

San Jose State University
Department of Computer Science
CS 185C, Introduction to Machine Learning with Applications in
Information Security, Fall 2018

- **Course and Contact information**

- **Instructor:** Mark Stamp
- **Office Location:** MH 216
- **Telephone:** 408-924-5094
- **Email:** mark.stamp@sjsu.edu
- **Office hours:** Wednesday, 10:00am - noon
- **Class Days/Times:** Tuesday & Thursday, 10:30am - noon
- **Classroom:** MH 233
- **Prerequisites:** TBD
- **Note:** This class, in conjunction with CS 166, can be used to satisfy the "deep course" requirement. Also, this course has been approved as a graduate elective.

- **Course Description**

- Topics in machine learning. The following machine learning techniques are covered in detail: Hidden Markov Models (HMM), Profile Hidden Markov Models (PHMM), Principal Component Analysis (PCA), Support Vector Machines (SVM), and clustering. Illustrative applications of each of these major topics are provided, with most of the applications drawn from the field of information security. In addition, the course will include an overview of each of the following topics: k-Nearest Neighbor, Neural Networks, Boosting/AdaBoost, Random Forests, Linear Discriminant Analysis, Naive Bayes, Regression Analysis, Conditional Random Fields, and Data Analysis. Prerequisite: TBD.

- **Learning Outcomes**

- The focus of this course will be machine learning, with illustrative applications drawn primarily from the field of information security. After completing this course students should have a working knowledge of a wide variety of machine learning topics, and have a good understanding of how to apply such techniques to real-world problems.

- **Required Texts/Readings**

- The primary text will be [Machine Learning with Applications in Information Security](https://www.crcpress.com/Introduction-to-Machine-Learning-with-Applications-in-Information-Security/Stamp/p/book/9781138626782) (<https://www.crcpress.com/Introduction-to-Machine-Learning-with-Applications-in-Information-Security/Stamp/p/book/9781138626782>), by Mark Stamp, published by Chapman Hall/CRC in 2017. This book covers several machine learning techniques in detail, and includes a large number of illustrative applications. Many of the applications are from information security, including a variety of topics related to malware, intrusion detection (IDS), spam, and cryptanalysis, among others.
- Additional relevant material:
 - [PowerPoint slides](http://www.cs.sjsu.edu/~stamp/ML/powerpoint) at <http://www.cs.sjsu.edu/~stamp/ML/powerpoint>
 - Current semester [lecture videos](http://www.cs.sjsu.edu/~stamp/ML/lectures/CS185C_Fall18/) are available at http://www.cs.sjsu.edu/~stamp/ML/lectures/CS185C_Fall18/. If you are asked to login to access the videos, both the username and password are "infosec". **Note:** The instructor hereby gives students permission to record his lectures (audio and/or video). At least with respect to this class, your instructor has nothing to hide.

- Class-related discussion will be posted on [Piazza](https://piazza.com/class/jkpence3iqu12p) at <https://piazza.com/class/jkpence3iqu12p>. You are strongly encouraged to participate by asking questions, as well as by responding to questions that other students ask. At the start of the semester, you should receive an email asking you to join this discussion group—if not, contact your instructor via email.
 - The applications parts of this course are essentially self-contained, but for additional background information on the security-related topics, the following resources are recommended.
 - *Computer Viruses and Malware*, John Aycock, Springer 2006. Many of the applications we discuss are related to malware. Aycock's book is easy to read and in spite of being fairly old, it provides a good foundation for malware research.
 - *Information Security: Principles and Practice*, Mark Stamp, Wiley 2011. If you have not taken CS 265, you should do so. You can refer to this fine book if you have questions about security-related topics during this course.
 - [Open Malware](http://www.offensivecomputing.net/) (at <http://www.offensivecomputing.net/>) includes a large collection of samples of live malware.
 - [VX Heavens](http://vx.netlux.org/) (at <http://vx.netlux.org/>) is a source for "hacker" type of information on viruses. Malware samples are also available.
 - [Journal of Computer Virology and Hacking Techniques](http://www.springer.com/computer/journal/11416) (at <http://www.springer.com/computer/journal/11416>) is a journal for malware-specific research papers. There are also several good conferences that focus on malware and/or machine learning applications in information security.
 - [Recent masters project reports](http://www.cs.sjsu.edu/~stamp/cv/mss.html#masters) (at <http://www.cs.sjsu.edu/~stamp/cv/mss.html#masters>). Most of these projects involve applications of machine learning to malware or other topics in information security.
- **Course Requirements and Assignments**
 - SJSU classes are designed such that in order to be successful, it is expected that students will spend a minimum of forty-five hours for each unit of credit (normally three hours per unit per week), including preparing for class, participating in course activities, completing assignments, and so on. More details about student workload can be found in [University Policy S12-3](http://www.sjsu.edu/senate/docs/S12-3) at <http://www.sjsu.edu/senate/docs/S12-3.pdf>.
 - Schedule
 - Week 1 --- Introduction and overview
 - Week 2 --- Hidden Markov Models
 - Week 3 --- Data Analysis
 - Week 4 --- Applications of Hidden Markov Models
 - Week 5 --- Profile Hidden Markov Models
 - Week 6 --- Applications of Profile Hidden Markov Models
 - Week 7 --- Principal Component Analysis
 - Week 8 --- Applications of Principal Component Analysis
 - Week 9 --- Support Vector Machines
 - Week 10 --- Applications of Support Vector Machines
 - Week 11 --- Clustering
 - Week 12 --- Clustering Applications
 - Week 13 --- k-Nearest Neighbor, Neural Networks, Boosting/AdaBoost, Random Forests
 - Week 14 --- Linear Discriminant Analysis, Naive Bayes, Regression Analysis, Conditional Random Fields
 - Week 15 --- Project presentations
 - Homework is due *typewritten* (include source code, but not executable files) by class starting time on the due date. Each assigned problem requires a solution and an explanation and work detailing how you arrived at your solution. Cite any outside sources used to solve a problem. When grading an assignment, I may ask for additional information. Note that a *subset* of the assigned problems

