



Introduction to Machine Learning

Inas A. Yassine, PhD

Assoc. Prof. , Systems and Biomedical Engineering Department, Cairo University

inas.yassine@eng.cu.edu.eg



Outlines

- What is machine learning
 - By example
 - Definition
- Medical Informatics
 - Definition
 - Applications
 - Clinical decision making
 - Content based Medical Image Retrieval



What is Machine Learning

(By Example)



Classification

From Data to Discrete Classes

Spam Filtering

data

Osman Khan to Carlos [show details](#) Jan 7 (6 days ago) [Reply](#)

sounds good
+ok

Carlos Guestrin wrote:
Let's try to chat on Friday a little to coordinate and more on Sunday in person?

Carlos

Welcome to New Media Installation: Art that Learns

Carlos Guestrin to 10615-announce, Osman, Miche [show details](#) 3:15 PM (8 hours ago) [Reply](#)

Hi everyone,

Welcome to New Media Installation:Art that Learns

The class will start tomorrow.
Make sure you attend the first class, even if you are on the Wait List.
The classes are held in Doherty Hall C316, and will be Tue, Thu 01:30-4:20 PM.

By now, you should be subscribed to our course mailing list: 10615-announce@cs.cmu.edu.
You can contact the instructors by emailing: 10615-instructors@cs.cmu.edu

Natural _LoseWeight SuperFood Endorsed by Oprah Winfrey, Free Trial 1 bottle, pay only \$5.95 for shipping mfw rik [Spam](#) | [x](#)

Jaquelyn Halley to nherrlein, bcc: thehorney, bcc: anç [show details](#) 9:52 PM (1 hour ago) [Reply](#)

=== Natural WeightLOSS Solution ===

Vital Acai is a natural WeightLOSS product that Enables people to lose wieght and cleansing their bodies faster than most other products on the market.

Here are some of the benefits of Vital Acai that You might not be aware of. These benefits have helped people who have been using Vital Acai daily to Achieve goals and reach new heights in there dieting that they never thought they could.

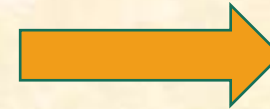
- Rapid WeightLOSS
- Increased metabolism - BurnFat & calories easily!
- Better Mood and Attitude
- More Self Confidence
- Cleanse and Detoxify Your Body
- Much More Energy
- BetterSexLife
- A Natural Colon Cleanse

Prediction

Spam

vs.

Not Spam



Face Recognition



Example training images
for each orientation



Weather Prediction





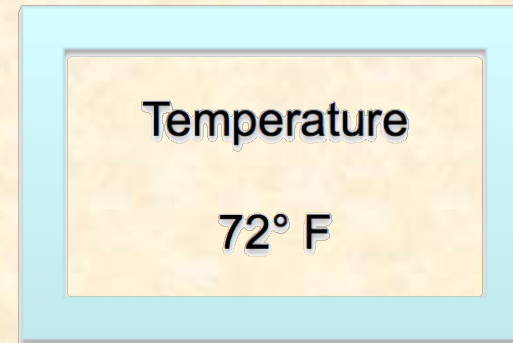
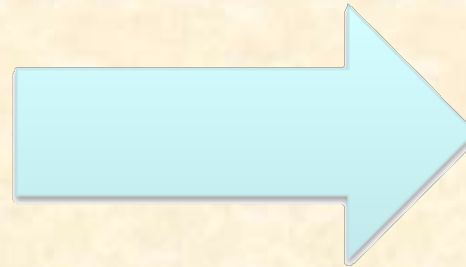
Regression

Predicting a numeric Value

Stock Market



Weather Prediction ...Revisited





Ranking










Comparing Items



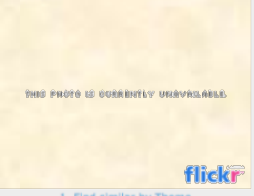


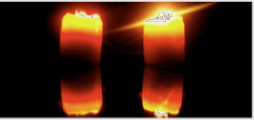
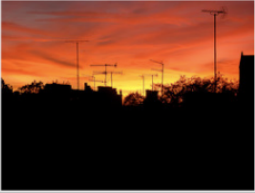
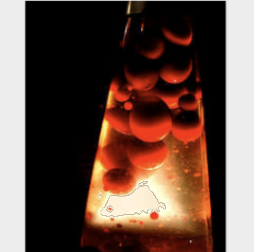

Web Search



The screenshot shows a Google search interface. The search bar contains the text "learning to rank". Below the search bar, a dropdown menu displays several suggestions: "learning to rank", "learning to rank for information retrieval", "learning to rank using gradient descent", and "learning to rank tutorial". To the right of the search bar is a blue search button with a magnifying glass icon. Below the search bar, the word "Search" is written in red. On the left side, there is a vertical navigation menu with options: "Web", "Images", "Maps", "Videos", "News", "Shopping", and "More". The "Web" option is selected. Below the menu, there is a location indicator for "Manhattan, NY 10012" and a "Change location" link. The main search results are displayed on the right. The first result is from Wikipedia, titled "Learning to rank - Wikipedia, the free encyclopedia", with a URL "en.wikipedia.org/wiki/Learning_to_rank". The second result is from Yahoo!, titled "Yahoo! Learning to Rank Challenge", with a URL "learningtorankchallenge.yahoo.com/". The third result is a PDF document titled "[PDF] Large Scale Learning to Rank" from Tufts University, with a URL "www.eecs.tufts.edu/~dsculley/papers/large-scale-rank.pdf". The fourth result is from Microsoft Research, titled "Microsoft Learning to Rank Datasets - Microsoft Research", with a URL "research.microsoft.com/en-us/projects/mslr/". The fifth result is also from Microsoft Research, titled "LETOR: A Benchmark Collection for Research on Learning to Rank ...", with a URL "research.microsoft.com/~letor/".

