# Lahore University of Management Sciences
## CS 535/EE 514 Machine Learning
Fall 2023-24

| Course description | |
|---|---|
| Machine learning (ML) techniques allow computers to adapt to data and solve new problems related to previously encountered problems more efficiently. Such methods enable machines to perform practical exploratory and predictive tasks without being explicitly programmed. ML finds applications in speech recognition and synthesis, machine translation, object recognition, chatbots, question-answering, natural language understanding, anomaly detection, medical diagnosis and prognosis, autonomous vehicles and robots, time series forecasting, and much more. This introductory course covers the theoretical foundations and practical applications of ML and the design, implementation, and analysis of various ML algorithms. Students will learn to compare across and choose the most appropriate algorithms for multiple problem types and be able to design and implement their solutions. Students will be prepared for industry and academia and for pursuing advanced courses. | |

| Course distribution | |
|---|---|
| Elective | This is an elective course. |
| Open for Student Category | Juniors, seniors, and graduates. |
| Close for Student Category | Please see the prerequisites below. |

| Course prerequisites |
|---|
| <ul><li>Undergrads (Seniors/Juniors) must have passed:<ul><li>An Ugrad/Grad course in Probability (MATH230 (Probability) OR DISC203 (Probability & Statistics) OR CS501 (Applied Probability)) OR ECON230 (Statistics and Data Analysis))</li><li>And a programming course (CS200/EE201 (Intro. to Programming)</li><li>And a course on Linear Algebra (MATH120 (LA with Diff. Equations))</li></ul></li><li>Grads are strongly advised to brush up their programming skills and take CS501 (Applied Probability), may be in parallel with ML</li><li>All students must possess strong programming skills and proficiency in algorithm implementation in JAVA/C/Python/MATLAB</li></ul> |

| Course Offering Details | | | | |
|---|---|---|---|---|
| Credit Hours | 3 hours | | | |
| Lecture(s) | Nbr of lec(s) per week | 2 | Duration | 75 minutes |
| Recitation/Lab (per week) | Nbr of lec(s) per week | | Duration | |
| Tutorial (per week) | Nbr of lec(s) per week | 1 (optional) | Duration | 50 minutes |

| | |
|---|---|
| Instructor | Agha Ali Raza |
| Room No. | SBASSE 9-G49A |
| Office Hours | TBA |
| Email | agha.ali.raza@lums.edu.pk |
| Telephone | 3336 |
| Secretary/TA | Zohaib Khan (24100074@lums.edu.pk), Nida Tanveer (24100091@lums.edu.pk), Saad Hasan Iqbal (24020320@lums.edu.pk), Syeda Mah Noor Asad (24100243@lums.edu.pk), Alina Faisal (24100314@lums.edu.pk), Muawiz Feroze Khan (24100052@lums.edu.pk), Fatima Ali (24100048@lums.edu.pk), Rafey Rana (24100103@lums.edu.pk), Zain Ali Khokhar (24100130@lums.edu.pk), Muhammad Haad Zahid (24100134@lums.edu.pk), Mughees Ur Rehman (24100086@lums.edu.pk), Syed Kabir Ahmad (24100249@lums.edu.pk) |
| TA Office Hours | TBA |
| Course URL (if any) | From last offering: https://www.c-salt.org/courses/machine-learning-f2021 |

**Course Teaching Methodology (Please mention the following details in plain text)**

- **Lectures:** In-person.
- **TA Sessions:** TAs will conduct asynchronous and synchronous sessions (in-person and online) to cover tutorials related to assignments.
- **Exams:** Exams will be conducted in person in pre scheduled sessions.
- **Quizzes:** Quizzes will be conducted during announced class timings.
- **Class discussions:** There will be a slack channel for all discussions (general, assignments, quizzes, etc.)

**PROGRAM EDUCATIONAL OBJECTIVES (PEOs)**

| | |
|---|---|
| PEO-01 | Demonstrate excellence in the profession through in-depth knowledge and skills in the field of Computing. |
| PEO-02 | Engage in continuous professional development and exhibit a quest for learning. |
| PEO-03 | Show professional integrity and commitment to societal responsibilities. |

**Course Objectives**

The goal of this course is to get the students excited about Machine Learning and to enable them to:
- Develop a firm grip on the theory behind statistical learning
- Understand and rigorously go through the phases of the design, implementation, and evaluation of fundamental ML algorithms
- Choose the appropriate algorithm for each problem type and be able to compare the strengths and weaknesses of the algorithms
- Appreciate the end-to-end organic integration of ML in its application areas, from data sources, annotation pipelines, and choice of algorithms to societal biases, explainability of models, and potential to impact and even disrupt existing processes

**COURSE LEARNING OUTCOMES (CLOs)**

| | By the end of the course, students should be able to: |
|---|---|
| CLO1: | • Develop an appreciation for what is involved in learning models from data, and integrating ML in existing real-world processes |
| CLO2: | • Thoroughly understand the ML pipeline from design and data gathering to meaningful and relevant evaluation |
| CLO3: | • Learn a wide variety of learning algorithms, and formulate and implement solutions to machine learning problems |
| CLO4: | • Apply algorithms to real-world problems, optimize the trained models and report on the expected performance |

| CLO | CLO Statement | Bloom's Cognitive Level | PLOs/Graduate Attributes (Seoul Accord) |
|---|---|---|---|
| CLO1 | Develop an appreciation for what is involved in learning models from data, and integrating ML in existing real-world processes | C2, C3 | PLO2 |
| CLO2 | Thoroughly understand the ML pipeline from design and data gathering to meaningful and relevant evaluation | C3, C4, C5 | PLO2, PLO3, PLO4 |
| CLO3 | Learn a wide variety of learning algorithms, and formulate and implement solutions to machine learning problems | C2, C5, C6 | PLO4, PLO5 |
| CLO4 | Apply algorithms to real-world problems, optimize the trained models and report on the expected performance | C5, C6 | PLO5 |

