

CS 592: AI for Scientific Discovery

Course Information

When: Mon/Wed/Fri, 12:30 pm -- 1:20 pm.

Where: LWSN 1106.

Instructor: Yexiang Xue, yexiang@purdue.edu.

Office Hour: Mon. 11:30 pm -- 12:30 pm. *Where:* LWSN 2142V (or virtual meeting). Appointments need to be filed (via emails) at least 24 hours in prior.

Course website:

<https://www.cs.purdue.edu/homes/yexiang/courses/24spring-cs592/index.html>.

Notifications and slides will be via Brightspace
(<https://purdue.brightspace.com/>).

Discussions will be via emails, office hours, (virtual) meetings.

Course project submission at CMT (<https://cmt3.research.microsoft.com/>).

Course Description

Despite the latest progress of Artificial Intelligence (AI), fundamental knowledge gaps need to be addressed before AI can be proven useful to accelerate scientific discovery and design. AI-driven scientific discovery aims at discovering new science knowledge automatically from experiment

data, while science-based design searches for better designs guided by scientific knowledge. The main difficulty for AI to be applied in both fields lies in the disconnection between the knowledge learned by neural networks in the form of parameter weights, and human knowledge such as physics laws and constraints. In scientific discovery, such disconnection results in black-box models, where one can hardly verify if any new knowledge has been learned and how the model extrapolates to unseen environments. In science-based design, such disconnection leads to useless designs violating physical rules from pure data-driven approaches.

This class explores ways to advance scientific discovery and science-based design via novel AI technologies. The goal of this course has two folds. First, this course intends to expose students with fundamental computational tools to address challenges in scientific discovery and design, such as mathematical programming, planning, constraint satisfaction, multi-agent modeling and statistical machine learning. Second, this course intends to motivate students with successful applications of artificial intelligence on real-world problems in scientific discovery and design. We intend to cover successful applications of artificial intelligence in discovering new materials with human computation, ecological monitoring through citizen science programs, etc.

Classes will consist of instructor presentations, student presentations, and group discussions. The first few lectures consist of introductions to basic computational tools, such as constraint programming, probabilistic inference, supervised and unsupervised learning, etc. Then the course moves to discussing successful applications of AI on scientific discovery and design. Students are expected to (1) read, discuss, and present research papers, (2) complete a semester-long class project in groups, (3) review and comment on one class project proposal from other groups.

Prerequisites

Basic knowledge of linear algebra and calculus, a basic course in probability and statistics (e.g., STAT301/350/416, CS373), and basic programming skills (e.g., CS381) are required. Students without this background should have a discussion with the instructor prior to registering the course.

Target Students

Graduate or senior undergraduate students with interest in machine learning, data mining and artificial intelligence in general. This course also welcomes students from related fields, such as agriculture, economics, applied math, physics, chemistry, and engineering, who are interested in using computational tools to solve real-world problems in their domain of interest.

Course Activities and Evaluation

	% final score	Due (exam) date (tentative)
Attendance:	10%	
Course presentations:	25%	Topic, papers to read, and slides to use due 2 weeks prior to presentation.
Course project proposal:	20%	
Course project reviews:	10%	
Course project mid-term progress report:	10%	
Course project final report (and slides):	15%	

Course project final presentation:	10%	
---	------------	--

Attendance: Attendance is highly encouraged. The instructor will include random in course small quizzes to account for attendance.

Course presentations: Each student needs to present at least one topic in the class. Students are encouraged to form teams in the presentation. The available topics can be seen in the syllabus section.

Course project: MOST important part of this course. AI is a **practical field**, so it cannot be emphasized more the importance of completing a project *yourself!* In addition, because this is a graduate-level course, one important aspect is **basic scientific training**, including asking the right questions, commenting others' work, literature review, experimental design, coding, scientific writing, and presentation skills. This can **ONLY** be trained with a real project. Each student is encouraged to lead in one research topic.

We provide a few research thrusts and datasets for your reference (see below). You are encouraged to choose a specific project within the overarching theme of one research thrust in the list, although you are free to choose any project at your will as long as it relates to AI for scientific discovery and design. The goal is to nurture **GROUND-BREAKING** course projects, which have potentials to be developed into innovative research papers in the future. Course projects outside of the suggested thrusts will receive less mentoring from the instructors and the TAs, and therefore are less preferred. We encourage you to *combine your domain of expertise with AI*. To guide you through the project, we split the entire process into five parts: proposal, peer review, mid-term report, final report and presentation.

