problems. It is marked by complex cognitive skills and high levels of motivation and self-awareness.

Artificial intelligence: Systems or machines that can mimic human intelligence to perform specific tasks that can iteratively improve themselves based on collected information.

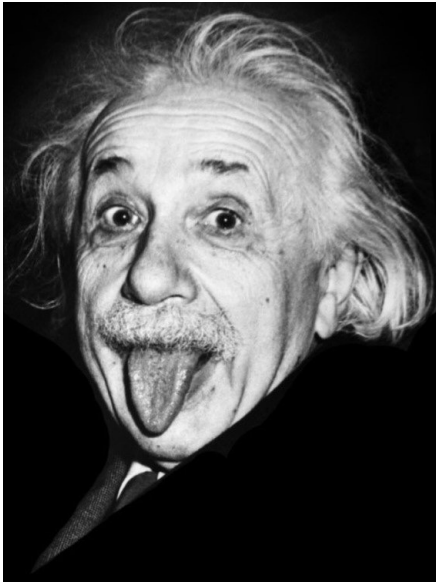
Machine learning: Branch of artificial intelligence and computer science that focuses on developing algorithms that imitate the way humans learn

Deep learning: Branch of machine learning that uses neural networks to leverage large amounts of data



Introduction to machine learning

Human intelligence



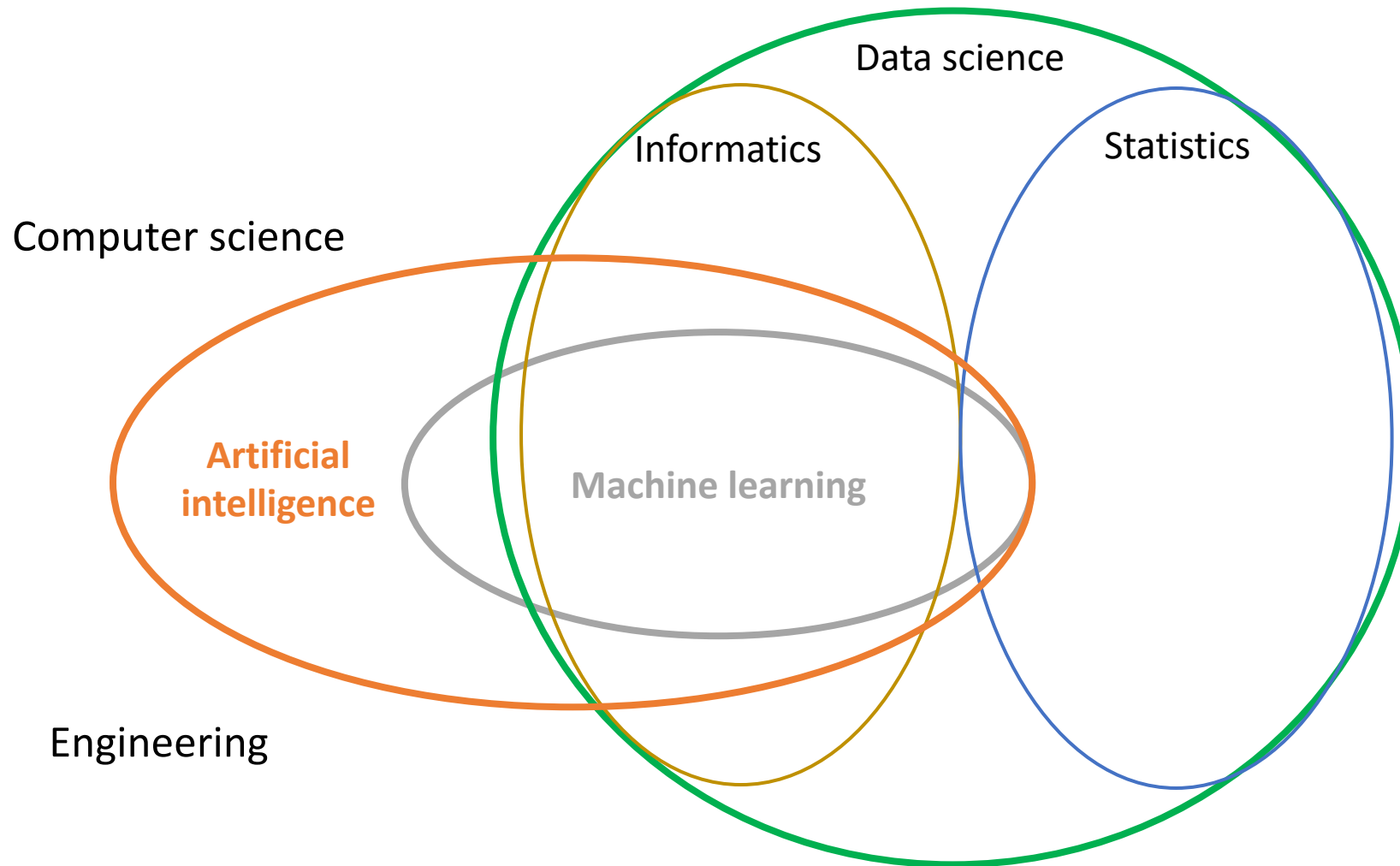
- Fast learning
- Can learn millions of highly complex tasks
- Creativity and originality
- Conscious
- Self-aware
- Power-efficient
- Influenced by emotions
- Inexact
- Slow
- Forgetful

