

Topics in Sequence-Based Machine Learning for Bioinformatics Section 01

CS 225

Spring 2024 3 Unit(s) 01/24/2024 to 05/13/2024 Modified 01/21/2024

Contact Information

Instructor(s): William "Bill" Andreopoulos

Office Location: Online (former MacQuarrie Hall 416)

Telephone: (408) 924 5085

Email: william.andreopoulos@sjsu.edu

Office Hours: Friday 12:00-2:00 pm Online via Zoom

Class Days/Time: Monday and Wednesday 10:30-11:45am

Classroom: MQH 225

Course Information

Course Format

This course adopts an in-person classroom delivery format.

Faculty Web Page and MYSJSU Messaging

Course materials such as syllabus, handouts, notes, assignment instructions, etc. can be found on Canvas Learning Management System course login website at <http://sjsu.instructure.com>. You are responsible for regularly checking with the course messaging system to learn of any updates. You should modify the Canvas settings for notifications of announcements and Slack messages to be sent to you.

Course Description and Requisites

A study of recent advances in machine learning methods with applications to solving sequence analysis problems in molecular biology. The methods examined include word embeddings, vector space representations, language models, and deep learning architectures. A substantial course project is required.

Prerequisite(s): BIOL 123B and MATH 162, or CS156, or CS171, or instructor consent. Graduate standing. Allowed Declared Major: Computer Science MS, Bioinformatics MS, and Data Science MS.

Letter Graded

* Classroom Protocols

Communication with the instructor

Students should follow the correct channels for communication. Questions should preferably be done during the regular class meeting time (in person or via Zoom) or office hours. For course-related electronic communication students should use the Discord channel:

1) We will be using the course Discord channel for class discussion. The system is catered to getting you help efficiently from classmates, the TA, embedded tutor, and the instructor. Rather than emailing redundant questions to the teaching staff, students should post questions on the Discord channel where the entire class can read and benefit from the responses. The professor may re-post questions that are of general interest to the general channel or discuss them in class. The professor may ask students to reveal their real name if they are making special requests on Discord (e.g. deadline extensions) to prevent abuse.

2) Students are invited to join the office hours.

Private messages sent to the instructor's other email addresses get lost due to the large volume of emails received.

The instructor does not write messages after normal business hours, on weekends or holidays.

Reviewing code for the homework and technical trouble-shooting should be done during the office hours.

Never email your entire code for an assignment to the instructor. The instructor will not fix all the bugs in your code. Limit the code you post to 20 lines or less.

Announcements that concern everyone, such as reminders about due dates or class policy, will be posted.

Class Attendance

Attendance (in-person or via Zoom) is highly recommended. Classes will be recorded as Zoom screencasts and posted on Canvas. Students are responsible for all material presented in all classes.

Regrading Procedure

Grades assigned are final, unless there was an error in the grading. If a student wants to request a regrade of a homework or test, please follow instructions on the "Regrade request" page on Canvas. A request for a regrade is not a technique to drum up a few more points. If the course instructor thinks a component was scored too generously the first time, it may be lowered in a regrade. Thus, regrading may result in a lower grade.

Classroom Protocol

Students on Zoom should be muted when not speaking, and dressed appropriately when their camera is on.

Course material developed by the instructor is the intellectual property of the instructor. Students can not publicly share or upload instructor generated material for this course such as exam questions, lecture notes, hands-on exercises or homework solutions without instructor permission.

Program Information

Diversity Statement - At SJSU, it is important to create a safe learning environment where we can explore, learn, and grow together. We strive to build a diverse, equitable, inclusive culture that values, encourages, and supports students from all backgrounds and experiences.

Course Learning Outcomes (CLOs)

Upon successful completion of this course, students will be able to:

1. Use machine learning and deep learning in bioinformatics sequence analysis to answer biological questions and to generate biological hypotheses.
2. Comprehend the nature, scope and limits of using machine learning and deep learning in the field of bioinformatics.
3. Develop machine learning and deep learning solutions for sequence data.
4. Compare different machine learning algorithms and choose a solution based on suitability for a particular data set.
5. Compare biomolecular analysis with machine learning to analysis with classical bioinformatics tools.
6. Appreciate some of the most challenging problems in life sciences that use machine learning methods, possess insight into how to solve those problems.