will typically be graded.

Homework must be submitted via email before the start of class on the due date. Be sure to have an extra copy of your homework with you in class, and be prepared to discuss your solutions. Your written solutions must be in a pdf file. Submit any source code or other attachments in separate files (i.e., no code in the solution itself). You must provide enough discussion of your solution so that the grader can understand your solution, and so that the grader can be sure that you understand your solution. Put your written solution and any relevant source code in a folder named "yourlastname". Then zip your homework folder and submit the file yourlastname.zip via email to CS185C.FALL18@gmail.com. The subject line of your email *must* be of the form:

CS185CHMK assignmentnumber yourlastname last4digitofyourstudentnumber

The subject line must consist of the four identifiers listed. There is no space within an identifier and each identifier is separated by a space.

- Assignment 0: Due **Tuesday, August 28**
 For this assignment, turn in a hardcopy of your solutions at the start of class.
 1. Read [A Revealing Introduction to Hidden Markov Models](https://www.cs.sjsu.edu/~stamp/RUA/HMM.pdf) (at <https://www.cs.sjsu.edu/~stamp/RUA/HMM.pdf>) and do the following.
 - a. Briefly (1 paragraph) summarize how an HMM is trained.
 - b. How is a trained HMM used to score a sequence?
 - c. Very briefly explain how an HMM and dynamic program differ.
 - d. Why is it necessary to scale the values of the matrices when training an HMM?
 2. Read the article "Models will rule the world" (I will give you a hardcopy of the article on the first day of class) and do the following.
 - a. In one paragraph, summarize the authors' main points.
 - b. Write second paragraph discussing what you most agree with and anything that you disagree with in this article.

- Assignment 1: Due **Tuesday, September 4**
 Chapter 2, problems 1, 2, 3, 10. For problem 10 you must use HMM code that you have written entirely on your own, following the algorithms given in your textbook.

- Assignment 2: Due **Thursday, September 20**
 Chapter 2, problems 11, 14, and 15. You will receive 10 points extra credit if you use your own HMM code to solve at least 2 of these 3 problems. If your code is too slow and/or buggy, use this [reference HMM implementation](https://www.cs.sjsu.edu/~stamp/RUA/HMM_ref.zip) (at https://www.cs.sjsu.edu/~stamp/RUA/HMM_ref.zip)
 Chapter 8, problems 1, 6, 7, 8, 9, 10.

- Assignment 3: Due ~~Thursday, September 27~~ **Friday, September 28**
 Chapter 3, problems 3, 4, 5a, 11. Graduate students must also do problem 7.

- Assignment 4: Due **Thursday, October 11**
 Chapter 4, problems 3, 4, 9, 10, 11, 13. Graduate students must also do problem 17.

- Assignment 5: Due **Thursday, October 25**
 Chapter 5, problems 1, 3, 4, 6, 9, 12, 15. Extra credit: 16.

