

Lecture 1

Introduction

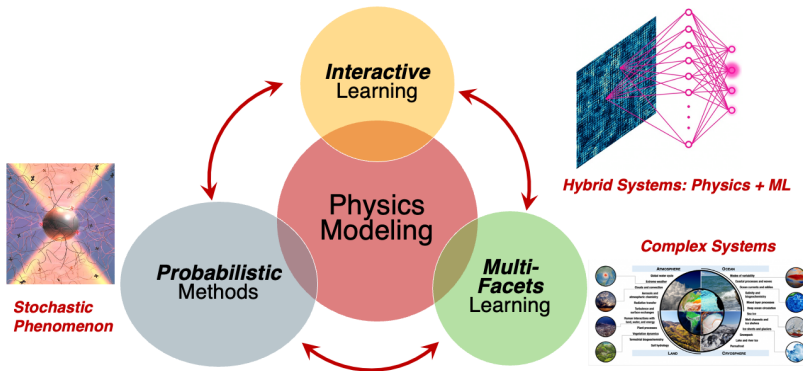
Instructor: Shibo Li

shiboli@cs.fsu.edu



Department of Computer Science
Florida State University

- **Shibo Li**, shiboli@cs.fsu.edu
- Probabilistic Machine Learning
- Assistant Professor, Department of Computer Science



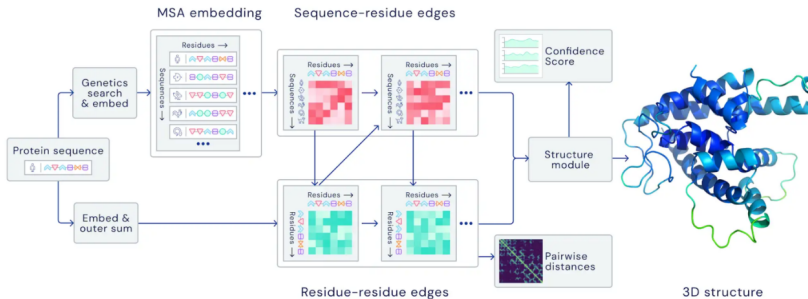
*"A computer program is said to learn from experience **E** with respect to some class of tasks **T** and performance measure **P**, if its performance at tasks in **T**, as measured by **P**, improves with experience **E**."*

– Tom M. Mitchell, *Machine Learning* (1997)

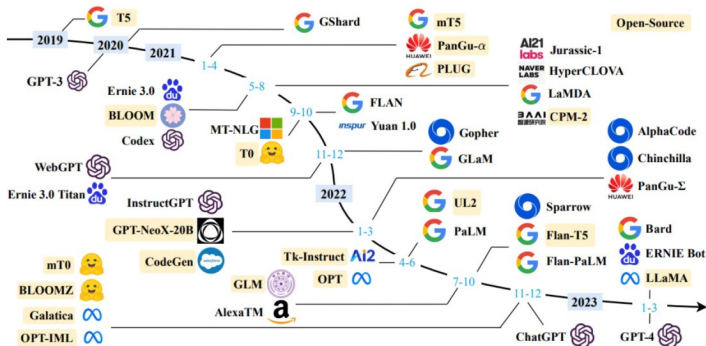


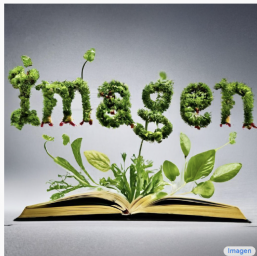
Machine Learning is the driving force of AI!





WTA





Sprouts in the shape of text 'Imagen' coming out of a fairytale book.



A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.



A high contrast portrait of a very happy fuzzy panda dressed as a chef in a high end kitchen making dough. There is a painting of flowers on the wall behind him.



Teddy bears swimming at the Olympics 400m Butterfly event.



A cute corgi lives in a house made out of sushi.



A cute sloth holding a small treasure chest. A bright golden glow is coming from the chest.

Machine Learning is EVERYWHERE!