Given Image, Find Similar Images

 <p>1. Search mode: Theme OR</p> <p>2. Find similar by Color / Texture</p>	 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>	 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>
 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>	 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>	 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>
 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>	 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>	 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>

 <p>1. Find similar by Theme OR</p> <p>2. Search mode: Color / Texture</p>	 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>	 <p>THIS PHOTO IS CURRENTLY UNAVAILABLE. flickr</p> <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>
 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>	 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>	 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>
 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>	 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>	 <p>1. Find similar by Theme OR</p> <p>2. Find similar by Color / Texture</p>



Collaborative Filtering

Recommendation Systems



amazon Try Prime David's Amazon.com Today's Deals Gift Cards Sell Help Daily Lightning Deals Back-to-School Savings Shop now

Shop by Department Books Hello, David Your Account

Your Amazon.com Your Browsing History Recommended For You Amazon Betterizer Improve Your Recommendations Your Profile Learn More

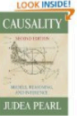
Your Amazon.com > Recommended For You > Books > Subjects > Science & Math > History & Philosophy

These recommendations are based on [items you own](#) and more.

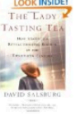
view: All | [New Releases](#) | [Coming Soon](#)

Just For Today
[Browse Recommended](#)


Recommendations
History & Philosophy
[History of Science](#)
[Philosophy of Biology](#)
[Philosophy of Medicine](#)

- 

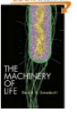
Causality: Models, Reasoning and Inference
by Judea Pearl (September 14, 2009)
Average Customer Review: ★★★★★ (19)
In Stock
List Price: \$60.00
Price: \$32.49
[61 used & new from \\$28.00](#)

I own it Not interested ★★★★★ Rate this item
Recommended because you purchased [Probabilistic Graphical Models](#) and more (Fix this)
- 

The Lady Tasting Tea: How Statistics Revolutionized Science in the Twentieth Century
by David Salsburg (May 1, 2002)
Average Customer Review: ★★★★★ (78)
In Stock
List Price: \$48.99
Price: \$13.88
[81 used & new from \\$9.00](#)

I own it Not interested ★★★★★ Rate this item
Recommended because you added [The Theory That Would Not Die](#) to your Wish List (Fix this)
- 

The Eighth Day of Creation: Makers of the Revolution in Biology, 25th Anniversary Edition
by Horace Freeland Judson (November 1, 1996)
Average Customer Review: ★★★★★ (10)
In stock on September 4, 2013
List Price: \$66.00
Price: \$36.09
[59 used & new from \\$26.95](#)

I own it Not interested ★★★★★ Rate this item
Recommended because you purchased [Molecular Biology of the Cell](#) (Fix this)
- 

The Machinery of Life
by David S. Goodsell (April 28, 2009)
Average Customer Review: ★★★★★ (41)
In Stock
List Price: \$25.00
Price: \$17.49
[92 used & new from \\$12.00](#)

Recommendation Systems

Machine learning competition with a \$1 million prize

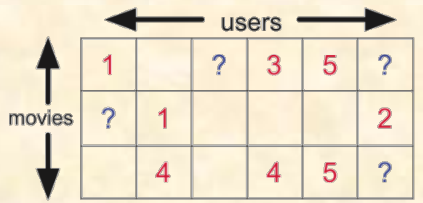
Leaderboard

Display top 20 leaders.

Rank	Team Name	Best Score	Improvement	Last Submit Time
1	The Ensemble	0.8553	10.10	2009-07-26 18:38:22
2	BellKor's Pragmatic Chaos	0.8554	10.09	2009-07-26 18:18:28
Grand Prize - RMSE <= 0.8563				
3	Grand Prize Team	0.8571	9.91	2009-07-24 13:07:49
4	Opera Solutions and Vandalis United	0.8573	9.89	2009-07-25 20:08:52
5	Vandalis Industries I	0.8579	9.83	2009-07-26 02:49:53
6	PanamaInTheatre	0.8582	9.80	2009-07-12 15:09:53
7	BellKor in BigChaos	0.8590	9.71	2009-07-26 12:57:25
8	Dane	0.8603	9.58	2009-07-24 17:18:43
9	Opera Solutions	0.8611	9.49	2009-07-26 18:02:08
10	BellKor	0.8612	9.48	2009-07-26 17:19:11
11	BigChaos	0.8613	9.47	2009-06-23 23:06:52
12	Feedz	0.8613	9.47	2009-07-24 20:06:46
Progress Prize 2009 - RMSE = 0.8616 - Winning Team: BellKor in BigChaos				
13	Xianfeng	0.8633	9.26	2009-07-21 02:04:40
14	Gronk	0.8634	9.25	2009-07-26 15:58:34
15	Coc	0.8642	9.17	2009-07-25 17:42:38
16	Invisible Ideas	0.8644	9.14	2009-07-20 03:26:12
17	Just a qu in a game	0.8650	9.08	2009-07-22 14:10:42
18	Coca Comichael	0.8656	9.02	2009-07-26 16:00:54
19	J.Damian Su	0.8658	9.00	2009-03-11 09:41:54
20	samsill	0.8659	8.99	2009-04-16 06:29:35
Progress Prize 2007 - RMSE = 0.8712 - Winning Team: KorBell				
Closest to score on quiz subject - RMSE = 0.9514				



The image shows a screenshot of the Netflix website. The top navigation bar includes 'Browse', 'Recommendations', 'Friends', 'Queue', and 'Buy DVDs'. Below this, there are links for 'Home', 'Genres', 'New Releases', 'Previews', 'Netflix Top 100', and 'Critics'. The main content area is titled 'Movies For You' and features a recommendation for 'Bowling for Columbine' based on the user's interest in 'Carnivale: Season 1' and 'Eisenstein 9/11'. A sidebar on the right shows a promotional offer: 'You really liked it... Now only for just \$5.99'. In the foreground, there are silhouettes of two people looking at a screen, with a background of green code snippets.