- Assignment 6: Due ~~Thursday, November 8~~ **Friday, November 9**
[Deep Learning Problems](https://www.cs.sjsu.edu/~stamp/ML/files/ann.pdf) (<https://www.cs.sjsu.edu/~stamp/ML/files/ann.pdf>), 1, 2, 6, 10a. For

problem 10a, let $C(O) = \sum(\log(c_t))$, rather than $\sum(c_t)$.

- Assignment 7: Due **Friday, December 7**
Chapter 6, problems 4, 5, 7, 8, 13, 15, 16.
Chapter 7, problems 1, 4, 7, 15, and problem 1 from [here](https://www.cs.sjsu.edu/~stamp/ML/files/ada.pdf) (at <https://www.cs.sjsu.edu/~stamp/ML/files/ada.pdf>)
- Assignment 8: Due **TBD**
- Assignment 9: Due **TBD**

- NOTE that [University policy F69-24](http://www.sjsu.edu/senate/docs/F69-24.pdf) at <http://www.sjsu.edu/senate/docs/F69-24.pdf> states that "Students should attend all meetings of their classes, not only because they are responsible for material discussed therein, but because active participation is frequently essential to insure maximum benefit for all members of the class. Attendance per se shall not be used as a criterion for grading."

• Grading Policy

- Test 1, 100 points. Date: **Tuesday, October 30**.
- Homework, quizzes, class participation and other work as assigned, 100 points. A subset of the assigned problems will be graded.
- [Machine Learning Project](#), 100 points. You must send your project proposal to me (via email) by **Monday, September 24**. A written project report is due ~~Tuesday, November 27~~ **Friday, November 30**. The project presentations will begin on **Tuesday, November 27**.
- Final, 100 points. Date: **Wednesday, December 12 at 9:45 am**. The official finals schedule is here: <http://info.sjsu.edu/static/catalog/final-exam-schedule-fall.html>
- Semester grade will be computed as a weighted average of the major scores listed above.
- **No** make-up tests or quizzes will be given and **no** late homework or project (or other work) will be accepted.
- Grading Scale:

Percentage	Grade
92 and above	A
90 - 91	A-
88 - 89	B+
82 - 87	B
80 - 81	B-
78 - 79	C+
72 - 77	C
70 - 71	C-
68 - 69	D+
62 - 67	D
60 - 61	D-
59 and below	F

- Note that "All students have the right, within a reasonable time, to know their academic scores, to review their grade-dependent work, and to be provided with explanations for the determination of

their course grades." See [University Policy F13-1](http://www.sjsu.edu/senate/docs/F13-1.pdf) at <http://www.sjsu.edu/senate/docs/F13-1.pdf> for more details.

- **Guest Lectures**

- Sravani Yajamanam, Automated Driving Research Engineer at Ford Motor Company
 - Date: ~~November 15~~ Tuesday, November 13
 - Time: 10:30am
 - Location: MH 233
 - Topic: Deep Learning & Big Data Analytics
 - Abstract: Typically in machine learning, we get a single number representing a model (like 91% accuracy) but we don't have a good representation of the data that failed. For example, for a traffic sign classification problem, how many times was a 'do not enter' sign misclassified as a stop sign? In other words, how can we evaluate the performance of deep learning models? How can we pull interesting statistics about large datasets and visualize them efficiently? My summer internship project at Ford tackled this problem of evaluating the performance of deep learning models by using big data analytics.
- Paco Guzman, Facebook
 - Date: Tuesday, November 20
 - Time: 10:30am
 - Location: E 189 (Engineering Auditorium)
 - Topic: TBD
 - Abstract: TBD

- **Classroom Protocol**

- Keys to success: Do the homework, complete a good project, and attend class
- **Wireless laptop is required.** Your laptop must remain closed (preferably in your backpack and, in any case, not on your desk) until I inform you that it is needed for a particular activity
- **Cheating** will not be tolerated, but working together is encouraged
- Student must be respectful of the instructor and other students. For example,
 - No disruptive or annoying talking
 - Turn off cell phones
 - Class begins on time
 - Class is not over until I say it's over
- Valid picture ID required at all times
- The last day to drop without a "W" grade is **Friday, August 31**, and the last day to add is **Monday, September 10**

- **University Policies**

- Office of Graduate and Undergraduate Programs maintains university-wide policy information relevant to all courses, such as academic integrity, accommodations, etc. You may find all syllabus related University Policies and resources information listed on GUP's [Syllabus Information web page](http://www.sjsu.edu/gup/syllabusinfo/) at <http://www.sjsu.edu/gup/syllabusinfo/>