What Other Items Do Customers Buy After Viewing This Item?



Wasabi Power Battery (2-Pack) and Dual Charger for GoPro HERO4 and GoPro AHD8T-401, AHB8P-401

★★★★★ (238)

\$23.99



SanDisk Extreme 64GB UHS-I/U3 Micro SDXC Memory Card Up To 60MB/s Read With Adapter- ...

★★★★★ (443)

\$79.99



EEEKit 8-in-1 Accessories Kit for Gopro Hero4 Black/Silver Hero HD 3+/3/2/1 Camera, Head Belt Strap...

★★★★★ (299)

\$29.99



SanDisk Ultra 32GB UHS-I/Class 10 Micro SDHC Memory Card Up to 48MB/s With Adapter- ...

★★★★★ (2,719)

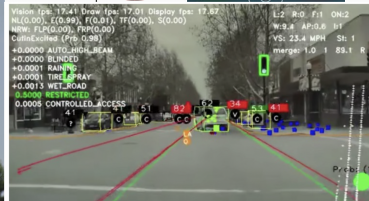
\$19.44

Explore similar items

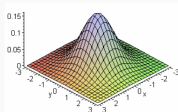
Translate

English Spanish French

Jan de kinde



In a nutshell, probabilistic learning is branch of ML that uses probabilistic (or Bayesian) principles for model design and algorithm development.!



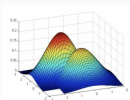
Prior Distribution

$$p(\theta)$$



Data Likelihood

$$p(\mathcal{D}|\theta)$$



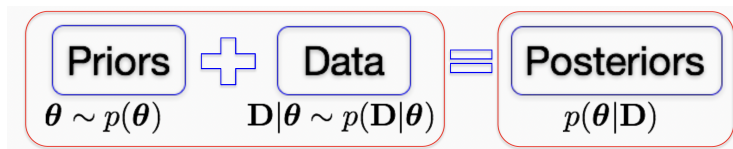
Posterior Distribution

$$p(\theta|\mathcal{D})$$

Baye's Rule:

$$p(\theta|\mathcal{D}) = \frac{p(\theta, \mathcal{D})}{p(\mathcal{D})} = \frac{p(\theta)p(\mathcal{D}|\theta)}{\int p(\theta)p(\mathcal{D}|\theta)d\theta}$$

- Unified, principled mathematical framework



- Uncertainty reasoning

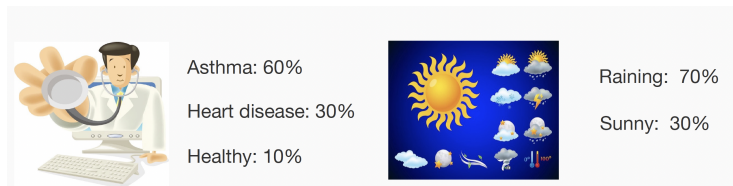


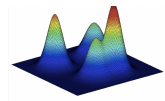


Figure: “Neither Autopilot nor the driver noticed the white side of the tractor trailer against a brightly lit sky, so the brake was not applied. The high ride height of the trailer combined with its positioning across the road and the extremely rare circumstances of the impact caused the Model S to pass under the trailer, with the bottom of the trailer impacting the windshield,” said Tesla in a statement after the crash last year.

► Modeling



Complex
Knowledge/Assumptions



Valid Prior
Distribution

► Calculation

$$p(\boldsymbol{\theta}|\mathbf{D}) = \frac{p(\boldsymbol{\theta})p(\mathbf{D}|\boldsymbol{\theta})}{\int p(\boldsymbol{\theta})p(\mathbf{D}|\boldsymbol{\theta})d\boldsymbol{\theta}}$$

High dimensional integration

MCMC sampling

Variational approximations

Belief propagation

We will cover both the classical and state-of-the-art approaches to deal with these challenges!

