

Course Logistics and Introduction to Machine Learning

Piyush Rai

Introduction to Machine Learning (CS771A)

July 31, 2018



Course Logistics

- **Course name:** Introduction to Machine Learning or “ML” (CS771A)
- **Timing and Venue:** Tue/Thur 6:00-7:30pm, L-16
- **Course website:** <https://tinyurl.com/cs771-a18w> (slides/readings etc will be posted here)
- **Piazza discussion site:** <https://tinyurl.com/cs771-a18p> (use it actively and **responsibly**)
- **Gradescope (for assignment submission):** <https://tinyurl.com/cs771-a18g>
- Course-related announcements will be sent on the **class mailing list** (and also on Piazza)
- **Instructor:** Piyush Rai (e-mail: piyush@cse.iitk.ac.in, office: RM-502)
 - Prefix email subject by **CS771A** (better alternative: Piazza private message to instructor)
 - Office Hours: Wed 6:00-7:30pm (by appointment)
- **Auditing:** Don't need formal permission from me. Send me email to be added to the mailing list.
 - Will have access to all the course material; can participate in Piazza discussions
 - However, we are unable to grade your assignments/exams. Can't form project groups with creditors.



The TA Team



Shivam Bansal



Dhanajit Brahma



Sunabha Chatterjee



Prerit Garg



Gopichand Kotana



Neeraj Kumar



Pawan Kumar



Kranti Parida



Kawal Preet



Prem Raj



Utsav Singh



Samik Some



Vinay Verma



Project Mentors



Homanga Bhardwaj



Aadil Hayat



Ankit Jalan



Varun Khare



Sarthak Mittal



Gurpreet Singh

.. and some more..

Assignments, Exams, and Grading Policy

- Homework (4-5): 30%, Midsem Exam: 20%, Endsem Exam: 30%, Term Project: 20%
- Homeworks will usually be a mix of
 - **Pen-paper based questions:** Derivations/analysis/improvements of ML algos studied in class, designing “new” ML algos for some problem scenarios (using techniques studied in class)
 - **Programming questions:** Implement ML algos from scratch or using existing software tools, improve existing ML algos, apply ML algos to analyze data. Can use Python/MATLAB
- Homework solutions must be prepared using LaTeX. We will provide a LaTeX template.
 - Many resources to learn LaTeX available online. Must have skill. Learn it now!
- Homeworks must be submitted via Gradescope (no email submissions)
- Late homeworks will receive 10%/20%/30% penalty for 24/48/72 hour late submission, resp.
- Exams will be closed-book (an A4-sized cheat-sheet allowed)

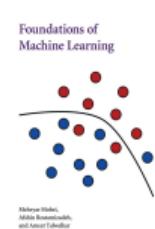
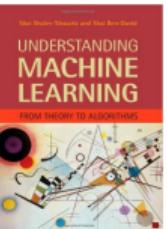
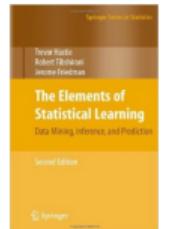
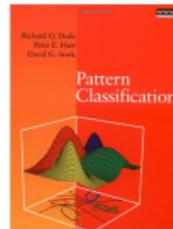
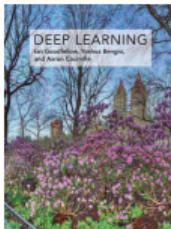
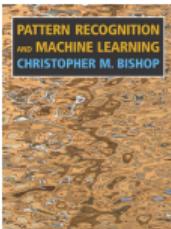
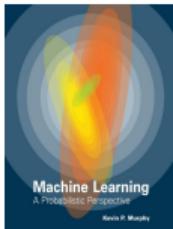
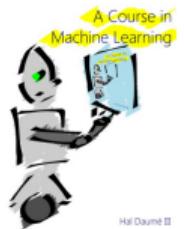


Term Project (Worth 20%)

- To be done in groups of 5-6 students. Must form the groups by August 18 and inform us
- Avoid smaller/larger groups (unless there is a very strong reason for that)
- We will float a list containing several project ideas that you can choose from
- You may also propose your own project idea
 - ... but please discuss the scope/feasibility with me and/or project mentors
- Many types of projects possible (theoretical/applied/mixed); even building a cool/useful app/portal for something using ML as its backend would be a nice project. Explore around for ideas.
- **Important:** Don't (re)use a project from another course you've done before (or doing currently)
- Need to submit a formal project proposal by Sept 7 (not graded but mandatory)
 - Should contain a problem description, tentative plan of action, etc (I'll provide more guidelines later)
- Please don't wait until Sept 7 to finalize the project idea

Textbook and References

- Many excellent texts but none “required”. Some of them include (list not exhaustive)



- Different books might vary in terms of
 - Set of topics covered
 - General approach taken e.g., classical statistics, deep learning, probabilistic/Bayesian, theory
 - **Terminology and notation (beware of this especially)**
- Avoid using too many sources until you have developed a reasonable understanding of a concept
- We will provide you the reading material from the relevant sources

Collaboration vs Cheating

- Collaboration is encouraged. Cheating/copying will lead to strict punishments.
- Feel free to discuss homework assignments with your classmates.
- However, your own solution must be in your own words (same goes for coding assignments)
- Plagiarism from other sources (for assignments/project) will also lead to strict punishment unless you duly credit the original source(s)
- Other things that will lead to punishment
 - Use of unfair means in the exams
 - Fabricating experimental results in assignments/project
- Important: Both copying as well as helping someone copy will be equally punishable

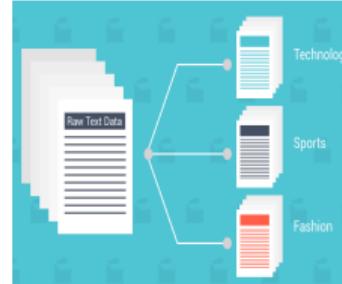
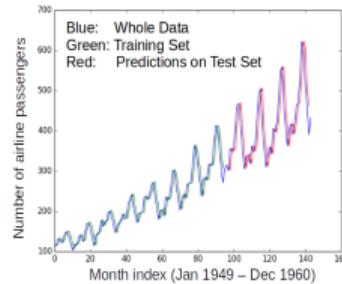


Intro to Machine Learning



Machine Learning (ML)