**Grading Breakup and Policy**

| Assessment | Weight (%) | Related CLOs | ACM Recommended Disposition |
|---|---|---|---|
| Programming assignment(s) | 25% | CLO2, CLO3, CLO4 | D3, D4, D7, D9, D10 |
| Quizzes | 25% | CLO1, CLO2 | D4, D7, D9, D10 |
| Project | 20% | CLO1 – CLO4 | D1, D3, D4, D5, D6, D7, D8, D9, D11 |
| Reading assignment(s)/homework(s)/Implementation of Research Paper(s)/viva | 15% | CLO1 – CLO4 | D3, D4, D7, D9, D10 |
| Final examination + viva | 15% | CLO1 – CLO4 | D4, D7, D9, D10 |

**Examination detail**

| | | |
|---|---|---|
| Midterm Exam | Yes/No: | No |
| | Duration: | |
| | Exam Specifications: | |

| Final Exam | Yes/No: | Yes |
| | Duration: | 2.5 – 3 hours |
| | Exam Specifications: | In-person exam |

| **SSE Council on Equity and Belonging** |
| --- |
| In addition to LUMS resources, SSE's **Council on Belonging and Equity** is committed to devising ways to provide a safe, inclusive and respectful learning, living, and working environment for students, faculty and staff. To seek counsel related to any issues, please feel free to approach either a member of the council or email at cbe.sse@lums.edu.pk. |

| **Mental Health Support at LUMS** |
| --- |
| For matters relating to counselling, kindly email student.counselling@lums.edu.pk, or visit  https://osa.lums.edu.pk/content/student-counselling-office for more information. You are welcome to write to me or speak to me if you find that your mental health is impacting your ability to participate in the course. However, should you choose not to do so, please contact the Counseling Unit and speak to a counsellor or speak to the OSA team and ask them to write to me so that any necessary accommodations can be made. |

| **Harassment Policy** |
| --- |
| SSE, LUMS and particularly this class, is a harassment free zone. Harassment of any kind is unacceptable, whether it be sexual harassment, online harassment, bullying, coercion, stalking, verbal or physical abuse of any kind. Harassment is a very broad term; it includes both direct and indirect behavior, it may be physical or psychological in nature, it may be perpetrated online or offline, on campus and off campus. It may be one offense, or it may comprise of several incidents which together amount to sexual harassment. It may include overt requests for sexual favors but can also constitute verbal or written communication of a loaded nature. Further details of what may constitute harassment may be found in the LUMS Sexual Harassment Policy, which is available as part of the university code of conduct.<br><br>LUMS has a Sexual Harassment Policy and a Sexual Harassment Inquiry Committee (SHIC). Any member of the LUMS community can file a formal or informal complaint with the SHIC. If you are unsure about the process of filing a complaint, wish to discuss your options or have any questions, concerns, or complaints, please write to the Office of Accessibility and Inclusion (OAI, oai@lums.edu.pk) and SHIC (shic@lums.edu.pk) —both of them exist to help and support you and they will do their best to assist you in whatever way they can. You can find more details regarding the LUMS sexual harassment policy here.<br><br>**To file a complaint, please write to harassment@lums.edu.pk.** |

| **Rights and Code of Conduct for Online Teaching** |
| --- |
| A misuse of online modes of communication is unacceptable. TAs and faculty will seek consent before the recording of live online lectures or tutorials. Please ensure if you do not wish to be recorded during a session to inform the faculty member in a timely manner. Please also ensure that you prioritize formal means of communication (email, LMS) over informal means to communicate with course staff. |

| **Course overview** | | | | |
| --- | --- | --- | --- | --- |
| **W** | **Topics** | **Recommended Readings** | **Related CLOs** | **ACM Comp Knowledge Landscape** |
| 1. | **Course overview**<br>● What is ML? Traditional CS vs. ML, history of ML, AI vs. ML<br>● Classification and Regression with examples. Training and Testing.<br>● Rules vs. Patterns, Deterministic vs. Probabilistic, Certainty vs. Uncertainty.<br>● Learning: Supervised, unsupervised, semi-supervised<br>● Labeled data sources: Expert annotators, crowd<br>● Example ML application areas: Speech and Language Technologies<br>● Challenges and Opportunities of ML:<br>   ○ Explainability<br>   ○ Fairness and Societal Biases<br>   ○ ML for Social Good, ML for Development (ML4D), Speech and Language Technologies for Development (SLT4D) | ● Murphy chapter 1<br>● Alpaydin, chapter 1 | CLO1, CLO2 | |
| 2. | **Supervised Learning**<br>● Features, Labels, Training, Testing, Classification, Regression.<br>● Formalizing the supervised learning setup<br>● Feature spaces and feature vectors<br>   ○ Sparse and dense feature vectors, one-hot vectors<br>   ○ Bag-of-word features | ● Murphy: 1.1, 1.2, 1.4.2, 1.4.3, 1.4.9<br>**Recommended topics:**<br>● Goals of Cross Validation: | CLO1, CLO2, CLO3 | |