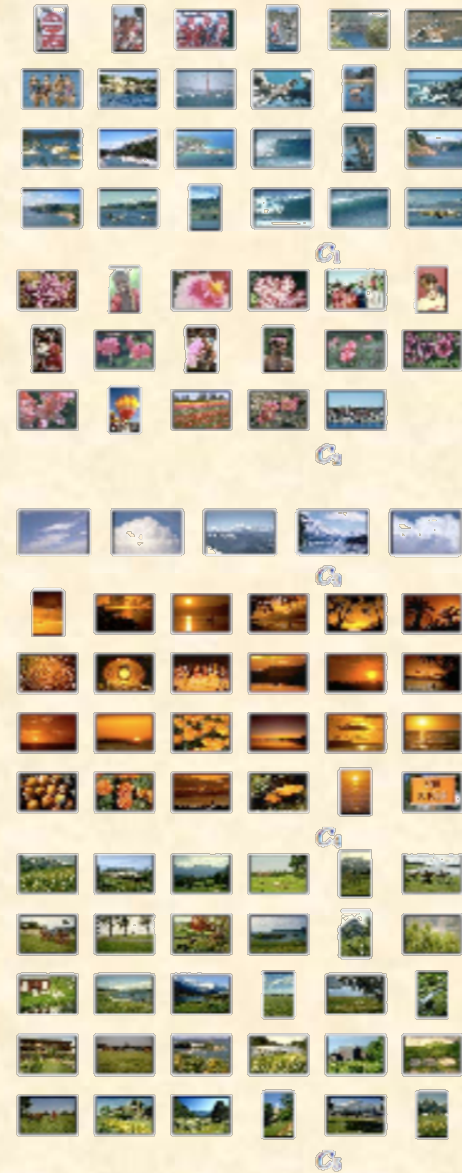
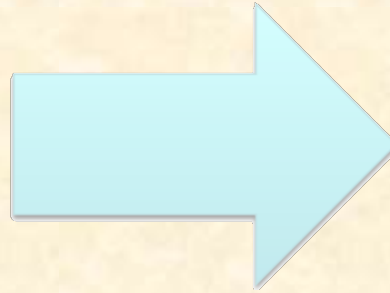
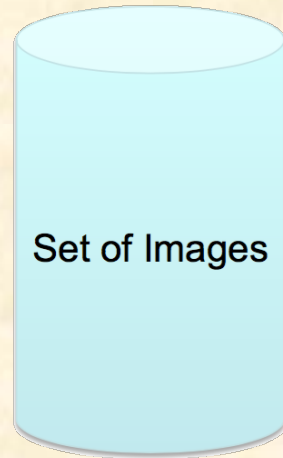




Clustering

Discovering Structures in Data

Clustering Images





Clustering Web Search Results

web news images wikipedia blogs jobs more »

Clusty

race Search advanced preferences

clusters sources sites

All Results (238) remix

- Car (28)
 - Race cars (7)
 - Photos, Races Scheduled (5)
 - Game (4)
 - Track (3)
 - Nascar (2)
 - Equipment And Safety (2)
 - Other Topics (7)
- Photos (22)
- Game (14)
- Definition (13)
- Team (18)
- Human (8)**
 - Classification Of Human (2)
 - Statement, Evolved (2)
 - Other Topics (4)
- Weekend (8)
- Ethnicity And Race (7)
- Race for the Cure (8)
- Race Information (8)

more | all clusters

find in clusters: Find

Cluster Human contains 8 documents.

Search Results

- [Race \(classification of human beings\) - Wikipedia, the free ...](#)
The term **race** or racial group usually refers to the concept of dividing **humans** into populations or groups on the basis of various sets of characteristics. The most widely used **human** racial categories are based on visible traits (especially skin color, cranial or facial features and hair texture), and self-identification. Conceptions of **race**, as well as specific ways of grouping **races**, vary by culture and over time, and are often controversial for scientific as well as social and political reasons. History · Modern debates · Political and ...
[en.wikipedia.org/wiki/Race_\(classification_of_human_beings\)](http://en.wikipedia.org/wiki/Race_(classification_of_human_beings)) - [cache] - Live, Ask
- [Race - Wikipedia, the free encyclopedia](#)
General. **Racing** competitions The **Race** (yachting **race**), or La course du millénaire, a no-rules round-the-world sailing event; **Race** (biology), classification of flora and fauna; **Race** (classification of **human** beings) **Race** and ethnicity in the United States Census, official definitions of "**race**" used by the US Census Bureau; **Race** and genetics, notion of racial classifications based on genetics. Historical definitions of **race**; **Race** (bearing), the inner and outer rings of a rolling-element bearing. **RACE** in molecular biology "Rapid ... General · Surnames · Television · Music · Literature · Video games
en.wikipedia.org/wiki/Race - [cache] - Live, Ask
- [Publications | Human Rights Watch](#)
The use of torture, unlawful rendition, secret prisons, unfair trials, ... Risks to Migrants, Refugees, and Asylum Seekers in Egypt and Israel ... In the run-up to the Beijing Olympics in August 2008, ...
www.hrw.org/background/usa/race - [cache] - Ask
- [Amazon.com: Race: The Reality Of Human Differences: Vincent Sarich ...](#)
Amazon.com: **Race: The Reality Of Human** Differences: Vincent Sarich, Frank Miele: Books ... From Publishers Weekly Sarich, a Berkeley emeritus anthropologist, and Miele, an editor ...
www.amazon.com/Race-Reality-Differences-Vincent-Sarich/dp/0813340861 - [cache] - Live
- [AAPA Statement on Biological Aspects of Race](#)
AAPA Statement on Biological Aspects of **Race** ... Published in the American Journal of Physical Anthropology, vol. 101, pp 569-570, 1996 ... PREAMBLE As scientists who study **human** evolution and variation, ...
www.physanth.org/positions/race.html - [cache] - Ask
- [race: Definition from Answers.com](#)
race n. A local geographic or global **human** population distinguished as a more or less distinct group by genetically transmitted physical
www.answers.com/topic/race-1 - [cache] - Live
- [Dopefish.com](#)
Site for newbies as well as experienced Dopefish followers, chronicling the birth of the Dopefish, its numerous appearances in several computer games, and its eventual take-over of the **human** **race**. Maintained by Mr. Dopefish himself, Joe Siegler of Apogee Software.
www.dopefish.com - [cache] - Open Directory



