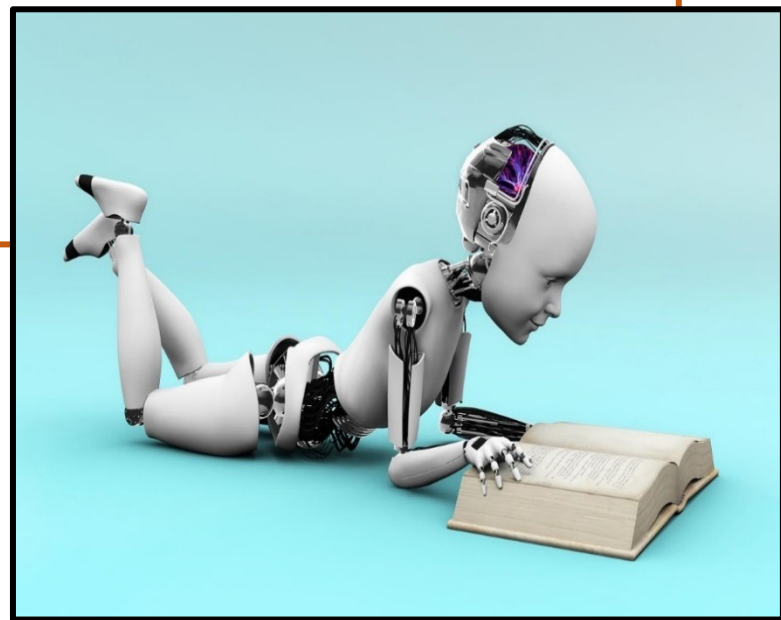


# Machine Learning from Data, Summary of topics

Zohar Yakhini & Ariel Shamir

Riechman University  
2022/TASHPAB



# Types of Learning Tasks

- Regression
  - + Given training data  $\{x_i, y_i\}$ , learn a function  $f$  to be used to predict the value  $y$  for a new data point  $x$ :  $\hat{y} = f(x)$
- Classification
  - + Given training data  $\{x_i, t_i\}$  where  $t_i \in \{0,1\}$ , learn a mechanism that determines, for a new data point  $x$ , its label  $t(x)$
  - + Also applies to multiple categories
- Parameter Estimation
  - + Given  $\{x_i\}$  find a PDF that best explains the data
- Unsupervised Learning
  - + Given  $\{x_i\}$  find regularities, such as clustering

<u>Techniques/topics we have learned and their properties</u>	Numerical or categorical features/attributes?	Algorithm Type	Related techniques and concepts; Comments
<b>Linear regression</b>	Numerical only	Regression	Normalization, Gradient descent, pseudo-inverse Ridge, LASSO
<b>Decision trees</b>	Any type of attributes. May need to scale/normalize	Both regression and Classification	Entropy, Gini
<b>Bayes classifiers (full and naïve)</b>	Both (Iris and playing tennis)	Classification	Probability distributions, prior and posterior probabilities
<b>MLE</b>		Estimation and inference	Poisson distribution, Binomial distribution, Normal distribution, <b>EM algorithm</b>
<b>KNN</b>	Mostly numerical but can be adapted	Both regression and Classification	
<b>Perceptron and dual perceptron</b>	Numerical	Classification	Mapping into higher dimensions, Cover's Thm, Kernels
<b>SVM</b> (soft margins, hard margins)	Numerical	Classification	Lagrange multipliers, kernels, slack variables

<b><u>Techniques/topics we have learned and their properties</u></b>	Numerical or categorical features/attributes?	Algorithm Type	Related techniques and concepts; Comments
<b>Logistic regression</b>	Numerical	Classification	Gradient descent
<b>Sample complexity</b>		Classification	Bounds on sample complexity, PAC learning
<b>VC dimension</b>		Classification	
<b>Clustering</b>	Numerical	Unsupervised	K-means, Naïve cluster growing, hierarchical clustering
<b>Performance evaluation</b>	Both	Both regression and Classification	Confusion matrix, Cost function, TPR/FPR and ROC curves, PR curves, conf intervals
<b>PCA, LDA</b>	Numerical	Unsupervised & Supervised	Dimensionality Reduction
<b>tSNE, MDS</b>	Numerical	Unsupervised & Supervised	Only introduced very briefly

# General techniques and principles

- Gradient descent
- Pseudo-inverse
- Stochastic, batch, mini-batch approaches (to running through training data)
- Cross validation
- Split to training and test data
- Multidimensional probability distributions; the covariance matrix
- Confusion matrix, ROC curves, PR curves, Conf Intervals

# General techniques and principles

- Vector geometry in  $\mathbf{R}^n$
- MLE, EM algorithm
- Confidence intervals for the total error
- Histograms for representing distributions
- Entropy, Gini index
- Conditional independence
- Cover's Thm
- Total variance principle

# General techniques and principles

- Lagrange multipliers for solving constrained optimization problems
- Kernels
- Sample complexity
- VC dimensions
- Feature selection techniques and how to work with them
- Dimensionality reduction & representation learning

# Happy Summer!!