| | | | | |
|---|---|---|---|---|
| | • Label spaces<br>   ◦ Label spaces for classification (binary and multiclass) and regression<br>• Hypothesis spaces<br>   ◦ The No Free Lunch theorem<br>   ◦ Choosing the hypothesis class $H$ and hypothesis $h \in H$<br>   ◦ Various Algorithms for traversing hypothesis classes:<br>      ■ Pick $h$ randomly<br>      ■ Try every $h$<br>      ■ Just output the label of the training data (memorizer)<br>• Evaluating hypotheses: Loss functions and goals of optimization<br>   ◦ Zero-One<br>   ◦ Squared<br>   ◦ Absolute<br>• Loss reduction and Generalization in Learning<br>   ◦ Memorizers<br>   ◦ Smoothing and Priors<br>   ◦ Tradeoff between Bias and Variance<br>• Sampling from the distribution $P(X, Y)$<br>   ◦ Representative datasets<br>   ◦ Training, validation, and testing<br>• How to split the dataset $D$?<br>   ◦ Time series data<br>   ◦ Independent and Identically Distributed (IID)<br>• The weak law of large numbers ($\epsilon_{TE} \to \epsilon \; as \; |D_{TE}| \to +\infty$)<br>• How to prevent overfitting to test data? Do's and Don'ts<br>• Validation sets (dev sets) and Cross Validation | Model selection, training, and performance estimation<br>• Types of Cross Validation and Pros and Cons<br>• Exhaustive<br>   • Leave-p-out<br>   • Leave-one-out<br>• Non-Exhaustive<br>   • k-fold<br>   • Holdout<br>   • Repeated random subsampling<br>• Nested<br>   • $k * l$ fold<br>   • $k$-fold with validation and test sets<br>   • Bootstrapping<br>   • Stratified cross validation<br>• Time series cross validation (forward chaining – Rolling origin) | | |
| 3. | **The K-Nearest Neighbor Classifier** (An instance-based, lazy, discriminative, non-linear, non-parametric classifier)<br>• KNN – The Basics<br>   ◦ Nearest neighbor classification rule<br>   ◦ KNN formal definition<br>   ◦ KNN decision boundaries and Voronoi Tessellations<br>   ◦ Properties of KNN: Non-parametric, used for classification and regression, instance-based, lazy<br>• KNN Similarity/Distance measures and constraints<br>   ◦ Minkowski Distances (Manhattan, Euclidean and Chebyshev)<br>• The KNN algorithm and implementation<br>   ◦ KNN regression and classification with examples<br>   ◦ Space and Time complexity<br>   ◦ Bias/variance tradeoff as<br>   ◦ $K \to 1$ and $K \to n$<br>   ◦ Tuning the hyperparameter: K<br>   ◦ KNN: The good, the bad and the ugly<br>• KNN error bounds as $n \to \infty$<br>   ◦ Bayes Error<br>   ◦ 1-NN error as $n \to \infty$<br>• KNN Enhancements<br>   ◦ Parzen Windows and Kernels<br>   ◦ K-D trees<br>   ◦ Inverted lists<br>   ◦ Locality sensitive hashing<br>• The Curse of Dimensionality<br>   ◦ Demonstration and examples<br>   ◦ Challenges and opportunities | • Murphy: 1.1, 1.2, 1.4.2, 1.4.3, 1.4.9<br>• Videos of different values of k<br>• Video describing nearest neighbors<br>• A nice explanation of nearest neighbors | CLO3, CLO4 | |

| | | | | |
|---|---|---|---|---|
| | ○      Lower dimensional subspaces and manifolds in higher dimensional ambient space<br><br>***Tutorial Topics:***<br>● Review of KNNs<br>● Review of Voronoi Tessellations<br>● Review (and Expansion) of Error Bounds<br>● Review (and Expansion) of the Curse of Dimensionality | | | |
| 4. | **Evaluation metrics**<br>● The Confusion matrix (contingency tables) – binary and multi-label<br>     ○ True and False – Positives and Negatives<br>     ○ Type I and type II errors<br>● Performance Metrics<br>     ○ Accuracy, Sensitivity (recall, TPR), Specificity (TNR), Precision (Positive Predictive Value), Negative predictive value<br>     ○ False acceptance rate, False rejection rate<br>     ○ Examples – pros and cons<br>● The need for a combined measure<br>     ○ Types of averages: AM, GM, HM<br>     ○ F-$\beta$-measure, F-1-measure<br>● Multiclass Classification<br>     ○ Any-of (multi-label) classification<br>     ○ One-of (multinomial) classification-<br>     ○ Micro and Macro averaging<br>● Gold labels and annotation of data<br>     ○ Inter-annotator agreements<br>     ○ Cohen's Kappa and Krippendorff's alpha<br>● Evaluation of Classifiers, thresholds, comparing classifiers, imbalanced classes<br>     ○ Receiver operating Characteristic (ROC) and Precision-Recall (P-R) Curves<br>     ○ ROC *Area Under the Curve* (AUC)<br>     ○ Equal Error Rate (EER) and Biometric Systems<br><br>***Tutorial Topics:***<br>● Review of Cohen's Kappa<br>● Review of ROC and Precision-Recall Curves<br>● Assignment 1 | ● SLP3: 4.7-4.9 | CLO2, CLO3 | |
| 5. | **Linear Regression**<br>● Motivation for linearity<br>● Revision<br>     ○ Lines, planes, hyperplanes, and vectors<br>         ■ Lines and planes: Normal form and slope-intercept form<br>         ■ Decision boundaries with perpendicular weight vectors<br>         ■ Distance between a hyperplane and a point<br>     ○ The Dot Product<br>     ○ The geometric interpretation of absorbing the bias term<br>     ○ Visualizing n dimensions<br>● Linear Regression<br>     ○ Intuition<br>     ○ Derivation and implementation<br>         ■ Linear regression with one variable<br>         ■ Cost function: Square Loss, Mean Square Loss (MSE),<br>         ■ $\frac{1}{2}$ Mean Square Loss<br>         ■ Motivating gradient descent<br>● The Gradient Descent Algorithm<br>     ○ Gradients<br>     ○ The step size,<br>     ○ $\alpha$<br>     ○ Convex and non-convex cost functions – Global and local optima<br>     ○ The batch gradient descent algorithm<br>     ○ Types of Gradient Descent: Batch, Minibatch, and Stochastic<br>● Multivariate Linear Regression<br>     ○ Multivariate Gradient Descent<br>     ○ Vectorizing the notation | ● ESLII Ch3<br>● Murphy 7-7.5.1, 7.5.4 | CLO2, CLO3, CLO4 | |