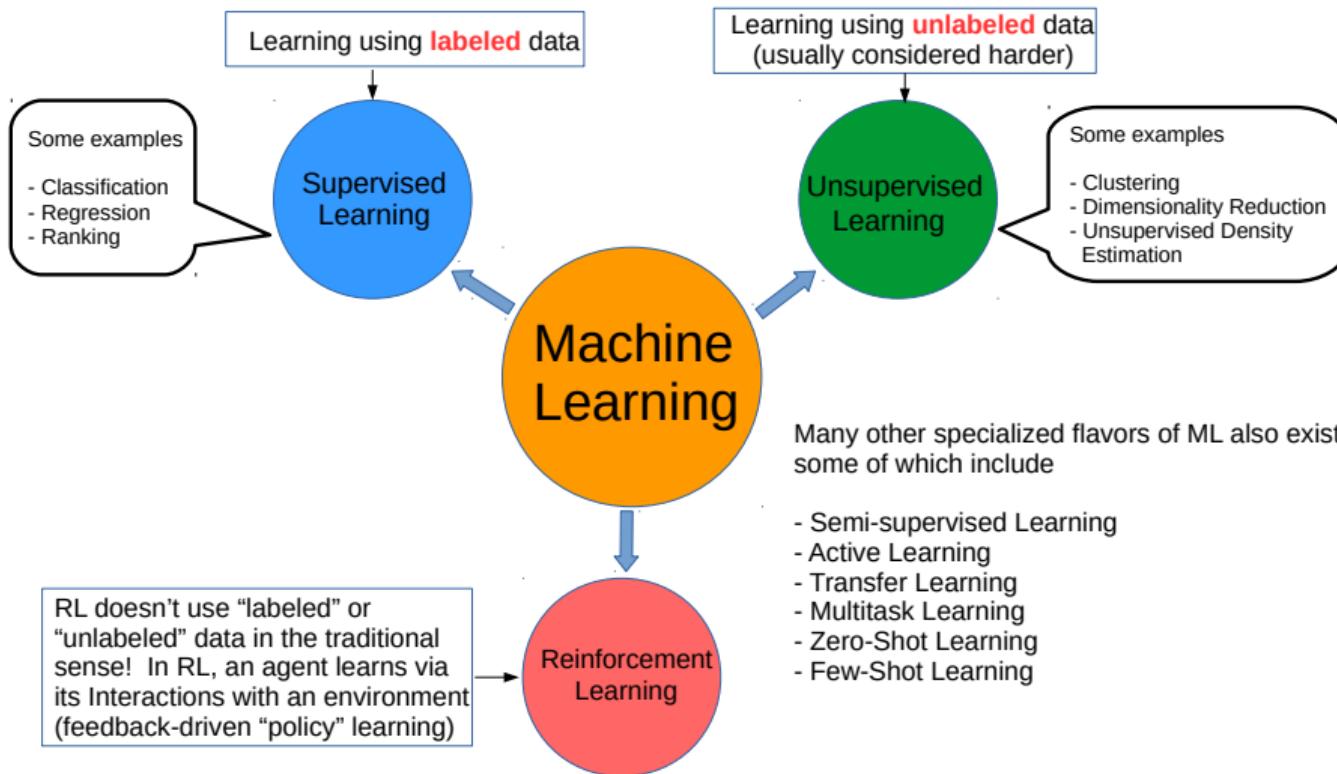
- Designing algorithms that ingest data and learn a (hypothesized) model of the data
- The learned model can be used to
 - Detect patterns/structures/themes/trends etc. in the data
 - Make predictions about future data and make decisions



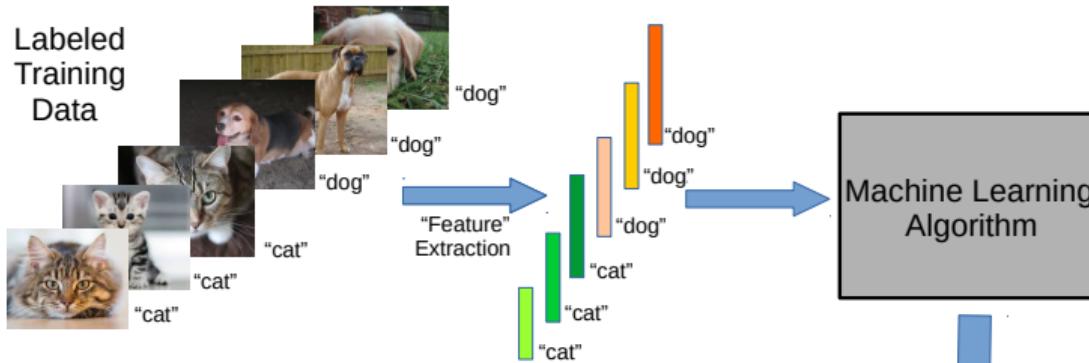
- Modern ML algorithms are heavily "data-driven"
 - No need to pre-define and hard-code all the rules (usually infeasible/impossible anyway)
 - The rules are **not "static"**; can adapt as the ML algo ingests more and more data