Course project proposal: the proposal will be evaluated by intellectual merit, broader impact, and tractability (same criteria for NSF proposals). The instructor **DO** respect that it is a course project, so the bar is much lower. However, the following three aspects are emphasized equally: (i)

intellectual merit: how does the project advance machine learning (or your understanding on machine learning); (ii) broad impact: how does the course project bring impact to a practical field via machine learning? (iii) tractability: is this proposal tractable (as a one-semester course project)? [grading rubrics will be posted on Brightspace.]

Course project reviews: Each student is asked to review at least three proposals of others. The student is asked to review proposals based on intellectual merit, broader impact, and tractability. Peer reviews are safety belts for other students. Unrealistic proposals should be flagged out. We will also host review panels in the class, where reviewers discuss project proposals and give panel summary reviews. Gaming does not work: the grading of the original proposal will NOT be affected by how other students review your proposal. [grading rubrics will be posted on Brightspace.]

Course project mid-term progress report: Each group is expected to submit a progress report by the deadline. This is to ensure that all projects are progressing on the right track. [grading rubrics will be posted on Brightspace.]

Course project final report / presentation: The final report and presentation will be graded in a similar way as conference papers (presentations) by TAs and the instructor jointly (although the bar is much lower). [grading rubrics will be posted on Brightspace.]

Grading Scale

The exact grading scale will be determined at the end of the semester, but use the following as a guideline for the course. The instructor promises that any alterations will be in favor of more generosity, not less.

Grade	Score	Grade	Score	Grade	Score
A+	100-96.0	A	95.9-93.0	A-	92.9-90.0
B+	89.9-86.0	B	85.9-83.0	B-	82.9-80.0
C+	79.9-76.0	C	75.9-73.0	C-	72.9-70.0
D+	69.9-66.0	D	65.9-63.0	D-	62.9-60.0

Syllabus (Tentative)

Time	Topic	Notes
1/8 Mon.	Introduction and logistics: Slides and announcements will be posted on Brightspace (https://purdue.brightspace.com/). Please let the instructor know if you cannot log into the course page on Brightspace.	
1/10 Wed.	Course overview: AI-driven scientific discovery and design.	
1/12 Fri.	Continue	
1/15 Mon.	Scientific Discovery using Scientific Approaches	
1/17 Wed.	Continue	
1/19 Fri.	Discussion	
1/22 Mon.	Symbolic and Statistical Reasoning Integration	
1/24 Wed.	Continue	
1/26 Fri.	Discussion	
1/29 Mon.	Satisfiability Modulo Counting for Symbolic and Statistical Reasoning Integration	
1/31 Wed.	Continue	
2/2 Fri.	Discussion	
	STUDENT PRESENTATIONS START	
2/5 Mon.	Review panel for proposals.	

2/7 Wed.	Continue	
2/9 Fri.	Continue	
2/12 Mon.		
2/14 Wed.		
2/16 Fri.		
2/19 Mon.		
2/21 Wed.		
2/23 Fri.		
2/26 Mon.		
2/28 Wed.		
3/1 Fri.		
3/4 Mon.		
3/6 Wed.		
3/8 Fri.		
3/11 Mon.	Spring break	
3/13 Wed.	Spring break	
3/15 Fri.	Spring break	

3/18 Mon.		
3/20 Wed.		
3/22 Fri.		
3/25 Mon.		
3/27 Wed.		
3/29 Fri.		
4/1 Mon.		
4/3 Wed.		
4/5 Fri.		
4/8 Mon.	Final project presentation I.	
4/10 Wed.	Final project presentation II.	
4/12 Fri.	Final project presentation III.	
4/15 Mon.	Final project presentation IV.	
4/17 Wed.	Final project presentation V.	
4/19 Fri.	Final project presentation VI.	
4/22 Mon.	Week for review.	
4/24 Wed.	Week for review.	

4/26
Fri.

Week for review.

Course Project Thrusts

Thrust 1: stochastic optimization: encoding machine learning for decision-making

In data-driven decision-making, we have to reason about the optimal policy of a system given a stochastic model learned from data. For example, one can use a machine learning model to capture the traffic dynamics of a road network. The decision-making problem is: given the traffic dynamics learned from data, what is the most efficient way to travel between a pair of locations? Notice that the solution can change dynamically, depending on the shift in traffic dynamics. As another example in Physics, machine learning models have been used to predict the band-gap of many metal alloy materials. The decision-making problem is: given the machine learning model, what is the best alloy, which is both cheap to synthesize and has a good band-gap property?

