AI & Machine Learning Intro

Complex Systems 530 - Marisa Eisenberg

What are machine learning (ML) and artificial intelligence (AI)?

- Machine learning (ML): a statistical algorithm or method that performs a task by "learning" from data
- Artificial intelligence (Al): a broader class of methods that enable computers to mimic human intelligence (e.g. vision, being able to understand and respond to language, classify objects, etc.)
- Both terms often used interchangeably, but ML is one type of AI—AI also includes deep learning, robotics, natural language processing, etc. (Although some of these are/make use of ML methods too!)



By Lollixzc - Own work, CC BY-SA 4.0, https:// commons.wikimedia.org/w/index.php?curid=122023216



Supervised, unsupervised, and reinforcement learning

- **Supervised learning**: the model learns from "labeled" data—i.e. a subset of data where the answers are given
 - E.g. Classification, ChatGPT
 - Usually more focused on prediction
- **Unsupervised learning**: the model uses unlabeled data to make decisions/understand the structure of the data
 - E.g. Clustering, PCA (principal component analysis)
 - Usually less focused on prediction, more focused on understanding structure in the data

Supervised, unsupervised, and reinforcement learning

- Reinforcement learning: the model/system interacts with an environment or data and learns as it goes based on feedback it receives
- E.g. Simulated Google Deepmind robots learning to walk
 - <u>https://www.youtube.com/watch?v=gn4nRCC9TwQ</u>
- Iterated game theory systems (e.g. iterated prisoner's dilemma)
 - <u>https://ncase.me/trust/</u>

- ML methods don't usually include the underlying mechanisms —they are focused on prediction rather than understanding why the system works the way it does (explanation)
- Many ML/AI methods are "black box"—we can test their performance and see how well they work, but we don't necessarily understand how the structure of the model relates to the predictions it makes (e.g. why does a neural network do what it does? How do the weights correspond to a decision?)
- Compare to complex systems models (e.g. Schelling model, coupled oscillators, etc., or even more standard linear models)

Do machine learning models include mechanism?

- A machine learning model like a linear model or a neural network could "learn" the gravity rule if given enough data on masses and distances
- But—what about if we were on Mars?

Why is mechanism important?

ML/AI models don't have to learn the correct mechanism from the data? Sometimes they learn confounding variables, or irrelevant things that relate to the training data they are given

Wilhelm von Osten and Clever Hans. https://en.wikipedia.org/wiki/Clever_Hans

"Clever Hans" predictors

Detecting skin cancer with AI/ML

- Neural network for skin lesion classification -Nature 2017 paper
- Built on large datasets of clinical and nonclinical images of benign and malignant lesions—performed equivalently to panel of dermatologists

Automated Classification of Skin Lesions: From Pixels to Practice, 2018 <u>https://</u> www.sciencedirect.com/science/article/pii/S0022202X18322930?via%3Dihub=

Detecting skin cancer with AI/ML

- But when tested on real data, performed worse —why?
- Without mechanism, we must be extremely careful how we set up the training data and optimization (cost function etc), otherwise the algorithm can learn the wrong thing!

Automated Classification of Skin Lesions: From Pixels to Practice, 2018 <u>https://</u> www.sciencedirect.com/science/article/pii/S0022202X18322930?via%3Dihub=

Detecting skin cancer with AI/ML

Al for pneumonia treatment decisions

- Caruana et al. 2015: Al system to predict likelihood of death in patients with pneumonia—to assist clinicians in deciding whether a certain person should be admitted or treated as an outpatient
- All predicted that patients with asthma are at lower risk of severe pneumonia than those without asthma
- This is known to be false! And could result in very problematic or fatal consequences if followed without thinking
- Why?