Artificial intelligence

- Slow learning process
- Can learn a limited amount of simple tasks
- Highly limited creativity
- Unconscious
- Not self-aware
- Power-inefficient
- Repeatable
- Exact
- Fast
- Persistent data



Introduction to machine learning

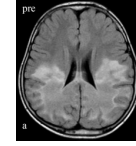
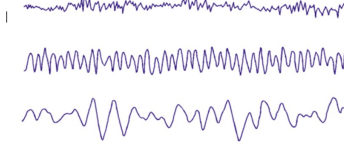


It's all math!

Introduction to machine learning for biomedical applications

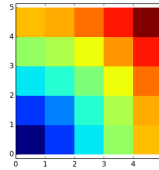
An overview of the machine learning approach in biomedicine

1. Data collection



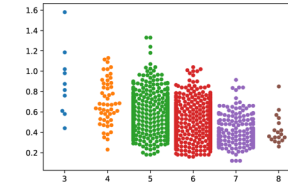
2. Data pre-processing

3. Data representation



10	20	30	40	50	60	70	80	90	100
----	----	----	----	----	----	----	----	----	-----

4. Data wrangling (and more pre-processing) and exploratory analysis



5. Feature selection and/or feature space transformation

6. Model construction

7. Model evaluation

8. Deployment

Machine learning?

Machine learning?

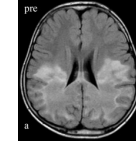
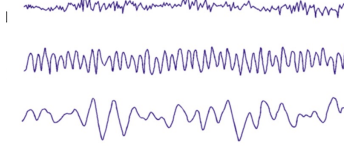
Introduction to machine learning for biomedical applications

An overview of the machine learning approach in biomedicine

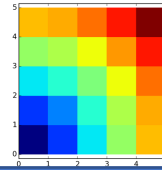
1. Data collection



2. Data pre-processing



3. Data representation



10	20	30	40	50	60	70	80	90	100
----	----	----	----	----	----	----	----	----	-----

4.

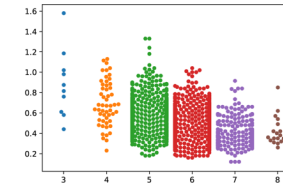
5.

6.

Deep learning

7. Model evaluation

8. Deployment



Machine learning?

Machine learning?

Introduction to machine learning

Why using machine learning in biomedical research?

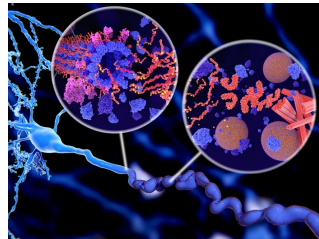
- Detailed analysis: machine learning methods can consider subtle quantitative variables that may be important to identify and/or predict the course of a disease or its treatment.
- Computational analysis: machine learning methods can consider large amounts of data and identify complex relationships between them to enable repetitive and reliable analysis.