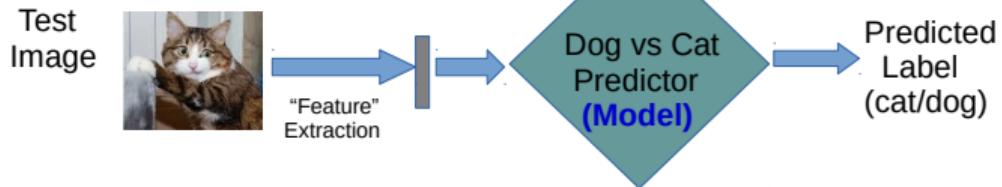
A Loose Taxonomy for ML



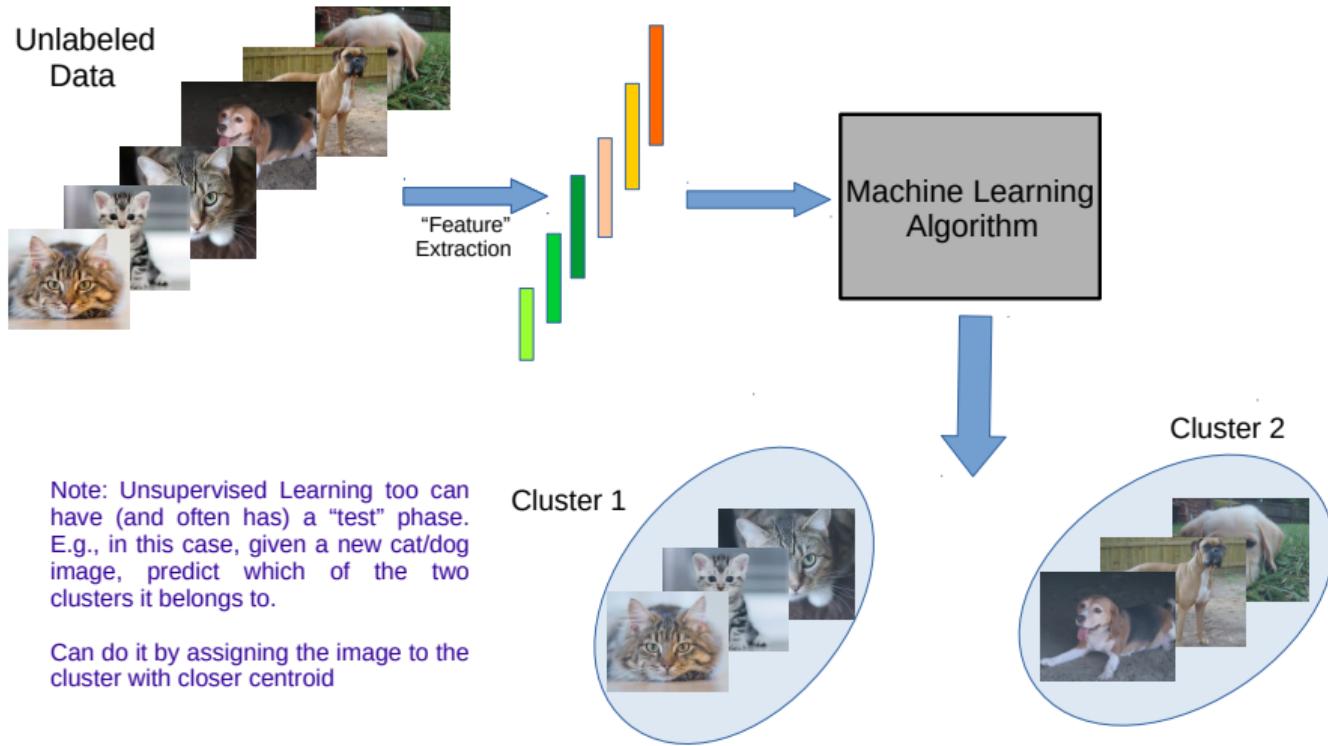
A Typical Supervised Learning Workflow (for Classification)



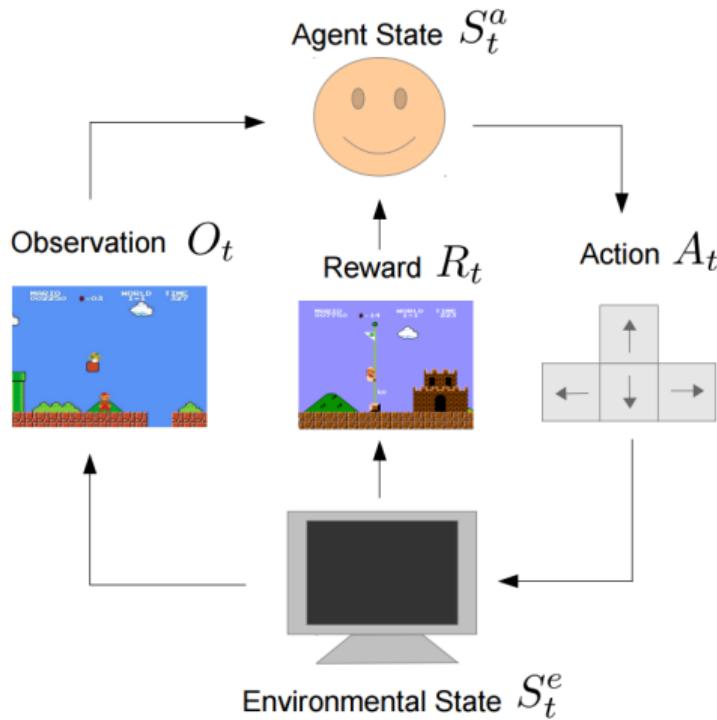
Note: The **feature extraction** phase may be part of the machine learning algorithm itself
(referred to as "feature learning" or "representation learning")
Modern "**deep learning**" algos do precisely that!



A Typical Unsupervised Learning Workflow (for Clustering)



A Typical Reinforcement Learning Workflow



Agent's goal is to learn a policy for some task

Agent does the following repeatedly

- Senses/observes the environment
- Takes an action based on its current policy
- Receives a reward for that action
- Updates its policy

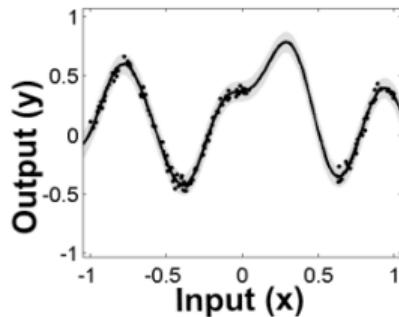
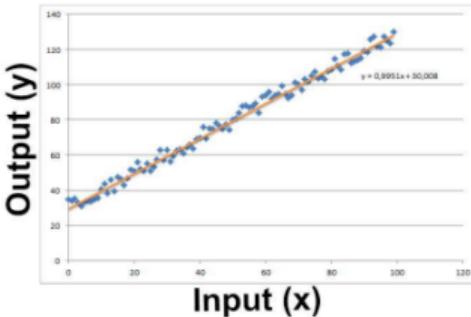
There IS supervision, not explicit (as in Supervised Learning) but rather implicit (feedback based)

Geometric View of Some Basic ML Problems

Regression

Supervised Learning: Learn a line/curve (the "model") using training data consisting of Input-output pairs (each output is a real-valued number)

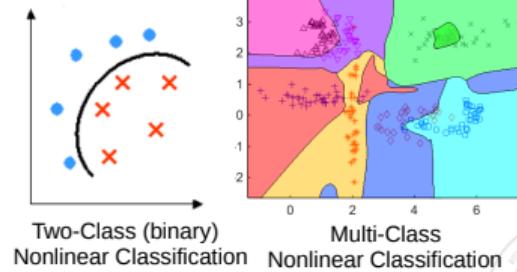
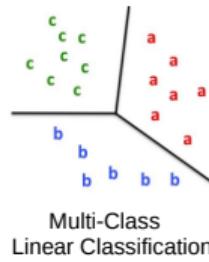
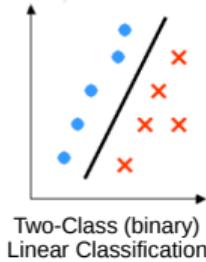
Use it to predict the outputs for new "test" inputs