| | | | | |
|---|---|---|---|---|
| | • Practical Issues of linear regression<br>  ○ Feature Scaling, local minima, ravines, saddle points, tracking progress in GD<br>  ○ Hyperparameters: Learning rate<br>• Polynomial regression<br><br>*Tutorial Topics:*<br>• Review of Hyperplanes<br>• Review of Gradient Descent<br>• Demo on Polynomial Regression<br>• Review (and Expansion) of Standardization for Linear Regression | | | |
| 6. | **Linear Regression: Bias and Variance**<br>• How to recognize high variance/high bias scenarios?<br>  ○ Underfitting and Overfitting<br>  ○ How to reduce bias and variance?<br>    ■ Cross validation<br>    ■ Feature selection<br>• Manual Feature Selection<br>  ○ Scatter Diagrams and Plots<br>  ○ Eyeballing those Correlations!<br>• Regularization<br>  ○ Motivation – The fitting problem<br>  ○ L2 Regularization or Ridge Regression<br>  ○ L1 Regularization or Lasso Regression<br>    ■ Automatic Feature selection<br>• Comparison of Ridge and Lasso regression<br>• Elastic Net Regression – Intuition<br><br>**Logistic Regression** (A linear, discriminative, parametric classifier)<br>• Intuition and derivation<br>  ○ Regression for classification<br>  ○ "Squishing" between 0 and 1 using a non-linear activation function: The sigmoid<br>• A simple sentiment classifier<br>• Visualizing the logistic regression decision boundary<br>• Hyperplanes, linear and non-linear decision boundaries<br>• Cost function: Derivation of the cross-entropy loss function (log loss)<br>• Learning algorithm: Batch, Stochastic and Mini-batch Gradient Descent<br>• Multiclass (multinomial) classification: One-vs-all (one-vs-rest), One-vs-one<br>• The SoftMax activation function and multivariate log loss<br><br>*Tutorial Topics:*<br>• Review of Bias and Variance: detection and techniques to deal with them<br>• Review (and Expansion) of Feature Selection<br>• Review (and Expansion) of Sigmoid and Softmax: derivatives, visualizations, dealing with overflow<br>• Review of using Binary Classifiers to setup Multiclass Classification<br>• Assignment 2 | • Ben Taskar's under- and overfitting<br>• MLaPP: 1.4.7<br>• Andrew Ng's lecture – ML debugging<br>• Ben Taskar's Notes on Bias Variance<br>• Notes by Scott Foreman-Roe<br>• ELSII Chapter 2.9<br>• Murphy: 6.2.2<br><br>• SLP3 Ch5, ESLII Ch4<br>• Murphy 8, 13.3, 13.5.3<br>• Ben Taskar's notes<br>• TM chapter: Naive Bayes and Logistic Regression<br>• Nice blogpost on Gradient Descent, Adagrad, Newton's method | CLO2, CLO3, CLO4 | |
| 7. | **The Perceptron** (A discriminative, linear, parametric classifier)<br>• The McCulloch-Pitts Neuron and its limitations<br>• The Perceptron and its limitations<br>  ○ The Heaviside step function<br>  ○ Boolean functions: AND, OR and XOR!<br>  ○ One perceptron, two perceptrons, …<br>• Linear separability in low and high dimensional spaces<br>• From the step function to other activation functions<br>• The perceptron learning algorithm and its geometric interpretation<br>• Proof of convergence<br>  ○ Relation between margin and rate of convergence<br><br>**Maximum Margin Classifiers: Support Vector Machines (SVMs)** (A discriminative, linear/non-linear classifier)<br>• Intuition and motivation<br>  ○ The perceptron and the optimal separating hyperplane<br>• Hard Margin Linear Support Vector Machines: Derivation | • The Perceptron Wiki page<br>• Murphy 8.5.4<br><br>• Murphy: 14.5 - 14.5.2.2<br>• Ben Taskar's Notes on SVMs<br>• Kernel Cookbook by David Duvenaud<br>• Laurent El Ghaoui's lectures on duality<br>• "Idiot's guide to SVM" | CLO2, CLO3, CLO4 | |

| | | | | |
|---|---|---|---|---|
| | ○    Constrained optimization and Lagrange Multipliers<br>● Soft Margin Linear Support Vector Machines: Derivation<br>● Hinge-loss<br>● Kernels and Kernel SVMs<br><br>***Tutorial Topics:***<br>● Decision Boundaries and Hyperplanes of SVMs<br>● Review (and Expansion) of Soft Margin Classifiers and their Hyperparameters<br>● Review of Margins, Hinge Loss, and optimizing SVMs<br>● Review of Kernels: Linear, Polynomial, RBFs<br>● Supplementary: the Kernel Trick and Quadratic Programming | | | |
| 8. | **Neural Networks**<br>● The Neuron and linear decision boundaries. Can we do better?<br>    ○ Review Logistic regression and gradient descent<br>● Changing the representation of the data<br>    ○ Kernels<br>    ○ Neural networks<br>● Non-linear activation<br>● Multi-layer Perceptron<br>● Neural Networks<br>● Deep Learning<br>● Intuition: How NNs work?<br>    ○ Non-linear and complex decision boundaries<br>    ○ Universal approximation and regression<br>● Layers: Depth vs. width<br>● Gradient descent for NNs<br>    ○ Overfitting and the importance of Stochastic Gradient Descent (SGD)<br>● Formal Notation for Logistic Regression and NNs<br>    ○ Vectorizing LR and NNs<br>    ○ Forward propagation (LR, NN)<br>        ■ Simple, vectored with a single instance<br>        ■ Vectored with m instances<br>● Activation functions, gradients, and pros and cons<br>    ○ Sigmoid<br>    ○ Tanh<br>    ○ ReLU and Leaky ReLU<br>● Backward propagation (LR, NN)<br>    ○ Simple, vectored with a single instance and m instances<br>    ○ Training the NN<br><br>***Tutorial Topics:***<br>● Review of the Forward Pass and Backward Pass (in equations)<br>● Using Decision Boundaries to infer the architecture of a NN<br>● Review of strategies to deal with Overfitting and Underfitting<br>● Assignment 3<br>● Supplementary: programming demo on using MLPs for character-level Language Modeling | ● SLP Ch7, ESLII Ch11 | CLO2, CLO3, CLO4 | |
| 9. | **Sequence Models**<br>● Notion of Sequence Modeling tasks<br>    ○ Text and Speech (One-Hot )<br>    ○ Time-Series<br>    ○ Videos as sequences of images<br>● Types of Sequence Modeling problems and specifications<br>    ○ 1-to-Many - Image Captioning<br>    ○ Many-to-1 - Video Classification, Sentiment Classification<br>    ○ Many-to-Many (Seq2Seq) - Machine Translation, Summarization etc.<br>● Why do Feedforward NNs perform poorly at Sequence Modeling tasks?<br>● Introducing RNNs<br>    ○ Idea of the Hidden State<br>    ○ The overall architecture<br>    ○ "Unrolling" a unit<br>    ○ Intuition of (truncated) BPTT<br><br>● Challenges of training RNNs, and how to deal with them<br>    ○ Exploding/Vanishing Gradients | Andrej Karpathy's [The Unreasonable Effectiveness of RNNs](#)<br><br>Jay Alammar's [Visualizing A Neural Machine Translation Model](#)<br><br>SLP3 | CLO2, CLO3, CLO4 | |