What is Machine Learning

(Definition)



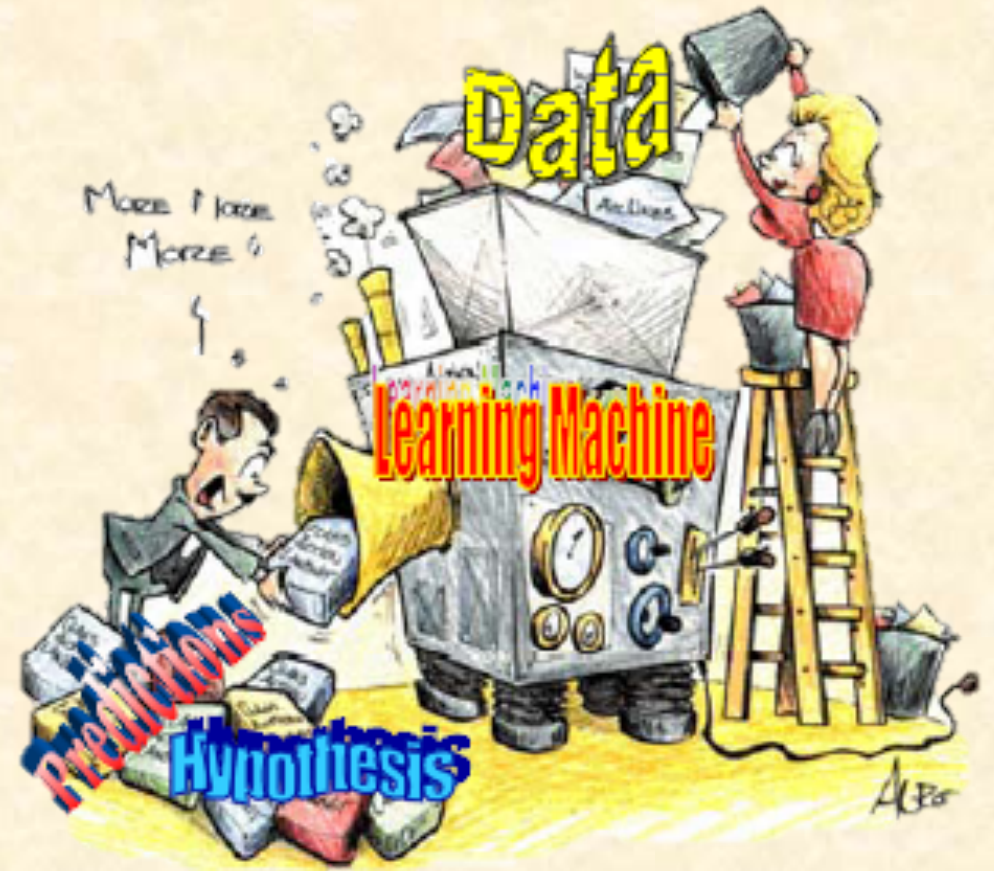
Why “Learn”?

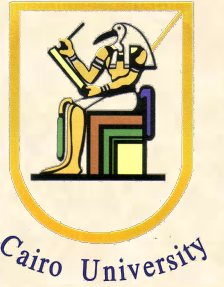
- Machine learning is programming computers to optimize a performance criterion using example data or past experience.
- There is no need to “learn” to calculate payroll
- It is about seeking a predictive and/or executable understanding of natural/
artificial subjects, phenomena, or activities from ...
- Learning is used when:
 - Human expertise does not exist (navigating on Mars)
 - Humans are unable to explain their expertise (speech recognition)
 - Solution changes in time (routing on a computer network)
 - Solution needs to be adapted to particular cases (user biometrics)



What is Machine Learning(ML)?

- Arthur Samuel: “It is a science that gives the computer the ability to learn without the need to write a program“.





What We Talk About When We Talk About “Learning”

- Learning general models from a data of particular examples
- Data is cheap and abundant (data warehouses, data marts); knowledge is expensive and scarce.
- Example in retail: Customer transactions to consumer behavior:

People who bought “Da Vinci Code” also bought “The Five People You Meet in Heaven” (www.amazon.com)

- Build a model that is *a good and useful approximation* to the data.



Data Mining

Is the non-trivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data

Applications:

- Retail: Market basket analysis, Customer relationship management (CRM)
- Finance: Credit scoring, fraud detection
- Manufacturing: Optimization, troubleshooting
- Medicine: Medical diagnosis
- Telecommunications: Quality of service optimization
- Bioinformatics: Motifs, alignment
- Web mining: Search engines
- ...