Classification

Supervised Learning: Learn a linear/nonlinear separator (the "model") using training data consisting of input-output pairs (each output is discrete-valued "label" of the corresponding input)

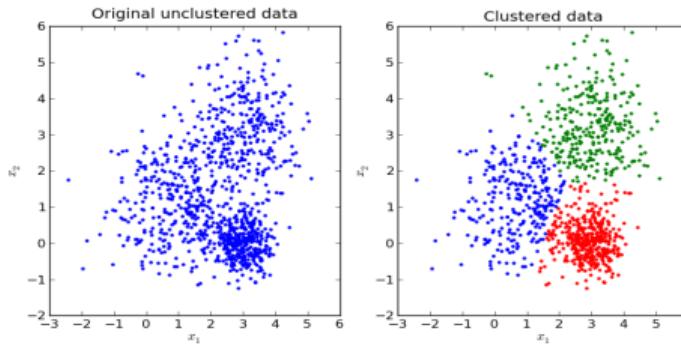
Use it to predict the labels for new "test" inputs



Geometric View of Some Basic ML Problems

Clustering

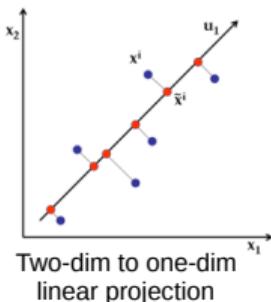
Unsupervised Learning: Learn the grouping structure for a given set of unlabeled inputs



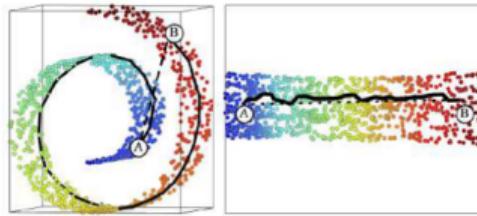
Dimensionality Reduction

Unsupervised Learning: Learn a Low-dimensional representation for a given set of high-dimensional inputs

Note: DR also comes in supervised flavors (supervised DR)



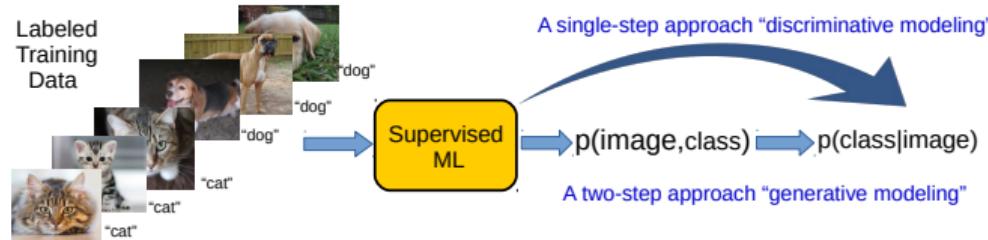
Two-dim to one-dim
linear projection



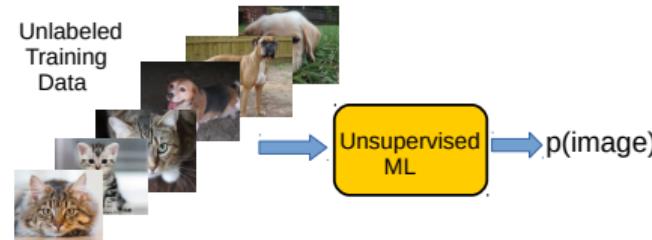
Three-dim to two-dim
nonlinear projection
(a.k.a. manifold learning)

Machine Learning = Probability Density Estimation

- Supervised Learning (“predict y given x ”) can be thought of as estimating $p(y|x)$



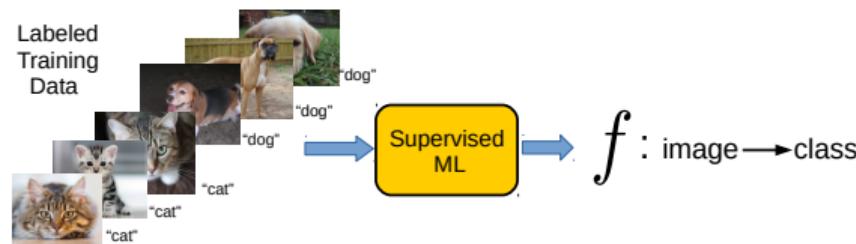
- Unsupervised Learning (“model x ”) can also be thought of as estimating $p(x)$



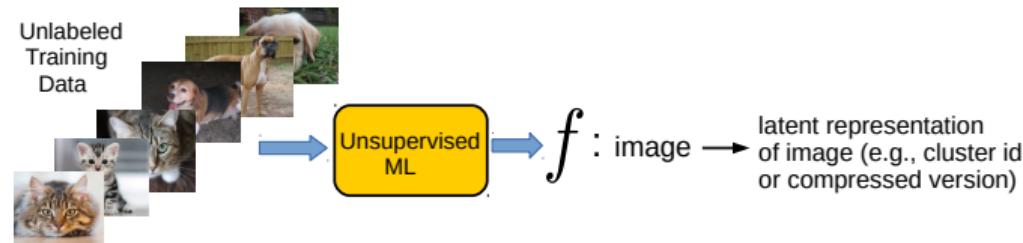
- Harder for Unsupervised Learning because there is no supervision y
- Other ML paradigms (e.g., Reinforcement Learning) can be thought of as learning prob. density

Machine Learning = Function Approximation

- Supervised Learning (“predict y given x ”) can be thought learning a function that maps x to y



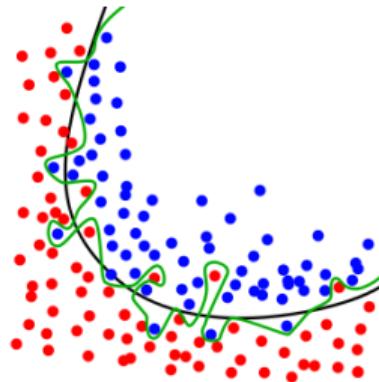
- Unsupervised Learning (“model x ”) can also be thought of as learning a function that maps x to some useful **latent representation** of x



- Harder for Unsupervised Learning because there is no supervision y
- Other ML paradigms (e.g., Reinforcement Learning) can be thought of as doing function approx.