The aforementioned examples are stochastic optimization problems, which make robust interventions that maximize the "expectation" of stochastic functions learned from data. It arises naturally in many applications ranging from economics, operational research, and artificial intelligence. Stochastic optimization combines two intractable problems, one of which is the inner probabilistic inference problem to compute the expectation across exponentially many probabilistic outcomes, and the other of which is the outer optimization problem to search for the optimal policy.

Research questions: (i) if the inner machine learning model is a decision tree, can you compute the optimal policy in polynomial time? How? (ii) What if the inner machine learning model is a logistic regression, a linear SVM, a kernelized SVM, a random forest, or a probabilistic graphical

model? (iii) What if the machine learning model is temporal, such as a recurrent neural network or a LSTM? (iv) In case the inner probabilistic inference problem is intractable, existing approaches to solve stochastic optimization problems approximate the intractable probabilistic inference sub-problems either in variational forms, or via the empirical mean of pre-computed, fixed samples. There is also a recent approach which approximates the intractable sub-problems with optimization queries, subject to randomized constraints (see following papers). Question: how does various approximation schemes of the inner machine learning models affect the overall solution quality of the stochastic optimization problem? (v) Suppose we are solving one stochastic optimization problem for a specific application, can we adapt existing approximation schemes in any way to fit the problem instance for better results?

Papers:

Yexiang Xue, Zhiyuan Li, Stefano Ermon, Carla P. Gomes, Bart Selman.
Solving Marginal MAP Problems with NP Oracles and Parity Constraints

In the Proceedings of the 29th Annual Conference on Neural Information Processing Systems (NIPS), 2016. [[pdf](#)] [[spotlight video](#)]

Anton J. Kleywegt, Alexander Shapiro, and Tito Homem-de Mello.
The sample average approximation method for stochastic discrete optimization.

SIAM Journal on Optimization, 2002. [[pdf](#)]

Miguel Á. Carreira-Perpiñán and Geoffrey E. Hinton.
On contrastive divergence learning.

AISTATS, 2005. [[pdf](#)]

Martin Dyer and Leen Stougie.

Computational complexity of stochastic programming problems.

Mathematical Programming, 2006. [[springer](#)]

John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira.

Conditional random fields: Probabilistic models for segmenting and

labeling sequence data. *In Proceedings of the Eighteenth International Conference on Machine Learning, ICML, 2001.* [[pdf](#)]

Stefano Ermon, Carla Gomes, Ashish Sabharwal, and Bart Selman.
Taming the Curse of Dimensionality: Discrete Integration by Hashing and Optimization

In Proc. 30th International Conference on Machine Learning (ICML) 2013. [[pdf](#)]

Carla P. Gomes, Ashish Sabharwal, Bart Selman.
Near-Uniform Sampling of Combinatorial Spaces Using XOR Constraints.

NIPS 2006. [[pdf](#)]

Carla P. Gomes, Willem Jan van Hoeve, Ashish Sabharwal, Bart Selman.
Counting CSP Solutions Using Generalized XOR Constraints.

AAAI 2007. [[pdf](#)]

Yexiang Xue*, Xiaojian Wu*, Bart Selman, and Carla P. Gomes.
XOR-Sampling for Network Design with Correlated Stochastic Events.

In Proc. 26th International Joint Conference on Artificial Intelligence (IJCAI), 2017. [[pdf](#)]

* indicates equal contribution.

Yexiang Xue, Xiaojian Wu, Dana Morin, Bistra Dilkina, Angela Fuller, J. Andrew Royle, and Carla Gomes.

Dynamic Optimization of Landscape Connectivity Embedding Spatial-Capture-Recapture Information.

In Proc. 31th AAAI Conference on Artificial Intelligence (AAAI), 2017.

[[pdf](#)] [[supplementary materials](#)]

Thrust 2: embedding physical constraints into deep neural networks

The emergence of large-scale data-driven machine learning and optimization methodology has led to successful applications in areas as

diverse as finance, marketing, retail, and health care. Yet, many application domains remain out of reach for these technologies, when applied in isolation. In the area of medical robotics, for example, it is crucial to develop systems that can recognize, guide, support, or correct surgical procedures. This is particularly important for next-generation trauma care systems that allow life-saving surgery to be performed remotely in presence of unreliable bandwidth communications. For such systems, machine learning models have been developed that can recognize certain commands and procedures, but they are unable to learn complex physical or operational constraints. Constraint-based optimization methods, on the other hand, would be able to generate feasible surgical plans, but currently, have no mechanism to represent and evaluate such complex environments. To leverage the required capabilities of both technologies, we have to find an integrated method that embeds constraint reasoning in machine learning.