Warning

- ▶ This course is **math intensive** and requires **a certain level of programming** (with Matlab, R or Python). Python components may require TensorFlow and/or PyTorch. The coding workload is not heavy, but requests **mathematical derivations and careful debugging**.
- ▶ The workload is heavy (6-10 hours per-week)

How will you learn?

- ▶ Attendance is **Required**. Take classes to follow the math, understand the models and algorithms
- ▶ **Derive the math details by yourself!**
- ▶ Finish the homework assignments to deepen your understanding
- ▶ **Implement and debug the models and algorithms by yourself!**
- ▶ Using **proper help is encouraged!** (ChatGPT, Claude, StackOverflow, MathStacks...). **But plagiarism is prohibited!**

This course focuses on the **mathematic foundations, modeling and algorithmic ideas in probabilistic learning**

This course is **not** about

- ▶ Applying ML to specific tasks (e.g., image tagging and autonomous driving)
- ▶ Using specific ML tools/libraries, e.g., scikit-learn and PyTorch
- ▶ How to program and debug, e.g., with Python, R or Matlab

This course is an **advanced** course for students who want to study ML **in depth** or quickly get to the **frontier** research of probabilistic learning

This course is **not** about a preliminary course, e.g., entry-level introduction of statistics. That means,

The content can be hard for some ones

- ▶ You are struggling with linear algebra, calculus or basic statistical concepts
- ▶ You are sick of mathematical symbols, derivations, proofs and calculations
- ▶ You do NOT feel good in programming and debugging

- ▶ You are **not** scared of math, statistics, calculus and calculations; you are happy with them!
- ▶ You are **comfortable** with abstract symbols and matrices operations
- ▶ You can pick-up Matlab/Python/R quickly (even if you have never used them before)
- ▶ You enjoy debugging, step in, step out, print, etc.
- ▶ You can quickly learn how to use TensorFlow or PyTorch or Jax by following the documentation and searching for the online examples

- ▶ You have planed for **enough efforts** for this class (e.g., 6-10 hours per-week)

You feel **NOT** right about any of these assumptions

- ▶ Seriously consider whether to take this course
- ▶ We want you to succeed. We do not want to make you feel tortured.

- ▶ The course website contains all the detailed information
- ▶ The course website is linked to my homepage
 - ▶ My homepage: `https://imshibo.com/`
 - ▶ Course Website: `https://cis5930.github.io/`

Note: this review is neither compressive nor in depth. Due to time limit, this review is just to point out **key concepts and computational rules as the guidance**. We list the references for you to check out details for future usage.

- Standard notations
 - non-bold letters: scalars

$$a, b, x, y, B, D, G, \alpha, \gamma, \dots$$

- Bold small letters: vectors

$$\mathbf{a}, \mathbf{b}, \mathbf{x}, \mathbf{y}, \boldsymbol{\gamma}, \boldsymbol{\eta}, \dots$$

- Bold capital: matrices

$$\mathbf{A}, \mathbf{X}, \mathbf{Z}, \mathbf{\Gamma}, \dots$$

- scalar input, scalar output

$$y(x + dx) = y(x) + a \cdot dx + (\text{high-order terms})$$

$$\frac{\partial y}{\partial x} = a \longleftrightarrow dy = a \cdot dx$$

- vector input, scalar output

$$\mathbf{x} = (x_1, \dots, x_n)^\top, \quad d\mathbf{x} = (dx_1, \dots, dx_n)^\top$$

$$y(\mathbf{x} + d\mathbf{x}) = y(\mathbf{x}) + \mathbf{a}d\mathbf{x} + (\text{high-order terms})$$

We use row-vector to represent gradient

$$\frac{\partial y}{\partial \mathbf{x}} = \left(\frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_n} \right) = \mathbf{a} \longleftrightarrow dy = \mathbf{a}d\mathbf{x}$$

- In general, vector input, vector output

$$\mathbf{y} = (y_1, \dots, y_m)^\top \quad \mathbf{x} = (x_1, \dots, x_n)^\top, \quad d\mathbf{x} = (dx_1, \dots, dx_n)^\top$$

$$\mathbf{y}(\mathbf{x} + d\mathbf{x}) = \mathbf{y}(\mathbf{x}) + \mathbf{A}d\mathbf{x} + (\text{high order terms})$$