Overfitting and Generalization

- Doing well on the training data is not enough for an ML algorithm

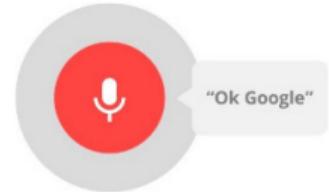
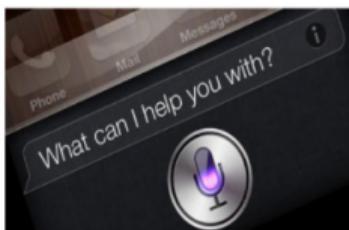


- Trying to do too well (or perfectly) on training data may lead to bad “generalization”
- Generalization: Ability of an ML algorithm to do well on future “test” data
- Simple models/functions tend to prevent overfitting and generalize well: A key principle in designing ML algorithms (called “regularization”; more on this later)

Picture courtesy: Wikipedia

Machine Learning in the real-world

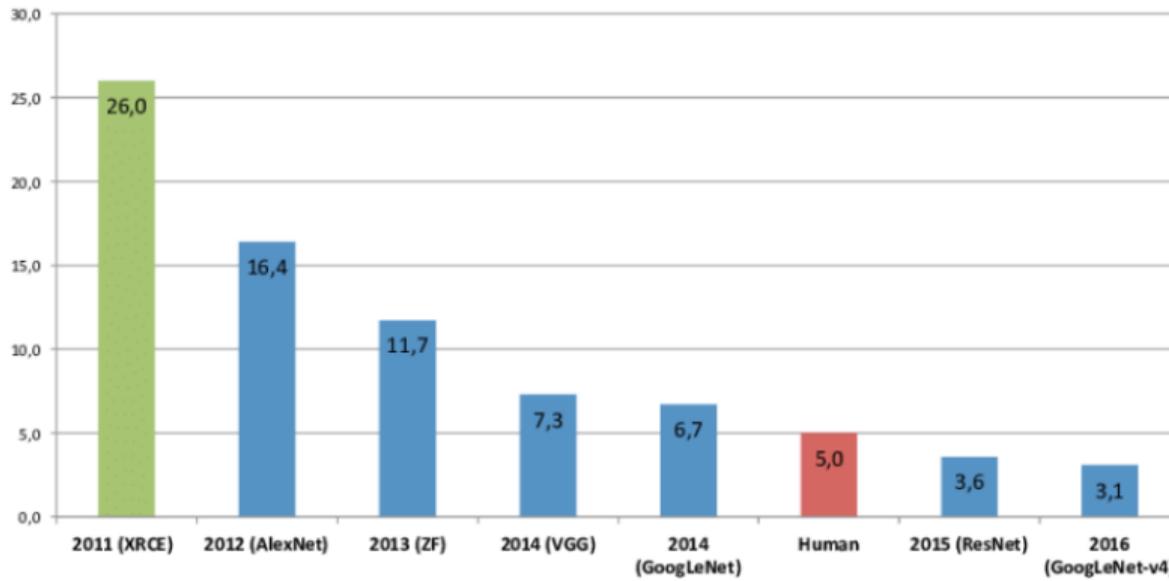
Broadly applicable in many domains (e.g., internet, robotics, healthcare and biology, computer vision, NLP, databases, computer systems, finance, etc.).



Picture courtesy: gizmodo.com, rcdronearena.com, www.wiseyak.com, www.charlesdong.com

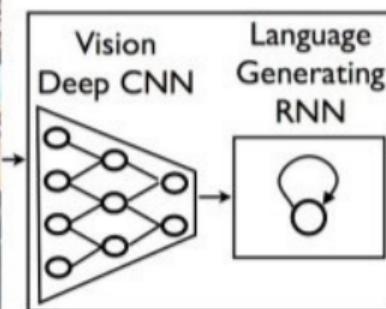
Machine Learning helps Computer Vision

ML algorithms can learn to recognize images better than humans!



Machine Learning helps Computer Vision

ML algorithms can learn to generate captions for images



**A group of people
shopping at an
outdoor market.**

**There are many
vegetables at the
fruit stand.**

<http://arxiv.org/abs/1411.4555> “Show and Tell: A Neural Image Caption Generator”

Machine Learning helps Computer Vision

ML algorithms can learn to answer questions about images (Visual QA)



What vegetable is on the plate?
Neural Net: broccoli
Ground Truth: broccoli



What color are the shoes on the person's feet ?
Neural Net: brown
Ground Truth: brown



How many school busses are there?
Neural Net: 2
Ground Truth: 2



What sport is this?
Neural Net: baseball
Ground Truth: baseball



What is on top of the refrigerator?
Neural Net: magnets
Ground Truth: cereal



What uniform is she wearing?
Neural Net: shorts
Ground Truth: girl scout



What is the table number?
Neural Net: 4
Ground Truth: 40



What are people sitting under in the back?
Neural Net: bench
Ground Truth: tent



Machine Learning helps NLP

ML algorithms can learn to translate text

English ▾



Hindi ▾



Welcome to this
course Edit

इस कोर्स में आपका स्वागत है

is kors mein aapaka svaagat hai

(even “transliterate”)