In a seminal paper, the authors proposed an approach, which provides a scalable method for machine learning over structured domains. The core idea is to augment machine learning algorithms with a constraint reasoning module that represents physical or operational requirements. Specifically, the authors propose to embed decision diagrams, a popular constraint reasoning tool, as a fully-differentiable layer in deep neural networks. By enforcing the constraints, the output of generative models can now provide assurances of safety, correctness, and/or fairness. Moreover, this approach enjoys a smaller modeling space than traditional machine learning approaches, allowing machine learning algorithms to learn faster and generalize better.

Research questions: (i) are there any other ways to enforce physical constraints other than using a decision diagram in the seminal work? (ii) What if the constraints are too complicated which cannot be fully captured by a decision diagram? (iii) In a specific applicational domain, is there a better way to encode constraints? (iv) Does enforcing physical constraints make machine learning easier or more difficult? Can you quantify the difference? (v) Can we apply this idea in natural language processing, computer vision, reinforcement learning, etc? (vi) Ethics and fairness in

machine learning are being discussed in our community. Can we use this technique to guarantee the ethics and/or the fairness of a machine learning model?

Papers:

Yexiang Xue, Willem-Jan van Hoeve.

Embedding Decision Diagrams into Generative Adversarial Networks.
In *Proc. of the Sixteenth International Conference on the Integration of Constraint Programming, Artificial Intelligence, and Operations Research (CPAIOR)*, 2019. [[springer](#)]

Md Masudur Rahman, Natalia Sanchez-Tamayo, Glebys Gonzalez, Mridul Agarwal, Vaneet Aggarwal, Richard M. Voyles, Yexiang Xue, and Juan Wachs.

Transferring Dexterous Surgical Skill Knowledge between Robots for Semi-autonomous Teleoperation.
In *ROMAN*, 2019. [[pdf](#)]

Naveen Madapana, Md Masudur Rahman, Natalia Sanchez-Tamayo, Mythra V. Balakuntala, Glebys Gonzalez, Jyothsna Padmakumar Bindu, L. N. Vishnunandan Venkatesh, Xingguang Zhang, Juan Barragan Noguera, Thomas Low, Richard M. Voyles, Yexiang Xue, and Juan Wachs
DESK: A Robotic Activity Dataset for Dexterous Surgical Skills Transfer to Medical Robots.
In *IROS*, 2019. [[pdf](#)]

Matt J. Kusner, Brooks Paige, José Miguel Hernández-Lobato.

Grammar Variational Autoencoder.

In *Proceedings of the 34th International Conference on Machine Learning, ICML, 2017*. [[pdf](#)]

Chenglong Wang, Kedar Tatwawadi, Marc Brockschmidt, Po-Sen Huang, Yi Mao, Oleksandr Polozov, Rishabh Singh
Text-to-SQL Generation with Execution-Guided Decoding
[[pdf](#)]

Kevin Lin, Ben Bogin, Mark Neumann, Jonathan Berant, Matt Gardner
Grammar-based Neural Text-to-SQL Generation
[\[ArXiv\]](#)

Thrust 3: machine learning for scientific discovery and/or social good

Machine learning models have defeated the brightest mind in this world (see the story of AlphaGo). Now, instead of using this technology for game playing, can we harness the tremendous progress in AI and machine learning to make our world a better place? In particular, I am curious at problems that have attracted the smartest minds of man kind historically -- the discovery of new science. Besides scientific discovery, can we use machine learning to create positive social impact?

If you think about it: in AlphaGo, machine learning is used to find a strategy in a highly complex space (all possible moves of Go), which beats all opponent's strategies. The problem is similar for scientific discovery, except that we are now playing Go with nature. For example, in materials discovery, we would like to find the best material in a highly complex space (all possible compositions) which enjoys the best properties. Should the strategy which was proven successful for Go work for scientific discovery (and/or AI for social good)?

I am listing a few example papers below in which machine learning are used successfully for scientific discovery and for social good. I hope this can motivate you to discover a good applicational area of machine learning. The key to the success is to combine your domain of expertise with machine learning.