What is this? What's size? $m \times n$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} = \mathbf{A} \quad \longleftrightarrow \quad d\mathbf{y} = \mathbf{A}d\mathbf{x}$$

$$y(x + dx) = y(x) + a \cdot dx + (\text{high-order terms})$$

$$y(\mathbf{x} + d\mathbf{x}) = y(\mathbf{x}) + a d\mathbf{x} + (\text{high-order terms})$$

$$\mathbf{y}(\mathbf{x} + d\mathbf{x}) = \mathbf{y}(\mathbf{x}) + \mathbf{A} d\mathbf{x} + (\text{high order terms})$$

In all the cases, $\{a, \mathbf{a}, \mathbf{A}\}$ are derivatives. We define (partial) gradient as the derivative

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{A}$$

This is consistent with the definition of Jacobian. However, we need to be aware if output is scalar, the gradient is a row vector

- What is the benefit of this notation? We can apply the chain-rule in a natural way

$$\mathbf{y} = \mathbf{f}(\mathbf{x}), \quad \mathbf{x} = \mathbf{g}(\mathbf{z})$$

$$\mathbf{y} : m \times 1 \quad \mathbf{x} : n \times 1 \quad \mathbf{z} : q \times 1$$

$$\begin{array}{c} \boxed{\frac{\partial \mathbf{y}}{\partial \mathbf{z}}} = \boxed{\frac{\partial \mathbf{y}}{\partial \mathbf{x}}} \cdot \boxed{\frac{\partial \mathbf{x}}{\partial \mathbf{z}}} \\ \swarrow \quad \uparrow \quad \nwarrow \\ m \times q \quad m \times n \quad n \times q \end{array}$$

- Some literature uses the notation of derivative transpose

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{A}^\top$$

The benefit is for scalar y , the gradient is a column vector. The cons is when doing the chain rule, you have to multiply from right to left. [Why?](#)

- In whichever case, the key to derive/compute the derivative!

$$\mathbf{y}(\mathbf{x} + d\mathbf{x}) = \mathbf{y}(\mathbf{x}) + \mathbf{A}d\mathbf{x} + (\text{high order terms})$$



$$d\mathbf{y} = \mathbf{A}d\mathbf{x}$$

- The general idea: recursively apply the chain rule to get the target derivative!

- Take a scalar case as an example

$$y = 3x + \frac{1}{x^2}$$

$$dy = d(3x) + d\left(\frac{1}{x^2}\right) = \dots$$

Let's do it together

- How to apply chain rule for matrices/vectors
- like scalar case, we have a set of basic rules in matrix world as well; just keep applying them recursively

$$\begin{aligned}\partial \mathbf{A} &= 0 && (\mathbf{A} \text{ is a constant}) \\ \partial(\alpha \mathbf{X}) &= \alpha \partial \mathbf{X} \\ \partial(\mathbf{X} + \mathbf{Y}) &= \partial \mathbf{X} + \partial \mathbf{Y} \\ \partial(\text{Tr}(\mathbf{X})) &= \text{Tr}(\partial \mathbf{X}) \\ \partial(\mathbf{X}\mathbf{Y}) &= (\partial \mathbf{X})\mathbf{Y} + \mathbf{X}(\partial \mathbf{Y}) \\ \partial(\mathbf{X} \circ \mathbf{Y}) &= (\partial \mathbf{X}) \circ \mathbf{Y} + \mathbf{X} \circ (\partial \mathbf{Y}) \\ \partial(\mathbf{X} \otimes \mathbf{Y}) &= (\partial \mathbf{X}) \otimes \mathbf{Y} + \mathbf{X} \otimes (\partial \mathbf{Y}) \\ \partial(\mathbf{X}^{-1}) &= -\mathbf{X}^{-1}(\partial \mathbf{X})\mathbf{X}^{-1} \\ \partial(\det(\mathbf{X})) &= \det(\mathbf{X})\text{Tr}(\mathbf{X}^{-1}\partial \mathbf{X}) \\ \partial(\ln(\det(\mathbf{X}))) &= \text{Tr}(\mathbf{X}^{-1}\partial \mathbf{X}) \\ \partial \mathbf{X}^T &= (\partial \mathbf{X})^T\end{aligned}$$