Al for pneumonia treatment decisions

- Asthmatic patients are high priority for care, and are more likely to seek care as soon as they show symptoms
- Asthma does NOT cause better pneumonia outcomes, but it correlates because those who are asthmatic are treated differently and earlier
- Al doesn't hypothesize a mechanism underlying the system
- Al learns correlations, not necessarily causation
- We have to be very careful about how we train these models —what optimization setup (e.g. cost function) and input data

"Black box" AI/ML can be problematic for medical decision-making

- Because of the risks of systematic bias, mismatch/ confounding/"Clever Hans" predictors, and cyberattacks, there are calls in the medical community to disclose the use of AI to patients as part of the discussion of risk (similar to how we discuss when doing genetic screening)
- Concerns that AI can propagate biases or potentially worsen medical decision-making if not implemented very carefully

Algorithmic bias, fairness, and justice

- Amazon hiring
- Healthcare decision-making
- Health systems currently use commercial prediction tools (AI/ ML) to identify and help patients with complex health needs
- Obermeyer et al. (2019) showed that a widely used algorithm, affecting millions of patients, exhibits significant racial bias
- "At a given risk score, Black patients are considerably sicker than White patients, as evidenced by signs of uncontrolled illnesses"

Algorithmic bias, fairness, and justice

- Quoting from their paper:
- Remedying this disparity would increase the percentage of Black patients receiving additional help from 17.7 to 46.5%
- The bias arises because the algorithm predicts health care costs rather than illness, but unequal access to care means that we spend less money caring for Black patients than for White patients.
- Thus, despite health care cost appearing to be an effective proxy for health by some measures of predictive accuracy, large racial biases arise.
- Be careful about your input data! Proxies like health care cost can be very problematic

Percentile of Algorithm Risk Score

Examples: Clustering

Cluster Analysis

- What is a cluster?
 - A set of objects/data points, such that the objects in the set are more similar to one another than they are to the objects outside the set/other clusters.

Uses of clustering in a modeling/ABM context

- Clustering methods can be used for a range of purposes
- Use for data analysis and understanding qualitative patterns in data set
 - Can use this to validate model outcomes when quantitative fitting is less possible
- Use clustering methods to understand the types of model behaviors observed across a sample of parameters (e.g. sample parameter space and then cluster the model results)

SIS network model example

- There are qualitative patterns differences (did the epidemic die out or persist) that aren't obvious from the quantile plot
- Clustering methods can be used to uncover these patterns in an automated way (and when higher dimensional outputs)
- Could use clustering approaches on the full trajectory or summary outputs like peak prevalence, epidemic duration

Cluster Analysis

- Broadly used in data analysis, including machine learning
- Clustering (unsupervised) vs. classification (supervised)
- Hard clustering (every element belongs to only one cluster) vs. fuzzy clustering (every element has various probabilities of belonging to a given cluster)
- Some methods find the number of clusters, others use a predefined number of clusters

Cluster Analysis

- Wide range of methods—which is best depends on the data to be clustered. Not really one 'best' method across all settings.
- In general, we want:
 - High intra-cluster similarity, low inter-cluster similarity (how to determine similarity?)
 - Potential to discover hidden features (especially in high dimensional data)

Some general classes (or clusters haha) of clustering methods:

- **Partitioning** methods (e.g. k-means clustering & other centroid methods)
- Hierarchical clustering methods
- Density-based methods
- Model or distribution-based methods (e.g. Gaussian mixture models, latent class analysis)
- Network clustering methods (community detection methods)
- & many others!

Partitioning methods

- General idea is often:
 - Construct a partition of the data into k clusters
 - Evaluate the resulting clusters and improve the partition
 - Repeat until optimal partition/clusters found
- Examples: k-means, k-medioids, k-modes (among many others)

K-means clustering

- Select k centroids (means), and each data point is assigned to the nearest centroid
- This partitions the space into Voronoi cells, which are our clusters
- For each cluster, calculate the centroid of all points
- These become the new cluster centroids
- Reassign points to nearest centroid and repeat

K-means clustering example

Randomly choose 3 cluster centers to start

The cluster centers partition the space based on which center is nearest

These are our starting **clusters**

What are the means of the data points in each cluster?

These are the new centers.

Now which data points are closest to each center?

Now which data points are closest to each center?

Redefine the clusters based on which center they're nearest

And repeat! Keep calculating the centers and redefining the clusters until they stop changing.