Course Materials

Texts/Readings

We don't use a specific textbook in this class as there exists a lot of relevant material on bioinformatics found in various references. The reading material will be the slides, references and handouts.

A copy of my slides will be available to the students enrolled in the class.

Additional handouts will be provided through Canvas.

Major references:

- Data Analytics in Bioinformatics: A Machine Learning Perspective, 1st Edition (2021). by Rabinarayan Satpathy, Tanupriya Choudhury, Suneeta Satpathy, Sachi Nandan Mohanty, Xiaobo Zhang (Editors). ISBN-13: 978-1119785538.
- Haoyang Li, Shuye Tian, Yu Li, Qiming Fang, Renbo Tan, Yijie Pan, Chao Huang, Ying Xu, Xin Gao. Modern deep learning in bioinformatics. *Journal of Molecular Cell Biology*, Volume 12, Issue 11, November 2020, Pages 823–7.
- Deep learning in bioinformatics. Edited by Xin Gao, Wei Wang. Elsevier Methods. Volume 166, 15 August 2019, Pages 1-120.
- Walsh, Ian; Pollastri, Gianluca; Tosatto, Silvio C. E. (September 2016). "Correct machine learning on protein sequences: a peer-reviewing perspective". *Briefings in Bioinformatics*. 17(5): 831–840.
- Chicco, D (December 2017). "Ten quick tips for machine learning in computational biology". *BioData Mining*. 10 (35): 35.
- Yang, Yuedong; Gao, Jianzhao; Wang, Jihua; Heffernan, Rhys; Hanson, Jack; Paliwal, Kuldeep; Zhou, Yaoqi (May 2018). "Sixty-five years of the long march in protein secondary structure prediction: the final stretch?". *Briefings in Bioinformatics*. 19 (3): 482–494.
- Wang, Sheng; Peng, Jian; Ma, Jianzhu; Xu, Jinbo (January 2016). "Protein secondary structure prediction using deep convolutional neural fields". *Scientific Reports*. 6: 18962.

Other technology requirements / equipment / material

Students will use colab.research.google.com and create Jupyter notebooks in Python to ensure their work is shareable and reproducible.

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.

Reading assignments: Readings will regularly be assigned for the next class (see schedule). Slides will be posted under the Canvas modules before the next class.

Hands-On Worksheets: We will have a number of hands-on worksheets. The hands-on worksheets will involve use of bioinformatics tools. The purpose of the hands-on exercises is to develop your understanding of the material and skills in using the tools.

The Hands-On worksheets will involve learning how to use machine learning and deep learning tools with the Python programming language for performing bioinformatics analysis. Students will use colab.research.google.com and create Jupyter notebooks in Python to ensure their work is shareable and reproducible.

Term Project and In-Class Presentation: There will be a term project. It is a group project. Each group consists of two students. A list of possible projects will be provided to you by the instructor.

Team Formation is due on Monday, February 7, 2024.

A Progress Report is due on Monday, April 8, 2024 (after Spring Recess).

The final project is due on Monday, May 13, 2024.

The in-class presentations will also take place from May 6-13, 2024.

A grading rubric will be provided.

All homework should be submitted on Canvas, not by e-mail.

Examinations

The midterm exams are each one hour and fifteen minutes long. The final exam is two hours and fifteen minutes long.

The exams will contain multiple choice questions, true/false and short answer questions. Exams are *open book*, *open notes*, and comprehensive. The exams should be done individually and are not group work. No make-up exams except in case of verifiable emergency circumstances.

Presentation of a research paper: Each student should present an influential research paper of his/her choice, which is related to their project topic, to one of the classes. Students should sign up in the given spreadsheet for a date to present a paper. The paper, chosen by the student, should either use machine learning/deep learning towards making a biological discovery or introduce a novel tool for Natural Language Processing or text mining or bioinformatics. The presentation should last for no more than 10 minutes followed by Q&A. A grading rubric will be provided.