- Let's do several examples

$$y = (\mathbf{x} + \mathbf{b})^\top (\mathbf{x} + \mathbf{b})$$

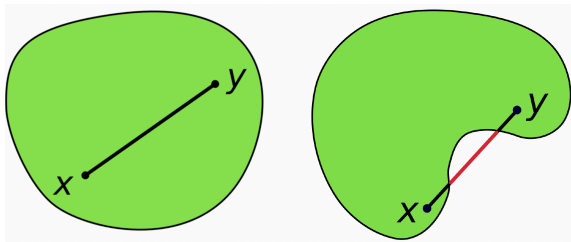
$$y = \text{tr} \left((\mathbf{I} + \mathbf{x}\mathbf{x}^\top)^{-1} \right)$$

- Commonly used references
 1. [Old and New Matrix Algebra Useful for Statistics](#), By Tom Minka, 2001
 2. [Matrix Cookbook](#)
- Strongly suggest the tutorial made by our TM for more examples

<https://www.youtube.com/watch?v=artvpNFSFgw>

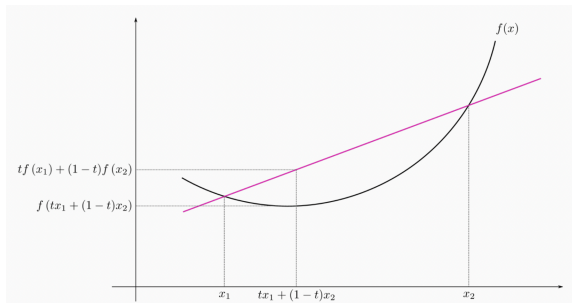
- Convex region/set in **Euclidean space** or more general vector space

$$\forall \mathbf{x}, \mathbf{y} \in \mathcal{S} \text{ and } \forall t \in [0, 1], (1 - t)\mathbf{x} + t\mathbf{y} \in \mathcal{S}$$



- The input domain $\mathbf{X} \subseteq \mathbb{R}^d$ is a convex region/set
- Convex function: $f : \mathbf{X} \rightarrow \mathbb{R}$

$$\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbf{X}, \forall t \in [0, 1] : f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$



► Examples of convex functions

Single variable

$$f(x) = e^x$$

$$f(x) = -\log(x)$$

Multivariable

$$f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x} + \mathbf{b}$$

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{x}$$

► How to determine a convex function?

- When differentiable, $\forall \mathbf{x}, \mathbf{y} \in \mathbf{X}$

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y})$$

- When twice differentiable, if Hessian matrix is positive semidefinite (PSD)

$$H_f(\mathbf{x}) = \nabla \nabla f(\mathbf{x}) \succeq 0, \quad \forall \mathbf{x} \in \mathbb{R}^d$$

PSD: All eigenvalues of the Hessian are non-negative

- **Jensen's inequality** (for convex functions), if X is a random variable: $X \sim p(X)$

$$\begin{aligned}f(\mathbb{E}[X]) &\leq \mathbb{E}[f(X)] \\f(\mathbb{E}[g(X)]) &\leq \mathbb{E}[f(g(X))]\end{aligned}$$

- Relates the value of a convex function of an integral to the integral of the convex function

- Convex conjugate (Fenchel's duality)

for an arbitrary convex function $f(\cdot)$, there exists a duality function $g(\cdot)$

$$f(x) = \max_{\lambda} \lambda x - g(\lambda)$$

$$g(\lambda) = \max_{\mathbf{x}} \lambda x - f(x)$$

Jensen's equality and convex conjugate plays the key role in approximate inference