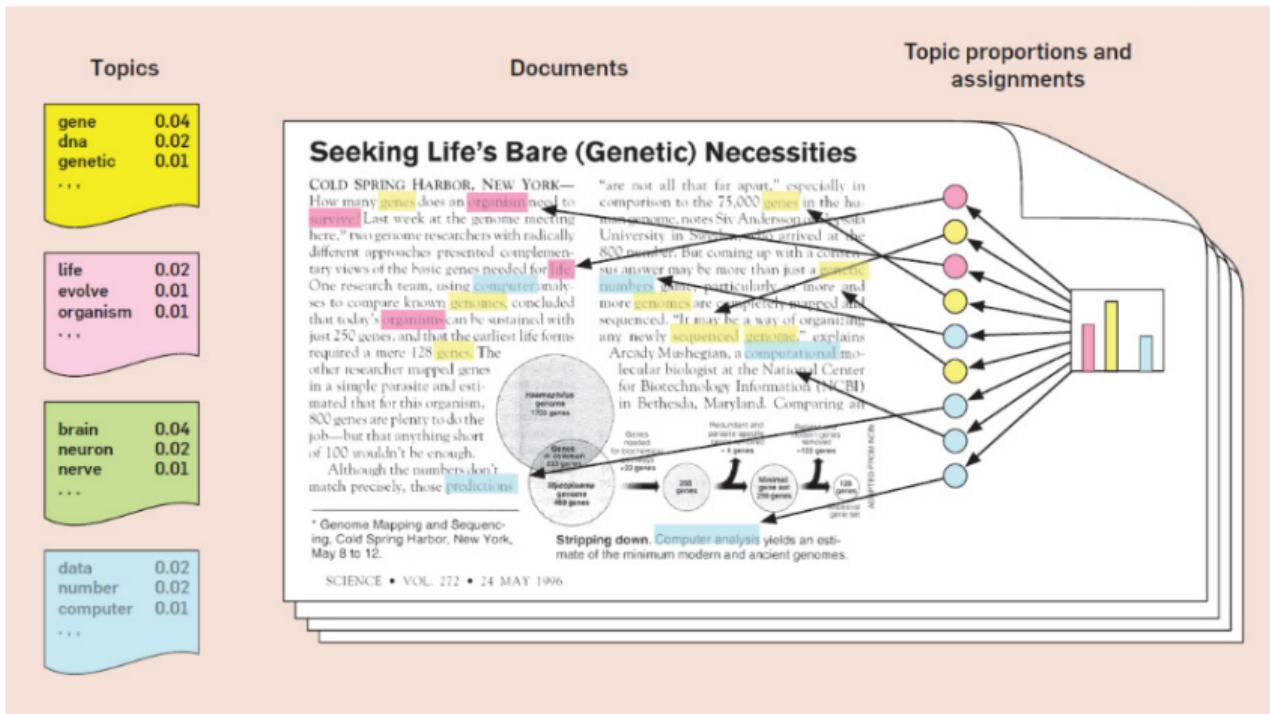
Machine Learning helps NLP

ML algorithms can learn to summarize text

Input: Article 1st sentence	Model-written headline
metro-goldwyn-mayer reported a third-quarter net loss of dlr\$ 16 million due mainly to the effect of accounting rules adopted this year	mgm reports 16 million net loss on higher revenue
starting from july 1, the island province of hainan in southern china will implement strict market access control on all incoming livestock and animal products to prevent the possible spread of epidemic diseases	hainan to curb spread of diseases
australian wine exports hit a record 52.1 million liters worth 260 million dollars (143 million us) in september, the government statistics office reported on monday	australian wine exports hit record high in september

Machine Learning helps NLP

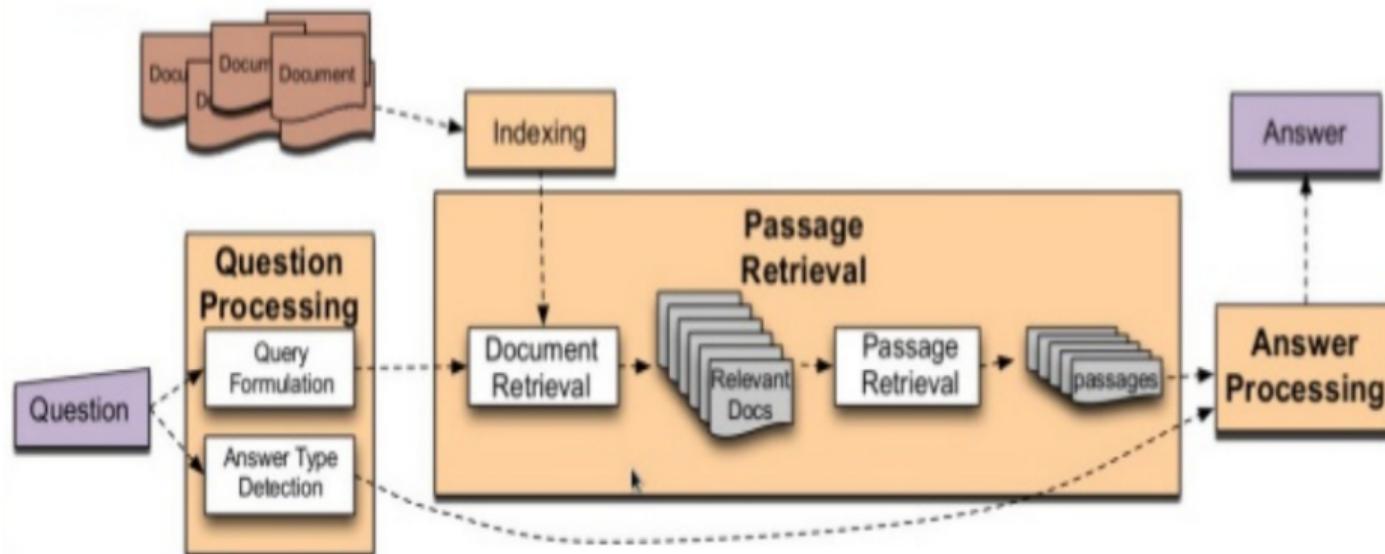
ML algorithms can learn the topics in a text corpus ("Topic Modeling")



Picture courtesy: David Blei

Machine Learning helps Search and Info Retrieval

ML algorithms can learn to search for the answer to a given question from a large database of documents



Machine Learning meets Speech Processing

ML algorithms can learn to translate speech in **real time**

PUTTING MACHINE LEARNING TO THE TEST
To provide a seamless user experience, Skype Translator uses machine learning to solve key challenges in interpreting human language, including:

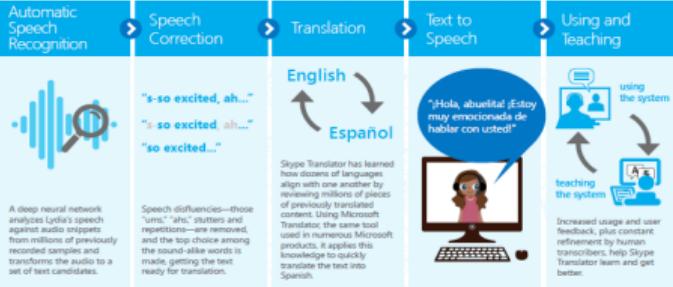
- Representing the different ways people really speak.
- Determining sentence boundaries, punctuation and case from speech.
- Disambiguating sound-alike words in context.
- Mapping words and phrases from one language to another.

NOW YOU'RE SPEAKING MY LANGUAGE (LITERALLY)



Skype has always been about making it easy to talk with family and friends all over the world. Now, by integrating advanced speech recognition and automatic translation into Skype, Skype Translator lets you speak with those you've always wished you could, even if they speak a different language.

HOW SKYPE TRANSLATOR WORKS



The process is as follows:

- Automatic Speech Recognition: A deep neural network analyzes Lydia's speech against audio snippets from tens of thousands of previously recorded samples and transforms the audio into a set of text candidates.
- Speech Correction: Speech disfluencies—those “ums,” “ahs,” stutters and repetitions—are removed, and the top ones are flagged. Then the sound-alike words are identified, and the sound-alike word is made, getting the text ready for translation.
- Translation: Skype Translator has learned how dozens of languages align with each other by reviewing millions of pieces of previously translated content. Using Microsoft’s neural machine learning used in numerous Microsoft products, it applies this knowledge to quickly translate the text into Spanish.
- Text to Speech: The text is converted into speech using Microsoft’s text-to-speech engine.
- Using and Teaching: Increased usage and user feedback, plus constant refinement by human translators, help Skype Translator learn and get better.