| | | | |
|---|---|---|---|
| | <ul><li>○ Choice of Activation functions</li><li>○ Gradient Clipping</li><li>○ Changing the architecture (intro to LSTMs)</li></ul><ul><li>Introducing LSTMs</li><ul><li>○ Changes to the architecture</li><li>○ Idea of storing the "memory" in a cell</li></ul></ul><ul><li>Exploring Machine Translation</li><ul><li>○ Setup as a Seq2Seq problem</li><li>○ Embeddings</li><li>○ Using the Encoder-Decoder framework</li><li>○ Information bottleneck: passing on only one hidden state</li><li>○ Improvements</li><ul><li>■ Passing all the hidden states</li><li>■ Intuition of the Attention Mechanism</li></ul></ul></ul>***Tutorial Topics:***<ul><li>Review of the representations of text: Label Encoding, One-Hot Encoding, Embeddings</li><li>Programming demo on Embeddings in Python</li><li>Review of RNNs and LSTMs</li><li>Review of the Attention Mechanism in Machine Translation</li></ul> | | |
| 10. | **Attention Mechanism and Transformers**<ul><li>The Attention Mechanism in Machine Translation</li><li>Self Attention</li><ul><li>○ Dot Product Attention</li><li>○ The idea of contextualized word embeddings</li></ul><li>Introducing the Transformer and its contributions</li><ul><li>○ Highly Parallelizable</li><li>○ Contextualized Embeddings vs. Plain Embeddings</li><li>○ Long-Term Dependencies</li></ul><li>Transformer Architecture in a nutshell</li><ul><li>○ Role of an Encoder</li><li>○ Role of a Decoder and Masked Attention</li><li>○ Query, Keys, Values</li></ul><li>The Transformer in equations</li><ul><li>○ Positional Embeddings</li><li>○ Projections to QKV</li><li>○ Self Attention as Dot Product Attention</li><li>○ Role of Feedforward layers</li><li>○ Multi Headed Self Attention: the role and its equations</li><li>○ (ignore Layer Norms for brevity)</li></ul><li>Case Studies: BERT, T5, and GPT-3/LLMs</li></ul>***Tutorial Topics:***<ul><li>Review (and Expansion) of Self-Attention</li><li>Differences between the Encoder and Decoder: masking, role, architecture</li><li>Demo on coding out GPT in Python</li></ul> | Jay Alammar's [The Illustrated Transformer](#)<br><br>Andrej Karpathy's [Let's Build GPT](#)<br><br>SLP3 | |
| 11. | **Decision (Classification/Regression) Trees** (Discriminative, non-linear, parametric/non-parametric)<ul><li>Clustering using K-D Trees and nearest neighbor methods</li><ul><li>○ Why do we still need nearest neighbors?</li></ul><li>Revision of Trees</li><ul><li>○ Graph: Nodes, edges, directed/undirected, path, cyclic/acyclic,</li><li>○ Tree: A rooted directed acyclic graph (Rooted DAG)</li><ul><li>■ Parent, children, siblings, root, leaves, degree, height, arity, various relationships</li><li>■ Forests</li></ul></ul><li>From K-D trees to decision trees</li><li>Decision tree examples: classification and regression</li><ul><li>○ Categorical and real-valued attributes</li></ul><li>Bias and variance</li><ul><li>○ Tree height</li></ul><li>Growing tree automatically – Decision tree training</li><ul><li>○ The ID3 and CART algorithms</li></ul></ul> | <ul><li>[Decision Tree wiki page](#)</li><li>Ben Taskar's old [notes](#)</li><li>Murphy: 16.2</li><li>ESLII 8.7, ESLII Ch10, 15, 16</li></ul> | CLO3, CLO4 | |