Most common types of data

Omics



Genomics

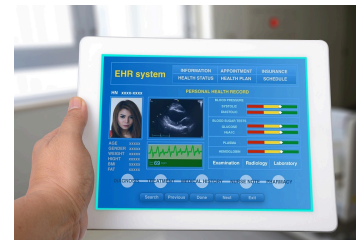


Proteomics

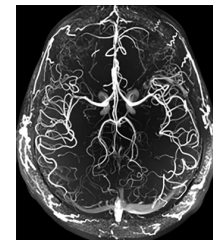


Microbiomics

Healthcare data



EHR



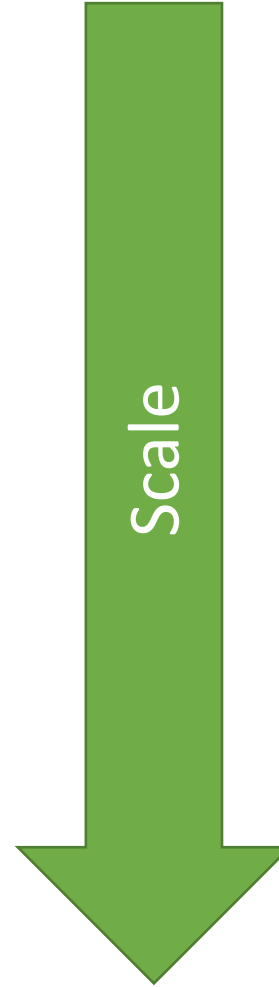
Imaging



Physiological
signals

Representing biomedical information

- ❑ Molecular information
- ❑ Cell information
- ❑ Tissue information
- ❑ Patient information
- ❑ Population information



Representing biomedical information

■ Molecular information

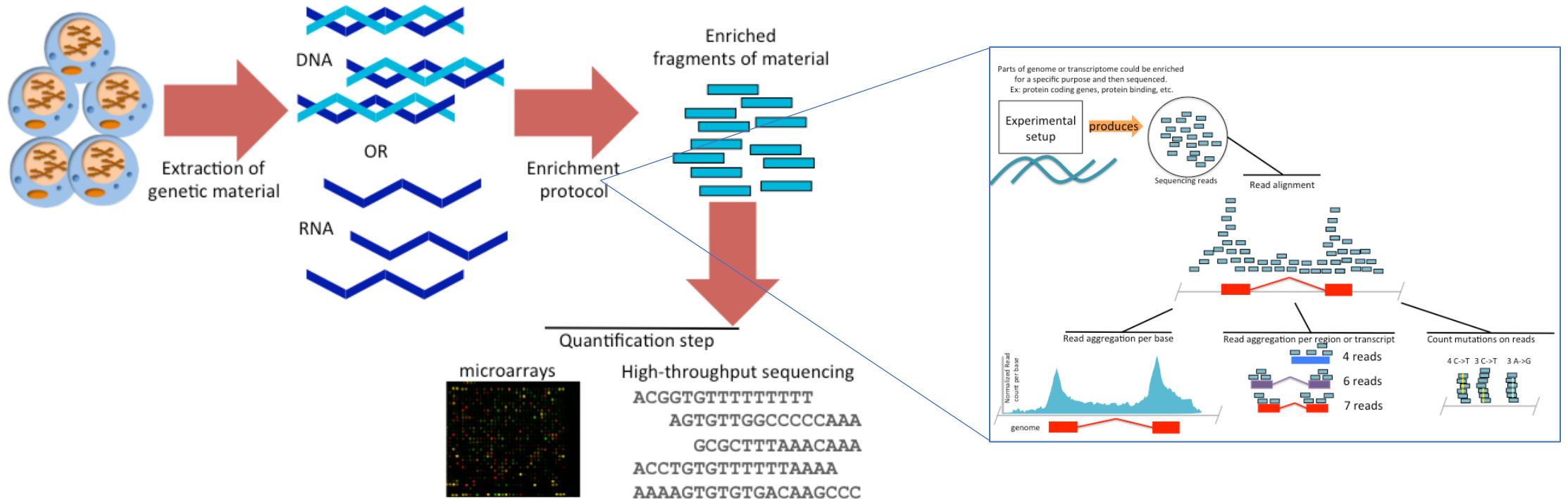
- Genomics: studies DNA molecules with two strands containing the genetic information coded as sequences of adenine, thymine, guanine and cytosine.
 - Structural: sequencing and mapping.
 - Functional: gene expression and their function – transcription, translation and interactions
- Transcriptomics: studies RNA transcripts produced by the genome and how they are affected by factors such as environment, drugs, etc.
- Proteomics: studies the structure and function of the proteome (set of all proteins).
- Epigenomics: studies the epigenetic modifications of genetic materials (e.g., DNA methylation, histone modification).
- Others: lipidomics, glycomics, metabolomics...