Supervised Learning : find f

- Given: Training set $\{(x_i, y_i) \mid i = 1 \dots N\}$
- Find: A good approximation to $f : X \rightarrow Y$

Examples: what are X and Y ?

- Spam Detection
 - Map email to {Spam, Not Spam}
- Digit recognition
 - Map pixels to $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$
- Stock Prediction
 - Map new, historic prices, etc. to \mathbb{R} (the real numbers)

A supervised Learning Problem

Dataset:

Example	x_1	x_2	x_3	x_4	y
1	0	0	1	0	0
2	0	1	0	0	0
3	0	0	1	1	1
4	1	0	0	1	1
5	0	1	1	0	0
6	1	1	0	0	0
7	0	1	0	1	0

- Our goal is to find a function $f : X \rightarrow Y$
 - $X = \{0, 1\}^4$
 - $Y = \{0, 1\}$
- Question1: How should we pick the *hypothesis space*, the set of possible functions f ?
- Question2: How do we find the best f in the hypothesis space?

Most General Hypothesis Space

- Consider all possible boolean functions over four input features!
- 2^{16} possible hypotheses
- 2^9 are consistent with our dataset
- How do we choose the best one?

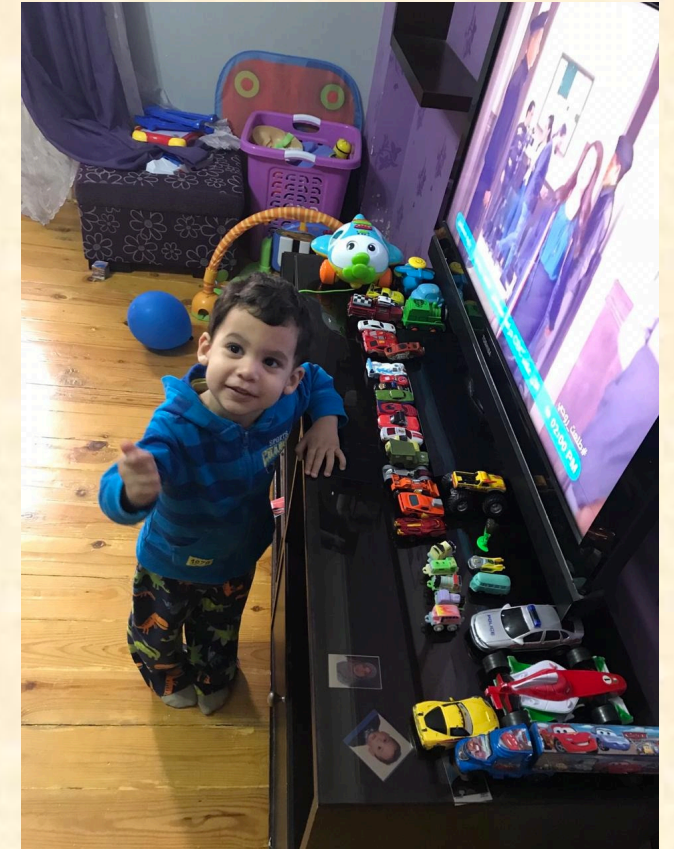
x_1	x_2	x_3	x_4	y
0	0	0	0	?
0	0	0	1	?
0	0	1	0	0
0	0	1	1	1
0	1	0	0	0
0	1	0	1	0
0	1	1	0	0
0	1	1	1	?
1	0	0	0	?
1	0	0	1	1
1	0	1	0	?
1	0	1	1	?
1	1	0	0	0
1	1	0	1	?
1	1	1	0	?
1	1	1	1	?

Dataset:

Example	x_1	x_2	x_3	x_4	y
1	0	0	1	0	0
2	0	1	0	0	0
3	0	0	1	1	1
4	1	0	0	1	1
5	0	1	1	0	0
6	1	1	0	0	0
7	0	1	0	1	0

UnSupervised Learning : find *groups*

- Given: Training set $\{(x_i) \mid i = 1 \dots N\}$
- Find: Learning “what normally happens”
→ No output
- Clustering: Grouping similar instances
 - Customer segmentation in CRM
 - Image compression: Color quantization
 - Bioinformatics: Learning motifs





Reinforcement Learning

- Topics:
 - Policies: what actions (Reward, punishment) should an agent take in a particular situation
 - Utility estimation: how good is a state (→ used by policy)
- No supervised output but delayed reward
- Credit assignment problem (what was responsible for the outcome)
- Applications:
 - Game playing
 - Robot in a maze
 - Multiple agents, partial observability, ...



Medical Informatics

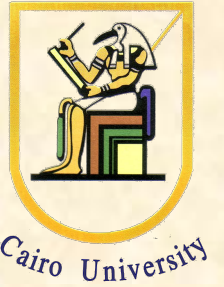


Medical informatics

Medical informatics is the application of computers, communications and information technology and systems to all fields of medicine *

- medical care
- medical education
- medical research.

**MF Collen, MEDINFO '80, Tokyo*



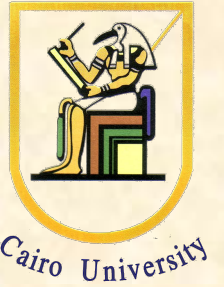
What is medical informatics?