And repeat! Keep calculating the centers and redefining the clusters until they stop changing.



The results once the clusters and centers are fixed are your final k-means clusters.



K-means clustering

- Relatively efficient
- Can converge to local optima (e.g. depending on starting points)
 k-Means Clustering
- Have to specify k (number of clusters)
- Cannot make clusters with non-convex shapes
- Tends toward equal sized clusters



• How to handle categorical data? (e.g. can use k-modes)

Hierarchical clustering methods

- Agglomerative approach to clustering
 - Starts with small clusters (e.g. individual points) and then merges based on distance
- Divisive approach does the reverse (all one cluster then split into smaller ones)
- Many different approaches with different distance measurements, etc.

Hierarchical clustering example





- Start with all single point clusters
- Merge the two nearest clusters—forms a new cluster
- Merge the next two nearest clusters, etc.
- How to decide cluster distances? (What metric, do we use nearest point distance, furthest, centroid?)
- Capture clusters as a dendrogram—can choose resolution of clusters as desired



Hierarchical clustering

- Slow for larger data sets
- Useful for finding substructures/subclusters in data
- Assumes every data point is relevant/part of the clusters
- How to choose level of granularity?

Density-based clustering

- Decides clusters based on density of points
- Not every point need be assigned a cluster—some can be considered noise or outliers
- One of the most commonly used algorithms is DBSCAN (Density-Based Spatial Clustering of Applications with Noise)



- Choose a radius *r* and a minimum number of points *m*
- Classify each point as a:
 - **Core point** has at least *m* other points within radius *r*
 - Border point does not have *m* points within radius *r*,
 but is reachable a core point *p* i.e. can be connected
 to data point *p* by a chain of core points each within
 radius *r* of the next point
 - Outlier neither core nor border

































DBSCAN



https://towardsdatascience.com/the-5-clustering-algorithms-data-scientists-need-to-know-a36d136ef68

DBSCAN

- Can find non-convex clusters
- Automatically determines number of clusters needed
- Not every point goes into a cluster (handles outliers/noise; however can be a drawback if you want to assign all points to a cluster)
- Tends to find/work best with clusters of similar density
- How to choose radius & min points? There are rules of thumb but can be tricky! Often use min points = 2 x dim, for radius, can us elbow plot of a k-distance graph, but harder to say)

Model based methods: Gaussian Mixture Models

- Assumes the data points come from a combination of multivariate gaussians
- This seems restrictive but is often no more so than other methods (e.g. k-means in some sense assumes a centroid and resulting Voronoi diagram govern the data)
- Each data point has a probability of belonging to each cluster
- Often fit via expectation maximization (a type of maximum likelihood approach)

Model based methods: Gaussian Mixture Models

- Select number of clusters (number of gaussians to fit)
- Randomly initialize them (or better yet, use a method to pick a good starting guess)
- Compute the probability that each data point is in each cluster (based on the value of the gaussian at that point)
- Compute new parameters (μ , σ) for each gaussian that maximize this probability
- Repeat last two steps above until convergence

Model based methods: Gaussian Mixture Models



https://towardsdatascience.com/the-5-clustering-algorithms-data-scientists-need-to-know-a36d136ef68

Clustering methods

- Many different approaches! These are just a few examples
- Different methods behave better/worse on different data sets
- Testing how well a clustering method behaves can be difficult, especially in high dimensions and/or without ground truth information

Ν	۹iniBatchKMeans	AffinityPropagation	MeanShift S	SpectralClustering	Ward Agg	lomerativeCluster	ing DBSCAN	Birch	GaussianMixture
	.01s	4.34s	.07s	1.48s	.23s	.12s	.01s	.04s	.01s
	.025	4.795	.05s	2.835	.225	.125	.01s	.045	.01s
	.025	2.87s	.105	.35s	.845	.64s	.01s	.05s	.02s
	.02s	2.40s	.085	.846	.425	.325	.01s	.04s	.035
	*			*			*		
	,02 s	2.185	. 05s	. 59s	.21s	.12s	. 02s	.04s	.01s
1	.025	2.08s	.085	.4/s	.125	.115	.01s	.045	.02s

https://towardsdatascience.com/the-5-clustering-algorithms-data-scientists-need-to-know-a36d136ef68

Examples: Classification

Classification

- A form of supervised learning
- Typically starts with a known set of classes (can learn this via unsupervised learning, e.g. clustering for example)
- Given training data that is labeled with which class it belongs to, predict which class new (unlabeled) data belongs to

Puppy or bagel?