Participation during class via Zoom: The polling questions are in the form of multiple-choice and true-false questions. All students are expected to participate with Zoom polling. Credit is given based on participation and it is not necessary to get the correct answer in polls to get credit. Please contact eCampus at ecampus@sjsu.edu with any questions or issues with the Zoom technology.

✓ Grading Information

The course grade is based on:

Hands-On Worksheets: 10%

Midterms: 20%

Final: 20%

Project: 40%

Presentation of a research paper: 10%

<i>Grade</i>	<i>Points</i>	<i>Percentage</i>
<i>A plus</i>	<i>960 to 1000</i>	<i>96 to 100%</i>
<i>A</i>	<i>930 to 959</i>	<i>93 to 95%</i>
<i>A minus</i>	<i>900 to 929</i>	<i>90 to 92%</i>
<i>B plus</i>	<i>860 to 899</i>	<i>86 to 89 %</i>
<i>B</i>	<i>830 to 859</i>	<i>83 to 85%</i>
<i>B minus</i>	<i>800 to 829</i>	<i>80 to 82%</i>
<i>C plus</i>	<i>760 to 799</i>	<i>76 to 79%</i>
<i>C</i>	<i>730 to 759</i>	<i>73 to 75%</i>
<i>C minus</i>	<i>700 to 729</i>	<i>70 to 72%</i>
<i>D plus</i>	<i>660 to 699</i>	<i>66 to 69%</i>

<i>Grade</i>	<i>Points</i>	<i>Percentage</i>
<i>D</i>	<i>630 to 659</i>	<i>63 to 65%</i>
<i>D minus</i>	<i>600 to 629</i>	<i>60 to 62%</i>

University Policies

Per [University Policy S16-9 \(PDF\)](http://www.sjsu.edu/senate/docs/S16-9.pdf) (<http://www.sjsu.edu/senate/docs/S16-9.pdf>), relevant university policy concerning all courses, such as student responsibilities, academic integrity, accommodations, dropping and adding, consent for recording of class, etc. and available student services (e.g. learning assistance, counseling, and other resources) are listed on the [Syllabus Information](https://www.sjsu.edu/curriculum/courses/syllabus-info.php) (<https://www.sjsu.edu/curriculum/courses/syllabus-info.php>) web page. Make sure to visit this page to review and be aware of these university policies and resources.

Course Schedule

Week	Topic
01/29-01/31	Introduction, overview of unsupervised and supervised ML in bioinformatics
02/05-02/07	Essentials of machine learning in bioinformatics, NLP and text mining
02/12-02/14	Sequence classification with Linear and Logistic Regression
02/19-02/21	Language models using k-mers and word embeddings
02/26-02/28	Vector space representations: clustering & visualization with PCA, t-SNE, UMAP
03/04-03/06	Hidden Markov Models and Markov chains

03/11-03/13	Review for midterm with problem-solving exercises / <i>Midterm 1</i>
03/18-03/20	Sequence classification with Naive Bayes
03/25-03/27	Deep Learning introduction, fundamentals and architectures
04/01-04/05	<i>Spring recess</i>
04/08-04/10	Deep Learning in bioinformatics: CNNs, LSTMs, Recurrent Neural Networks (RNNs), Long Short Term Memory (LSTM) neural networks for sequence modelling
04/15-04/17	Word embeddings and language models with neural networks, transformers, BERT, transfer learning
04/22-04/24	Review for midterm with problem-solving exercises / <i>Midterm 2</i>
04/29-05/01	Efficient sequence searching, min-hashing, locality-sensitive hashing, vector quantization
05/06-05/08	Case studies using deep learning in bioinformatics / Project discussion
05/13	Project presentations. Review, wrap-up
	Final exam on Wednesday, May 15, 9:45 AM-12:00 PM

The schedule is subject to change with fair notice.