|  | | | | |
|---|---|---|---|---|
| | <ul><li>Splits and purity/impurity<ul><li>Gini and Entropy</li><li>Information Gain and Information Gain Ratio</li><li>Multiclass</li><li>Regression: Variance</li></ul></li><li>Real-valued attributes</li></ul><ul><li>Strengths and weakness of CARTs<ul><li>Automatic feature selection</li><li>Generalizability of splitting on attributes with a large set of values</li><li>Missing data</li><li>Axis-aligned splits</li><li>Over fitting</li></ul></li><li>Side notes on Huffman coding</li><li>A deep dive into Entropy<ul><li>Quantifying expected surprise<ul><li>Events in isolation</li><li>Probability distributions</li></ul></li><li>Information, surprise, and uncertainty</li><li>Evaluating Language Models<ul><li>Comparing distributions using cross entropy</li></ul></li></ul></li><li>A discussion on Parametric/Non-parametric models</li></ul><br>**Ensemble methods: Bagging and Random Forests**<ul><li>Decomposition of Generalization Error<ul><li>Bias/Variance/Noise</li><li>Detecting high bias and high variance regimes</li></ul></li><li>Variance reduction<ul><li>The weak law of large numbers</li><li>$M$ datasets sampled from $P$</li><li>Bootstrapping<ul><li>Bootstrapped Aggregation (Bagging)</li></ul></li><li>Are we even drawing from the same $P$</li><li>Summary and advantages of bagging</li></ul></li><li>Random Forests<ul><li>Algorithm</li><li>Examples and benefits<ul><li>Out-of-box performance</li><li>The hyper parameters</li><li>$m$ and $k$</li><li>No need of normalization and feature-scaling</li><li>Resilience to the curse of dimensionality</li><li>Feature selection</li><li>Missing data and clustering</li></ul></li><li>Variants</li></ul></li></ul><br>*Tutorial Topics:*<ul><li>Review of Gini/Entropy and Decision Trees</li><li>Demo on Random Forests in Python: exploring Feature Importances for Feature Selection</li><li>Standardization in the context of Decision Trees</li><li>Decision Trees vs. Neural Networks</li></ul> | | | |
| 12. | **Ensemble methods: Boosting, Gradient Boosted Trees, and AdaBoost**<ul><li>Bias Reduction<ul><li>Intuition</li><li>Vectors</li><li>Gradient Descent in function space</li><li>Generic Boosting (AnyBoost)<ul><li>Algorithm and geometric interpretation</li></ul></li></ul></li><li>Gradient Boosted Regression Trees<ul><li>Algorithm</li><li>Detailed walk-through</li></ul></li><li>AdaBoost<ul><li>Setting</li><li>Odds Ratio and log-odds</li><li>Step size proportional to error reduction</li></ul></li></ul> | <ul><li>ESLII 8.7, ESLII Ch10, 15, 16</li></ul><ul><li>SLP Ch4</li><li>Murphy 2.2, 3.1-3.4</li><li>ESLII 6.6.3</li></ul> | CLO1, CLO2, CLO3, CLO4 | |

|  |  |  |  |  |
|---|---|---|---|---|
|  | ○    Instance weights proportional to (mis)classification and the say of the classifier<br>○    Algorithm and detailed walk-through<br>○    Properties and Summary<br><br>**Bayes Theorem**<br>●  Review of probability, joint and conditional probability, and derivation of the Bayes Theorem<br>●  Maximum a posteriori (MAP) and Maximum Likelihood Estimation (MLE)<br>○    Posterior, likelihood, prior and evidence<br>○    Classification using MAP and MLE<br>●  Example problems and solutions using Bayes Theorem<br>○    Binary and multiclass<br>○    Monty Hall problem, medical testing, Language Modeling<br>●  Generative and Discriminative classifiers<br>○    Solving SPAM vs. Not-SPAM |  |  |  |
| 13. | **The Naïve Bayes Classifier** (A linear, generative, parametric classifier)<br>●  Derivation and implementation<br>○    Classification using the Bayes Theorem<br>○    Learning by example: The SPAM vs. Not-SPAM problem<br>○    The "zeros" and how to get rid of them!<br>●  Independence, mutual exclusion, and conditional independence<br>○    The challenge of "how much of the context to use?" - Ngrams<br>○    Naïve Assumptions: Conditional Independence and Bag-of-Words<br>●  Data sparsity and Out-of-vocabulary (OOV) items<br>○    Laplace Add-1 smoothing<br>●  Another example: Sentiment analysis<br>●  Text generation using Naïve Bayes<br>○    Infinite monkeys on typewriters<br>○    The Shannon visualization method for Ngrams<br>○    Approximating Shakespeare and the Wall Street Journal<br>●  Real-valued features: Gaussian Naïve Bayes<br>●  Probability vs. likelihood<br>●  A worked example<br>●  Naïve Bayes decision boundary<br>●  Under assumptions and general case<br>●  Naïve Bayes: Strengths and weaknesses<br><br>*Tutorial Topics:*<br>●  Probability vs. Likelihood (and other terminologies)<br>●  Review of Naive Bayes<br>●  Review (and Expansion) of add-k smoothing<br>●  Demo on using Bigrams for Language Modeling | ●  Ben Taskar's notes on Naïve Bayes<br>●  TM chapter on Naive Bayes (ch 1-3)<br>●  Xiaojin Zhu's notes on Multinomial Naïve Bayes<br>●  Mannings' description of Multinomial Naive Bayes | CLO2, CLO3, CLO4 |  |
| 14. | **Unsupervised Learning**<br>●  The unsupervised learning setup<br>●  Use cases of unsupervised learning<br>○    Clustering<br>○    Anomaly detection<br>○    Feature selection and dimensionality reduction<br>●  Types of clustering<br>○    Monothetic and Polythetic<br>○    Hard and Soft<br>○    Flat and Hierarchical<br>●  Clustering<br>○    K-D Trees<br>■  Monothetic, hard boundary, hierarchical, divisive (top-down)<br>■  Motivation<br>■  Algorithm<br>○    Vector Quantization<br>■  Motivation and method<br>■  Codebook and distance metric<br>■  Euclidean and Mahalanobis distances<br>○    K-means<br>■  Polythetic, hard boundary, flat<br>■  Lloyd/Forgy method | ●  ESLII 8.7, ESLII Ch10, 15, 16 | CLO3, CLO4 |  |