Classification

- Many (many!) different kinds: neural networks, decision trees, Bayes classifiers, K-nearest neighbors, etc.
- 'Lazy' vs 'eager' learners
 - Lazy learners (related: instance-based learning) load all the training data and then use the training data when new inputs come in to decide (often fast to train, but slow to operate on new data)
 - Eager learners build the classifier on the training data first, then just apply the classifier (e.g. some model) on new input data (often slow to train, but fast on new data)

k-Nearest-Neighbors (kNN)

- Decide output (class) for new data point based on the most common class among the nearest k neighbors in the training set
- No explicit training phase, but the k nearest neighbors can be viewed as the training data
- Lazy learner/instance-based learner

Neural Networks

Neural networks

- Simulated "neurons" form a model to accomplish a wide range of tasks
- These neurons are more like a "toilet model" of a neuron—they take in inputs (numbers) from and if they cross a threshold send inputs to the next layer
- Training (i.e. model fitting) involves adjusting the weights



- **Neurons** each agent or node in the network. Neurons in are essentially functions: they take as inputs the weighted sum of the incoming connections and return some function of that number they receive. The function that they use is called the activation function.
- **Connections** the directed edges from one neuron to another in the network.
- Weights how the receiving neuron weights each edge as it calculates its total inputs
- Layers most neural networks are organized in layers (analogous to many brain regions)



Binary activation function

- When the weighted sum of the inputs is greater than some threshold, the neuron fires
- Mostly used for research rather than real-world Al/ ML in practice now, but early neural networks used this



Perceptron

- Simple neural network model that uses a binary activation function and a single layer to classify objects or do edge detection
- Developed in 1943
- Basically does a single linear classifier



Activation functions often represent firing rate rather than binary firing

- Sigmoid activation function
- In this case, the "neuron" should really be thought of as measuring a small patch or population of neurons



Activation functions often represent firing rate rather than binary firing

- Rectified linear unit or ReLU
- 0 until some threshold then linear
- In this case, the "neuron" should really be thought of as measuring a small patch or population of neurons



One neuron is basically linear or logistic regression (or similar regression models)



Feedforward neural networks: layers of neurons allow them to do more advanced processing



- 2-layer Neural Network
 - 1 hidden layer of 4 neurons
 - 1 output layer of 2 neurons
- We don't count the input layer as this doesn't do any processing, just represents the inputs (i.e. the data) as a set of numbers that feed into the subsequent layers

Feedforward neural networks: layers of neurons allow them to do more advanced processing



- 3-layer neural network
- 2 hidden layers (4 neurons each), 1 output layer

Feedforward neural networks: What does all this do?

- Each layer in a feedforward neural network acts like matrix multiplication (weight matrix times the vector of previous later outputs)
- The layers end up acting as repeated matrix operations (linear) interwoven with a nonlinear activation function



Universal approximators

- Many such theorems!
- Similar flavor to Stone-Weierstrass theorem for polynomials
- Neural networks with at least one hidden layer are universal approximators (under various conditions)
- I.e. they can approximate any continuous function to arbitrary precision
- (Note the perceptron has no hidden layer!)

Deep neural networks

- Most advanced neural networks (e.g. ChatGPT) are deep neural networks, meaning they have multiple (often many!) hidden layers
- Why?
- Just because something can approximate a classification/ function/data set, doesn't necessarily make it a good ML/ Al algorithm in practice (for example, connecting the dots also approximates well!)
- Deep neural networks often perform shallow ones



By Sven Behnke - Own work, CC BY-SA 4.0, https://commons.wikimedia.org/w/index.php?curid=82466022