TRANSLATE INSTANT MESSAGES IN OVER 40 LANGUAGES

Holding a translated IM conversation is super easy: Choose a contact, turn on the Translation switch for that person, and start typing. When you hit enter (or tap send), your original message will appear in the right-hand pane, followed by its translation. Your contact on the other end will see something very similar, albeit with the translated message in their preferred language presented first. While voice translation initially supports English and Spanish only, IM translation supports over 40 languages, so feel free to experiment with them all—even Klingon!

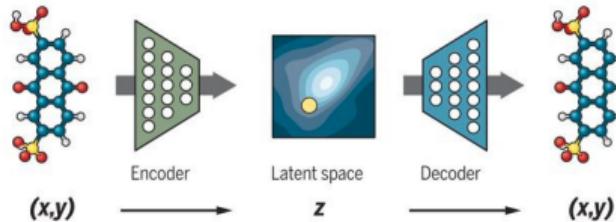


Register for the preview at www.skype.com/translator and wait for your invite.
Install the Skype Translator client.
Use Skype Translator to call someone who speaks Spanish. Or, if you speak Spanish, call someone who speaks English.
Every call you make helps Skype Translator get a little bit better. You won't see the improvement right away, but you will see gradual improvement over time.

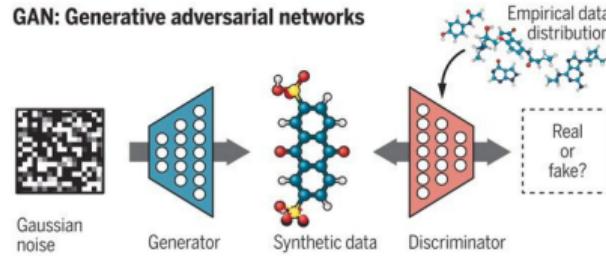
Machine Learning helps Chemistry

ML algorithms can understand properties of molecules and learn to synthesize new molecules

VAE: Variational autoencoders

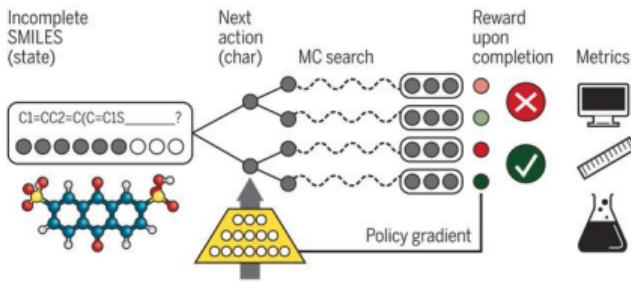


GAN: Generative adversarial networks

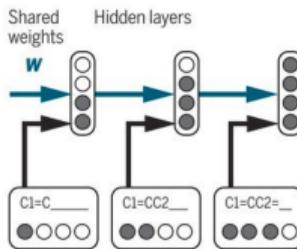


RL: Reinforcement learning

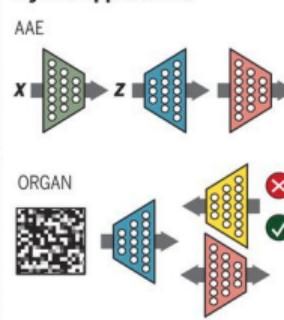
Policy gradient with Monte Carlo tree search (MCTS)



RNN: Recurrent neural network



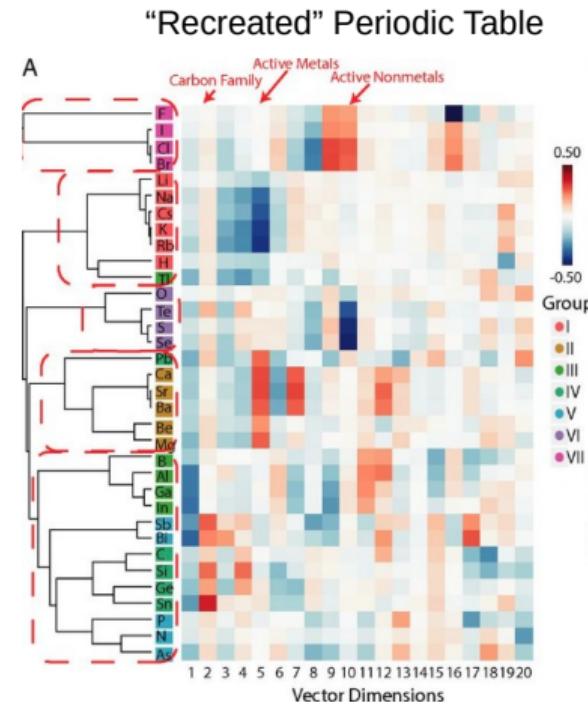
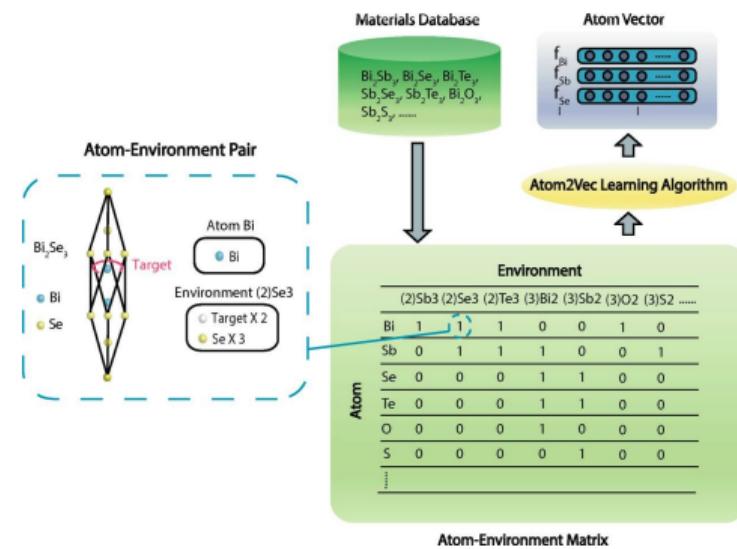
Hybrid approaches



Picture courtesy: Inverse molecular design using machine learning: Generative models for matter engineering (Science, 2018)