|  |  |  |  |  |
|--|--|--|--|--|
| | <ul><li>■ Expectation Maximization (EM) and K-means</li><li>■ The K-means objective</li><li>■ Optimal Number of Clusters</li><li>■ Categorical data and K-modes</li><li>■ Vector Quantization using K-means</li></ul><ul><li>○ Evaluating Clustering</li><ul><li>■ Extrinsic and Intrinsic evaluation</li></ul><li>○ Gaussian Mixture Models</li><ul><li>■ Polythetic, soft boundary, flat, probabilistic</li><li>■ K-means vs. GMMs</li><li>■ EM for GMMs</li><li>■ Mixture models in 1-dimension and n-dimensions</li><li>■ Likelihoods, cluster assignments, and cluster update rules</li><li>■ The covariance matrix</li><li>■ How many Gaussians?</li></ul><li>○ Hierarchical clustering</li><ul><li>■ Recursive K-means</li><ul><li>● Polythetic, hard boundary, hierarchical, top-down</li></ul><li>■ Agglomerative Clustering</li><ul><li>● Polythetic, hard boundary, hierarchical, bottom-up</li><li>● Examples</li><li>● Distances: Single link, complete link, average link, centroids, Ward's method</li></ul></ul></ul><br><ul><li>● Dimensionality Reduction</li><ul><li>○ Feature selection vs. feature reduction</li><li>○ Motivation</li><ul><li>■ Visualization</li><li>■ Redundant and correlated features</li><li>■ Real vs. apparent dimensionality</li><li>■ The curse of dimensionality</li></ul><li>○ Principal Component Analysis (PCA)</li><ul><li>■ The dimension of greatest variability</li><li>■ A side-note on Matrices, Linear Transformations, the determinant, Eigenvalues and Eigenvectors</li><li>■ The Eigenvectors of the Covariance Matrix</li><li>■ How many dimensions?</li><li>■ Strengths and weaknesses</li></ul><li>○ Linear Discriminant Analysis (LDA)</li><ul><li>■ Supervised setup</li><li>■ Discrimination vs. spread</li></ul></ul><li>● Anomaly Detection</li><ul><li>○ Anomalies</li><li>○ How do we define anomalies?</li><li>○ Why unsupervised?</li><li>○ Examples and challenges</li><li>○ Detection</li><ul><li>■ One Class Classification</li><li>■ Density estimation</li><li>■ One feature and multiple features</li><li>■ Algorithm</li><li>■ Example</li><li>■ Evaluation</li><li>■ Unsupervised vs. supervised</li></ul></ul></ul><br>**Big Challenges and Opportunities in AI and ML**<ul><li>● The case for Explainable AI</li><li>● The case for Fair AI</li><li>● Societal biases</li><li>● Imbalanced classification</li><li>● Machine Learning for Development (ML4D)</li></ul><br>*Tutorial Topics:*<ul><li>● Review of Clustering</li><li>● Review of PCA</li></ul> | | | |

| | | | |
|---|---|---|---|
| | ● Programming demo of PCA In Python as a preprocessing step | | |
| colspan="4" | **Other topics – to be covered if we have time** |
| 15. | **Bayes: Advanced topics** (supplementary)<br>● Hypothesis spaces<br>● Frequentist viewpoint<br>   ○ Intuition, derivation, pros and cons, extreme data<br>● Bayesian viewpoint<br>   ○ Intuition, derivation, MAP, pros and cons<br>   ○ Conjugate priors: The Beta and Dirichlet distributions<br>● Comparison of MAP and MLE<br>   ○ Laplace smoothing<br>● The Bayes Optimal Classifier | ● SLP Ch4<br>● Murphy 2.2, 3.1-3.4<br>● ESLII 6.6.3 | CLO2, CLO3, CLO4 | |
| 16. | **Graphical Sequence Processing Models**<br>● Hidden Markov Models (HMMs)<br>● Maximum Entropy Markov Models (MEMMs)<br>● Undirected Graphical Models (Markov Random Fields)<br>● Conditional Random Fields (CRFs)<br>● Directed Graphical Models (Bayes Nets) | ● SLP A, ESLII Ch17 | CLO3, CLO4 | |

---

**Textbook(s)/Supplementary Readings**

**Text Books**
• Machine Learning, Tom Mitchell, McGraw Hill, 1997 – TM
• The Elements of Statistical Learning: Data mining, Inference, and Prediction, Hastie, Trevor, Robert Tibshirani, and Jerome Friedman, Springer Science & Business Media, 2009 – ESLII

**Reference Books**
• Speech and Language Processing by Jurafsky and Martin, Ed 3 (online draft) – SLP
• Machine Learning: A Probabilistic Perspective, Murphy, Kevin P. MIT press, 2012 – Murphy.
• Pattern Recognition and Machine Learning, Christopher M. Bishop, Springer, 2006 – Bishop.
• Introduction to Machine Learning, Ethem Alpaydin, Ed 2, MIT Press, 2010 – Alpaydin.
• Deep Learning, Ian Goodfellow and Yoshua Bengio and Aaron Courville, 2016 – Goodfellow

---

**Course policies**

**Use of electronic devices (e.g., mobile phones and laptops) in the class is strictly forbidden.** A violation could result in deduction of marks and other strict penalties

**Late arrival:** You may not be allowed in the class 10 minutes after the start time

**Plagiarism:** All work MUST be done independently. In certain assignments students will be allowed to have discussions with peers, in which case they must mention the name and roll number of the student with whom the discussion took place and the nature of the discussion. Even in those assignments, all implementations need to be done independently. Any plagiarism or cheating of work from others or the internet will be immediately referred to the DC. If you are confused about what constitutes plagiarism, it is YOUR responsibility to consult with the instructor or the TA in a timely manner. No "after the fact" negotiations will be possible.

● Submitting someone else's assignment as your own "by mistake" would count as plagiarism. If this indeed happens accidentally, please let us know immediately (within minutes) along with an explanation and do not wait until we find it out on our own. In the latter case, it would be considered plagiarism.

**Quizzes:** Quizzes will be unannounced. We will be following an n-x (x=2) policy for the quizzes. <u>There is no makeup for a missed quiz</u>. If you have missed up to x quizzes, you will be covered only using the n-x policy (even if you have an approved petition with the OSA). If you have missed more than x quizzes, then you would be awarded the average marks (across all the quizzes that you attempted) for each missed quiz, provided that your case has been approved by the Office of Student Affairs.