Medical informatics is the branch of science concerned with the use of computers and communication technology **to acquire, store, analyze, communicate, and display medical information and knowledge** to facilitate understanding and **improve the accuracy, timeliness, and reliability of decision-making.***

**Warner, Sorenson and Bouhaddou, Knowledge Engineering in Health Informatics, 1997*

Clinical decision making

- Making sound clinical decisions requires:
 - right information, right time, right format
- Clinicians face a surplus of information
 - *ambiguous, incomplete, or poorly organized*
- **Rising tide of information**
 - Expanding knowledge sources
 - 40K new biomedical articles per month
 - Publicly accessible online health info
 - Hundreds of pictures per scan for one patient



Clinical decision making: What is the problem?

- Man is an imperfect data processor
 - We are sensitive to the **quantity and organization** of information
- Army officers and pilots commit ‘fatal errors’ when given too many, too few, or poorly organized data
- The same is true for clinicians who ‘watch’ for events
- Clinicians are particularly susceptible to errors of **omission**

Clinical decision making: What is the problem?

- Humans are “non-perfectable” data processors
 - Better performance requires more time to process
 - Irony
 - Clinicians increasingly face productivity expectations
 - Clinicians face increasing administrative tasks

What is medical imaging (MI)?

The study of **medical imaging** is concerned with the interaction of all forms of radiation with tissue

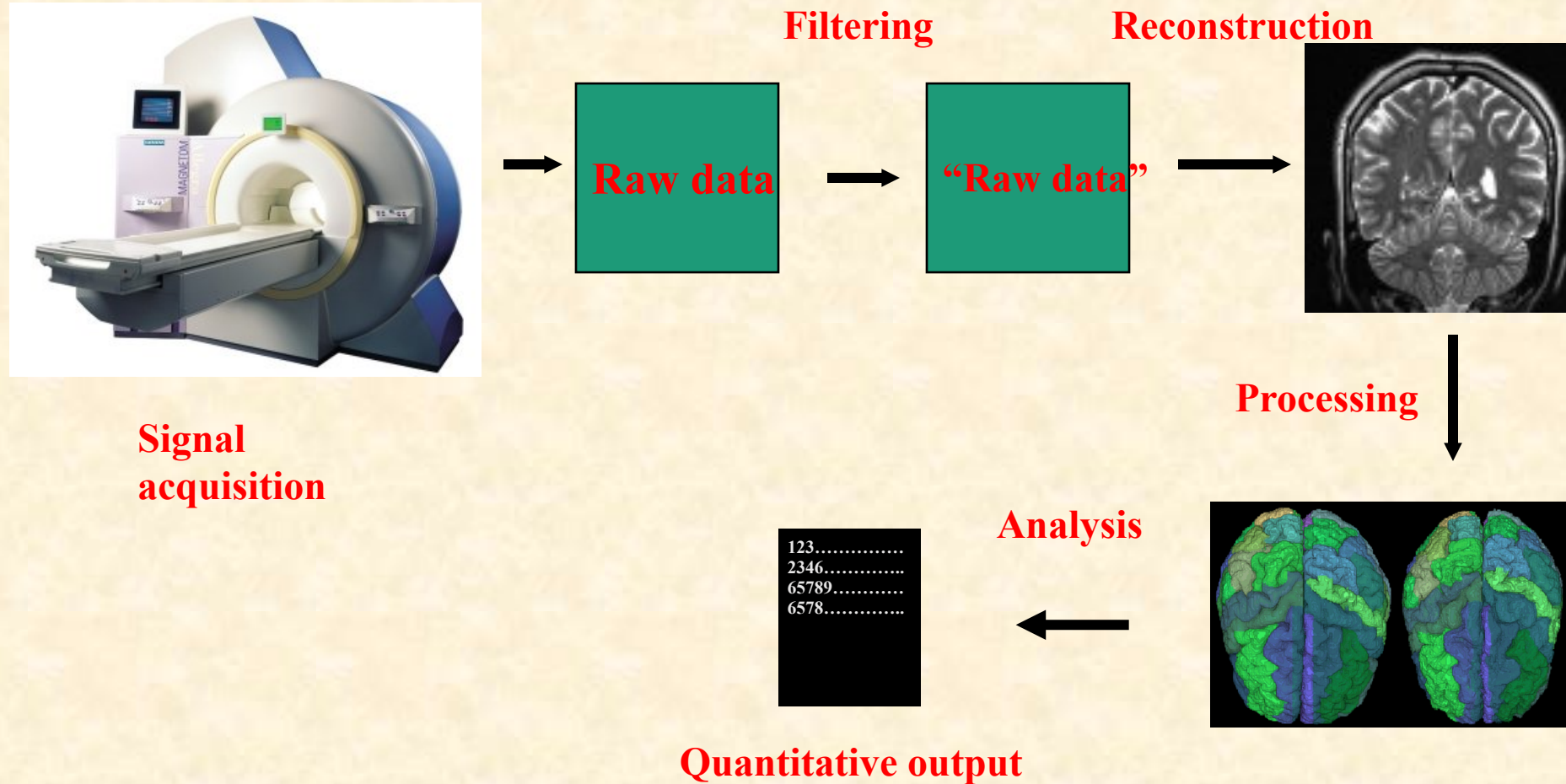
and

the development of appropriate technology to extract clinically useful information (usually displayed in an image format) from observation of this technology.

Sources of Images:

- **Structural/anatomical information** (CT, MRI, US) - within each elemental volume, tissue-differentiating properties are measured.
- **Information about function** (PET, SPECT, fMRI).

The imaging “chain”



Computer-Aided Diagnosis

- **Computed Aided Diagnosis (CAD)** is diagnosis made by a radiologist *when the output of computerized image analysis methods has been **incorporated** into his or her medical decision-making process.*
- CAD may be interpreted broadly to incorporate both
 - **the detection of the abnormality task** and
 - **the classification task**: likelihood that the abnormality represents a malignancy

Motivation for CAD systems

The amount of image data acquired during a CT scan is becoming overwhelming for human vision and the overload of image data for interpretation may result in oversight errors.