Machine Learning helps Chemistry

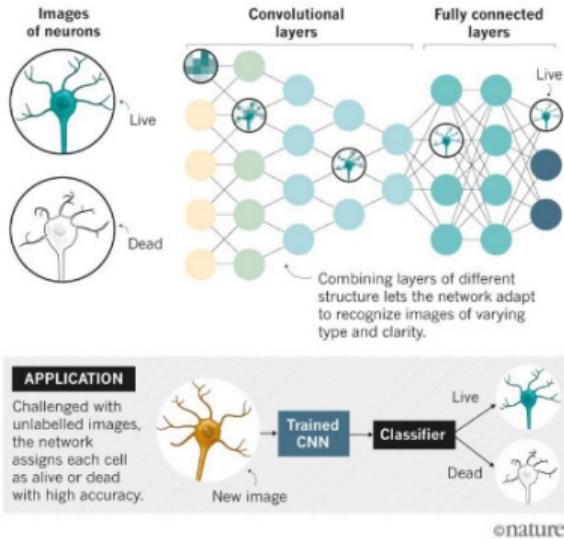
ML algorithms can “read” databases of materials and recreate the Periodic Table within hours



Picture courtesy: Learning atoms for materials discovery (PNAS, 2018)

Machine Learning helps Many Other Areas..

Biology



Finance



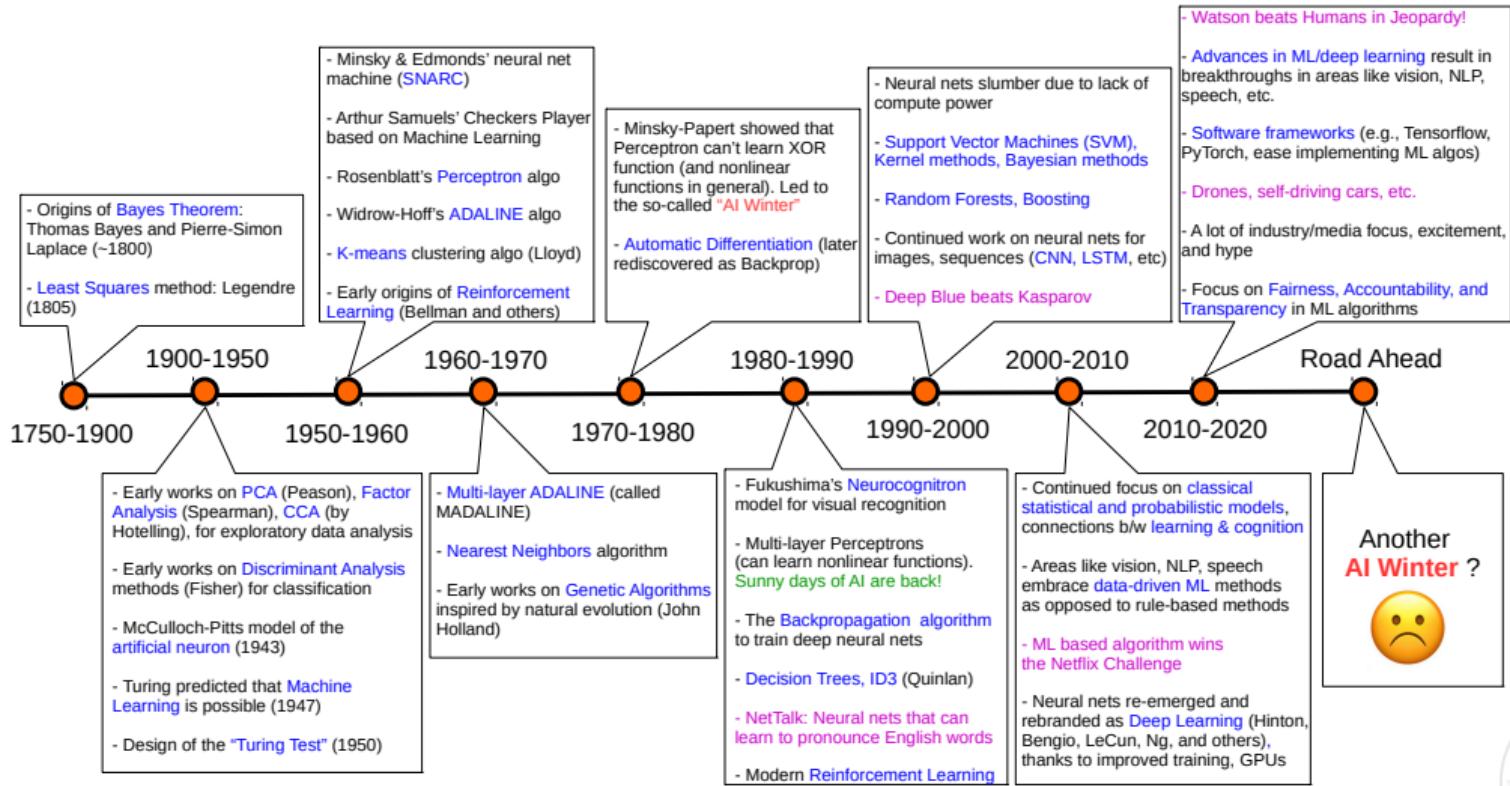
Source: Jeremy Linsley/Drew Linsley/Steve Finkbeiner/Thomas Serre

Picture courtesy: (1) <https://www.nature.com/articles/d41586-018-02174-z> (2) <https://responsiblefinanceforum.org>

Intro to Machine Learning (CS771A)

Course Logistics and Introduction to Machine Learning

Machine Learning: A Brief Timeline and Some Milestones



(Tentative) List of topics

- Supervised Learning
 - nearest-neighbors methods, decision trees
 - linear/non-linear regression and classification
- Unsupervised Learning
 - Clustering and density estimation
 - Dimensionality reduction and manifold learning
 - Latent factor models and matrix factorization
- Probabilistic Modeling
- Deep Learning
- Ensemble Methods
- Learning from sequential data
- Recent advances in ML



Course Goals

By the end of the semester, you should be able to:

- Understand how various machine learning algorithms work
- Implement them (and, hopefully, their variants/improvements) on your own
- Look at a real-world problem and identify if ML is an appropriate solution
- If so, identify what types of algorithms might be applicable
- Feel inspired to work on and learn more about Machine Learning :-)

Caution: There will be quite a bit of maths in this course (can't be avoided!). You are expected to be (or to make yourself) comfortable with multivariate calculus, linear algebra, probability and statistics. Please use the provided reference materials to brush up these concepts.