**Non-uniform weightage:** All subcomponents (e.g., quizzes, assignments) may not carry the same weight. These weights may not be announced prior to the submission of the components and will be determined by the course instructor based on factors including (but not limited to) the length, difficulty level, amount of help available, etc. for each subcomponent.

**Programming:** Strong programming skills are expected for this course. Please keep in mind that this is a programming intensive course, and you will be spending a lot of time designing and coding up your solutions.

**Assignments:** There is negative marking for skipped assignments and there is no n-x policy for assignments. Assignments are a basic building block of this course, and it will be ensured that students, who pass the course, have significant hands-on experience.

● You will be awarded 0 marks or investigated for plagiarism for submitting incorrect/corrupted files and/or older assignments. We will not accept resubmissions in these cases even if the system date shows that the file was not modified after the deadline.

● You are allowed 5 grace days for the entire semester. No late submission of assignments is allowed after your grace days have expired. We do not have any deduction policy for late submissions in addition to the grace days. All grace days must be utilized before the start of the dead week and any remaining grace days will expire as soon as the dead week begins.

● Please do not wait until the last moment to submit assignments and other components. Any requests to accommodate late submissions due to last minute issues (submission of partial or incorrect files, assignment server down-time, internet and power failures, personal problems, etc.) would not be accommodated.

| BLOOM's TAXONOMY* | |
|---|---|
| 1 - Remember<br>2 - Understand<br>3 - Apply<br>4 - Analyze<br>5 - Evaluate<br>6 - Create | ● Recall facts and basic concepts<br>● Explain ideas or concepts<br>● Use information in new situations<br>● Draw connection among ideas<br>● Justify a stand or decision<br>● Produce new or original work |

https://cft.vanderbilt.edu/guides-sub-pages/blooms-taxonomy/

## Appendix B

## ACM Dispositions Table - I

| ACM Dispositions | | | |
|---|---|---|---|
| Element | Elaboration | Element | Elaboration |
| D1 Adaptable: | Flexible; agile, adjust in response to change | D7 Professional: | Professionalism, discretion, ethical, astute |
| D2 Collaborative: | Team player; willing to work with others | D8 Purpose-driven: | Goal driven, achieve goals, business acumen |
| D3 Inventive: | Exploratory; Look beyond simple solutions | D9 Responsible: | Use judgment, discretion, act appropriately |
| D4 Meticulous: | Attentive to detail; thoroughness, accurate | D10 Responsive: | Respectful; react quickly and positively |
| D5 Passionate: | Conviction, strong commitment, compelling | D11 Self-directed: | Self-motivated, determination, independent |
| D6 Proactive: | With initiative, self-starter, independent | | |

## ACM Dispositions Table - II

| Class Assessments and Proposed Dispositions | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Assessment Type | D1 Adaptable | D2 Collaborative | D3 Inventive | D4 Meticulous | D5 Passionate | D6 Proactive | D7 Professional | D8 Purpose-driven | D9 Responsible | D10 Responsive | D11 Self-directed | Included |
| Quiz | | | | ✓ | | | ✓ | | ✓ | | | Yes |
| Assignment-Individual | | | ✓ | ✓ | | | ✓ | | ✓ | | | Yes |
| Assignment-Group | | ✓ | ✓ | ✓ | | | ✓ | | ✓ | ✓ | | Yes |
| Project-Individual | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | Yes |
| Project-Group | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | Yes |
| Presentation-Individual | | | | ✓ | | | ✓ | | ✓ | ✓ | ✓ | Yes |
| Presentation-Group | | ✓ | | ✓ | | | ✓ | | ✓ | ✓ | | Yes |
| Labs-Individual | | | ✓ | ✓ | | | ✓ | | ✓ | | | Yes |
| Labs- Group | | ✓ | ✓ | ✓ | | | ✓ | | ✓ | ✓ | | Yes |
| Exams | | | | ✓ | | | ✓ | | ✓ | | | Yes |
| Included | Yes | Yes | Yes | Yes | Yes | Yes | Yes | Yes | Yes | Yes | Yes | |

# Appendix C
## ACM Computing Knowledge Landscape Table

| ACM Computing Knowledge Landscape (CK) | | | |
|---|---|---|---|
| **1.** <br><br> **Users and Organizations** | CK1.1: Social Issues and Professional Practice <br><br> CK1.2: Security Policy and Management <br><br> CK1.3: IS Management and Leadership <br><br> CK1.4: Enterprise Architecture <br><br> CK1.5: Project Management <br><br> CK1.6: User Experience Design | **4.** <br><br> **Software Development** | CK4.1: Software Quality, Verification and Validation <br> CK4.2: Software Process <br> CK4.3: Software Modeling and Analysis <br> CK4.4: Software Design <br> CK4.5: Platform-Based Development |
| **2.** <br><br> **Systems Modeling** | CK2.1: Security Issues and Principles <br><br> CK2.2: Systems Analysis & Design <br><br> CK2.3: Requirements Analysis and Specification <br><br> CK2.4: Data and Information Management | **5.** <br><br> **Software Fundamentals** | CK5.1: Graphics and Visualization <br> CK5.2: Operating Systems <br> CK5.3: Data Structures, Algorithms and Complexity <br> CK5.4: Programming Languages <br> CK5.5: Programming Fundamentals <br> CK5.6: Computing Systems Fundamentals |
| **3.** <br><br> **Systems Architecture and Infrastructure** | CK3.1: Virtual Systems and Services <br><br> CK3.2: Intelligent Systems (AI) <br><br> CK3.3: Internet of Things <br><br> CK3.4: Parallel and Distributed Computing <br><br> CK3.5: Computer Networks | **6.** <br><br> **Hardware** | CK6.1: Architecture and Organization <br> CK6.2: Digital Design <br> CK6.3: Circuits and Electronics <br> CK6.4: Signal Processing |