Computed Aided Diagnosis for:

- **Breast Cancer**

- Mammography Images

- **Lung Cancer**

- A thoracic CT scan generates about 240 section images for radiologists to interpret.



Content Based Medical Image Retrieval

Content-based medical image retrieval (CBMS) systems



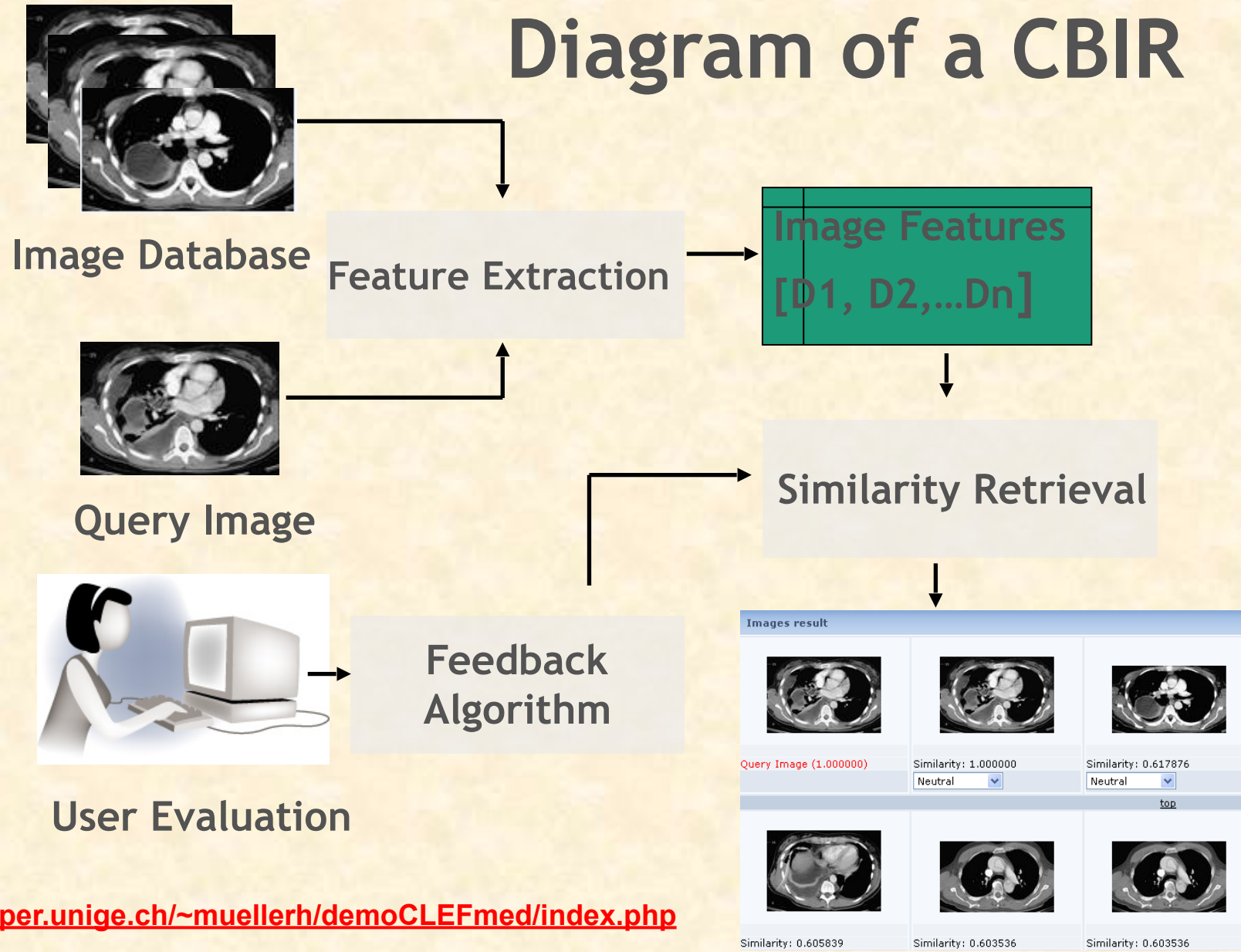
Definition of Content-based Image Retrieval:

Content-based image retrieval is a technique for retrieving images on the basis of **automatically** derived image features such as texture and shape.

Applications of Content-based Image Retrieval:

- **Teaching**
- **Research**
- **Diagnosis**
- **PACS and Electronic Patient Records**