Papers:

Yexiang Xue, Junwen Bai, Ronan Le Bras, Brendan Rappazzo, Richard Bernstein, Johan Bjorck, Liane Longpre, Santosh K. Suram, Robert B. van Dover, John Gregoire, and Carla Gomes.

Phase-Mapper: An AI Platform to Accelerate High Throughput Materials Discovery.

In *Proc. 29th Annual Conference on Innovative Applications of Artificial Intelligence (IAAI)*, 2017. [[pdf](#)][[video 1](#)][[video 2](#)][[video 3](#)]

Santosh K. Suram, Yexiang Xue, Junwen Bai, Ronan LeBras, Brendan H Rappazzo, Richard Bernstein, Johan Bjorck, Lan Zhou, R. Bruce van Dover, Carla P. Gomes, and John M. Gregoire.

Automated Phase Mapping with AgileFD and its Application to Light Absorber Discovery in the V-Mn-Nb Oxide System.

In *American Chemical Society Combinatorial Science*, Dec, 2016. [[DOI](#)]
[[pdf](#)][[video 1](#)][[video 2](#)][[video 3](#)]

Junwen Bai, Yexiang Xue, Johan Bjorck, Ronan Le Bras, Brendan Rappazzo, Richard Bernstein, Santosh K. Suram, Robert Bruce van Dover, John M. Gregoire, Carla P. Gomes.

Phase Mapper: Accelerating Materials Discovery with AI.

In *AI Magazine*, Vol. 39, No 1. 2018. [[paper](#)]

Yexiang Xue, Ian Davies, Daniel Fink, Christopher Wood, Carla P. Gomes.
Avicaching: A Two Stage Game for Bias Reduction in Citizen Science
In *the Proceedings of the 15th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, 2016. [[pdf](#)][[supplementary materials](#)][[video](#)]

Giuseppe Carleo and Matthias Troyer

Solving the quantum many-body problem with artificial neural networks.

In *Science*, 355, 2017. [[website](#)]

Ganesh Hegde and R. Chris Bowen

Machine-learned approximation to Density Functional Theory Hamiltons.

In *Scientific Reports*, 7, 2016. [[ArXiv](#)]

Graham Roberts, Simon Y. Haile, Rajat Sainju, Danny J. Edwards, Brian Hutchinson and Yuanyuan Zhu

Deep Learning for Semantic Segmentation of Defects in Advanced

STEM Images of Steels.

Scientific Reports, volume 9, 2019. [[website](#)]

Academic Policies

Late policy

Assignments are to be submitted by the due date listed. Each person will be allowed two days of extensions which can be applied to any combination of assignments (homework/projects only; exams excluded) during the semester without penalty. After that, a late penalty of 15% per day will be assigned. The use of a partial day will be counted as a full day. Use of extension days must be stated explicitly at the time of the late submission (by accompanying email to ALL TAs and the instructor), otherwise, late penalties will apply. Extensions cannot be used after the final day of classes. Extension days cannot be rearranged after they are applied to a submission. Additional no-penalty late days may be introduced in the later part of the semester conditioned on the completion of the course evaluations (details to be finalized). Assignments, project reports, etc, will NOT BE accepted if they are more than five days late (and receive zero points). Additional extensions will be granted only due to serious and documented medical or family emergencies. Use the late days wisely!

Classroom engagement is extremely important and associated with your overall success in the course. The importance and value of course engagement and ways in which you can engage with the course content even if you are in quarantine or isolation, will be discussed at the beginning of the semester. Student survey data from Fall 2020 emphasized students' views of in-person course opportunities as critical to their learning, engagement with faculty/TAs, and ability to interact with peers.

Only the instructor can excuse a student from a course requirement or responsibility. When conflicts can be anticipated, such as for many University-sponsored activities and religious observations, the student should inform the instructor of the situation as far in advance as possible.

For unanticipated or emergency conflicts, when advance notification to an instructor is not possible, the student should contact the instructor/instructional team as soon as possible by email, through Brightspace, or by phone. In cases of bereavement, quarantine, or isolation, the student or the student's representative should contact the Office of the Dean of Students via email or phone at 765-494-1747. Our course Brightspace includes a link to the Dean of Students under 'Campus Resources.'

Academic honesty

Please read the departmental [academic integrity policy](#). This will be followed unless we provide written documentation of exceptions.