Representing biomedical information

□ Molecular information

Which one is the data?

Every step depends on previous one

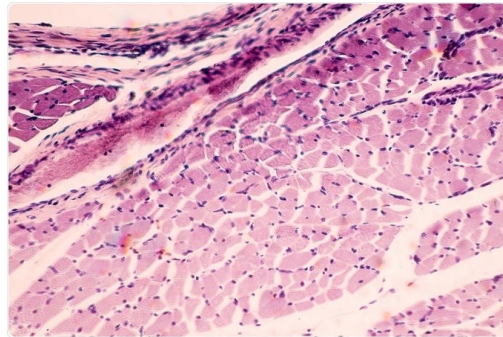


Representing biomedical information

□ Cell and tissue information:

- Highly driven by microscopy imaging

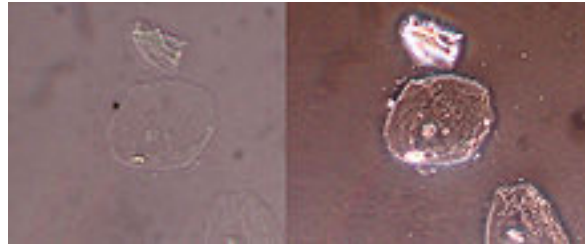
Optical
microscopy



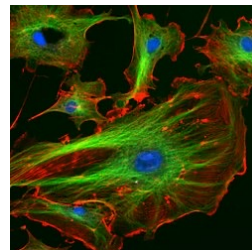
Scanning
electron
microscopy



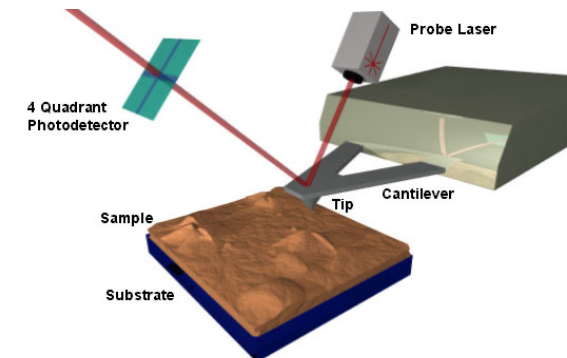
Phase contrast
microscopy



Fluorescent
microscopy



Atomic forces
microscopy



Representing biomedical information

□ Patient information:

- Anatomical
 - Imaging: computed tomography, magnetic resonance, etc.
- Functional:
 - Blood tests
 - Measured signals: electrocardiogram, electroencephalogram, electromyogram, etc.
- Other
 - Electronic health records: demographic, symptoms and history data

Image or signal data is not the same than image- and signal-derived data

□ Population data:

- Any of all previous information from many individuals
- Survey data

Next class

- ❑ Have a Python 3 installation
- ❑ Jupyter could be handy in block 1 (at your own risk)
- ❑ Install:
 - Pandas
 - Numpy
 - Scipy
 - Scikit-learn
 - Matplotlib
 - Statsmodels
 - Xlsxwriter
- ❑ “Play” with Numpy (data representation and matrix operations), Pandas (data I/O and representation) and Shelve (model and experiment persistent storage).