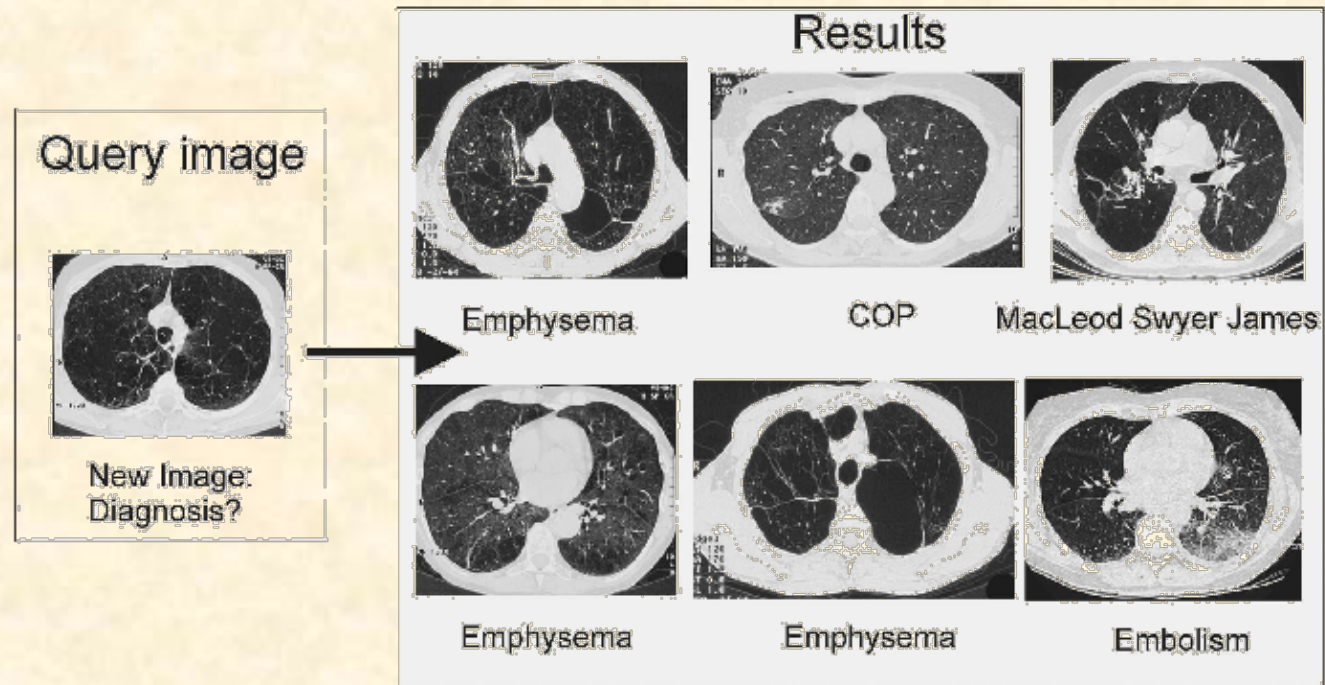
Diagram of a CBIR



<http://vipr.unige.ch/~muellerh/demoCLEFmed/index.php>

CBIR as a Diagnosis Aid

An image retrieval system can help when the diagnosis depends strongly on direct visual properties of images in the context of evidence-based medicine or case-based reasoning





Cairo University

CBIR as a Teaching Tool

An image retrieval system will allow students/teachers to browse available data themselves in an easy and straightforward fashion by clicking on “show me similar images”.

Advantages:

- stimulate self-learning and a comparison of similar cases
 - find optimal cases for teaching
- Teaching files:

- Casimage: <http://www.casimage.com>
- myPACS: <http://www.mypacs.net>

MyPACS.net
Teaching File Management System

Home About Support Solutions Cases
Welcome • Log in • Sign-up

MyPACS Log In
Login name:
Password:
Sign in
Sign up for your free account!

Create Your Own Teaching Files
Join radiologists from 400 institutions in 75 countries who are using MyPACS to create their cases online. Users have contributed **5000 radiology teaching files containing 20,000 images**, and new cases are being added every day. All you need is a free account and your web browser, and you can start creating cases right now.

Case of the Week
APPENDICITIS
by James Patrick, Resident, Medical College of Ohio, USA

This is a free service to the international community, funded in part by the National Institutes of Health. We also offer custom enterprise teaching file solutions.

Shared Cases		
Cranium and Contents	Heart	Vascular/Lymphatic
Face and Neck	Lung	Breast
Spine and Peripheral (GI)	Gastrointestinal	Other
Nervous System	(GI)	
Skeletal System	Genitourinary (GU)	

Highlighted Features:

CBIR as a Research Tool

Image retrieval systems can be used:

- to **complement text-based retrieval methods**
- for **visual knowledge management** where by the images and associated textual data can be analyzed together
 - multimedia data mining can be applied to **learn the unknown links between visual features and diagnosis or other patient information**
- for **quality control** to find images that might have been misclassified

In Summary ...

- Machine Learning is a Knowledge Discovery Process
- Machine learning is preferred approach to
 - Speech recognition, Natural language processing - Computer vision
 - Medical outcomes analysis
 - Robot control
 - Computational biology
 - Sensor networks
 - ...
- This trend is accelerating due to
 - Big data
 - Improved machine learning algorithms - Faster computers
 - Good open-source software
- Data Mining Tasks
 - classification, clustering, regression , ...



Course Content

- Introduction: Machine Learning for health informatics
 - Linear prediction: Regression
 - Learning strategies for basic classification algorithms: Linear classification, logistic regression
 - Probabilistic models (Bayes Theory, Decision trees)
 - Learning strategies for more advanced classification algorithms (SVM, ...)
 - Dimensionality Reduction (Principal component analysis)
 - Subspace Clustering
 - Neural Networks
 - Reinforcement Learning
 - Recommender Systems
 - Overfitting, underfitting, and the generalization gap
 - Machine Learning Algorithms testing and evaluation
-

Grading Criteria



Item	Course Grade Contribution
Final Exam	70%
Final Project	8%
Assignments/ Homework	15%
Midterm Exam	7%

Cheating Description

■ What is cheating?

- Sharing code: by copying, retyping, **looking at**, or supplying a file
- Describing: verbal description of code from one person to another.
- Coaching: helping your friend to write a line by line code
- Searching the Web for solutions
- Copying code from a previous course or online solution
- You are only allowed to use code we supply

■ What is NOT cheating?

- Explaining how to use systems or tools
- Helping others with high-level design issues

■ Penalty for cheating

- Zero in Assignments/ Exam
-



In Class/ Assignment General Code

- Electronic communications: **forbidden**
 - No email, instant messaging, cell phone calls, etc
- Presence in lectures: highly recommended
- No recordings of ANY KIND
- No side talks

.



Thank You ...