- Unless stated otherwise, each student should write up their own solutions independently. You need to indicate the names of the people you discussed a problem with; ideally you should discuss with no more than two other people.
- **NO PART OF THE STUDENT'S ASSIGNMENT (PROJECT, NOTES, ETC) SHOULD BE COPIED FROM ANOTHER STUDENT OR FROM OTHER RESEARCHERS OR FROM THE WEB (Plagiarism).** We encourage you to interact amongst yourselves: you may discuss and obtain help with basic concepts covered in lectures or the textbook, homework specification (but not solution), and general ideas of program implementation (but not the code). However, unless otherwise noted, work turned in should reflect your own efforts and knowledge. Sharing or copying solutions is unacceptable and could result in failure of this course. We use copy detection software, so do not copy code and make changes (either from the Web or from other students). You are expected to take reasonable precautions to prevent others from using your work.
- Any student not following these guidelines are subject to an automatic F (final grade).

Nondiscrimination Statement

Purdue University is committed to maintaining a community which recognizes and values the inherent worth and dignity of every person; fosters tolerance, sensitivity, understanding, and mutual respect among its members; and encourages each individual to strive to reach his or her potential. In pursuit of its goal of academic excellence, the University seeks to develop and nurture diversity. The University believes that diversity among its many members strengthens the institution, stimulates creativity, promotes the exchange of ideas, and enriches campus life. A hyperlink to Purdue's full Nondiscrimination Policy Statement is included [here](#).

Accessibility

Purdue University strives to make learning experiences as accessible as possible. If you anticipate or experience physical or academic barriers based on disability, you are welcome to let me know so that we can discuss options. You are also encouraged to contact the Disability Resource Center at: drc@purdue.edu or by phone: 765-494-1247.

Mental Health/Wellness Statement

If you find yourself beginning to feel some stress, anxiety and/or feeling slightly overwhelmed, try WellTrack. Sign in and find information and tools at your fingertips, available to you at any time.

If you need support and information about options and resources, please contact or see the Office of the Dean of Students. Call 765-494-1747. Hours of operation are M-F, 8 am- 5 pm.

If you find yourself struggling to find a healthy balance between academics, social life, stress, etc. sign up for free one-on-one virtual or in-person sessions with a Purdue Wellness Coach at RecWell. Student coaches can help you navigate through barriers and challenges toward your goals throughout the semester. Sign up is completely free and can be done

on BoilerConnect. If you have any questions, please contact Purdue Wellness at evans240@purdue.edu.

If you're struggling and need mental health services: Purdue University is committed to advancing the mental health and well-being of its students. If you or someone you know is feeling overwhelmed, depressed, and/or in need of mental health support, services are available. For help, such individuals should contact Counseling and Psychological Services (CAPS) at 765-494-6995 during and after hours, on weekends and holidays, or by going to the CAPS office on the second floor of the Purdue University Student Health Center (PUSH) during business hours.

Emergency Preparation

In the event of a major campus emergency, course requirements, deadlines and grading percentages are subject to changes that may be necessitated by a revised semester calendar or other circumstances beyond the instructor's control. Relevant changes to this course will be posted onto the course website or can be obtained by contacting the instructors or TAs via email or phone. You are expected to read your @purdue.edu email on a frequent basis.

Other general course policies can be found [here](#).

Resource

Datasets

[eBird citizen science dataset.](#)

[Synthetic and real datasets for materials discovery.](#)

[Dataset for the corridor-design problem and landscape optimization problem.](#)

[Remote sensing images](#) (a code repository which contains code to download from Google Earth engine).

[UCI Machine Learning Dataset.](#)

[Kaggle.](#)

Online Resources

Machine learning:

- [A First Encounter with Machine Learning](#) by Max Welling
- [Introduction to Machine Learning](#) by Alex Smola and S.V.N. Vishwanathan
- [A Course in Machine Learning](#) by Hal Daume III
- [Bayesian reasoning and machine learning](#) by David Barber
- [A tutorial](#) by Andrew Moore

Math references:

- [The Matrix Cookbook](#) by Kaare Brandt Petersen and Michael Syskind Pedersen
- [Calculus](#) by Gilbert Strang
- [Linear Algebra](#) by Gilbert Strang
- [Introduction to Probability and Statistics](#) by Jeremy Orloff and Jonathan Bloom

Learning Python

For those who are unfamiliar with Python, I strongly encourage you to spend one night learning it by following the official tutorial (see below). I did not know Python until my graduate school. It took me one night to learn it, so can you!

- [The official Python Tutorial](#)
- [LeetCode](#)

