

## Note

- Homework #2 submission closed.
- Start to work on choosing final project topic.
  - Read *LeCun et al.* if not done so yet.
  - Pick a subtopic in the paper that you like to go deeper.
  - Select main papers among those cited in the section describing your topic in *LeCun et al.*
  - Continue to select papers cited among the cited papers.
  - Finalize **at least 5 papers** on your subtopic.
  - **Submit your subtopic/summary/selected papers to Canvas forum** for my review (~~make your own thread~~).
  - **Due on March 18**

CSC872: PAMI – Kazunori Okada (C) 2025

1

1

PF

# Framework for Bayesian Probabilistic Reasoning

CSC 872

Pattern Analysis and Machine Intelligence

CSC872: PAMI – Kazunori Okada (C) 2025

2

2

# Logical Reasoning: Review

- *Knowledge based agent* Logical languages (PL, FOL) gave us formal ways to describe a world based on **Boolean truth**
  - World with *True/False*
  - PL's ontological commitment: **Facts**
  - FOL's ontological commitment: **Facts, Objects, Relations**
- **Inference**: *derive a new fact from known facts (KB) !!!*
  - Sound & Complete Inference = Entailment
- **BUT we live in a world that is very uncertain!**
  - The car accident example from the last lecture / tomorrow's weather / stock market / your and my life even! ...
- **How to represent facts with such uncertainty?**
- **How to do inference with such uncertainty?**

CSC872: PAMI – Kazunori Okada (C) 2025

3

3

# Probabilistic Reasoning

- One answer is ...
- **Bayesian Probability!**
  - OC: Proposition (Facts) *Probability*
  - OC: **Degree of Belief** ↙
  - Inference by **Bayes Rule**
- **Frequentist Approach**
  - Define a prob. as a limit of event's frequency in a large number of trials
  - More objective it seems: used for hypothesis testing ... **BUT**
  - Probability that "*it will rain tomorrow*"?
  - It rains **only once** tomorrow...



Sir Thomas Bayes (1763)  
An essay toward solving a problem in the doctrine of chances. Philosophical Transactions of the Royal Society of London, 53: 370-418  
[http://en.wikipedia.org/wiki/Bayesian\\_probability](http://en.wikipedia.org/wiki/Bayesian_probability)

CSC872: PAMI – Kazunori Okada (C) 2025

4

4

# Foundation of Probability

$X = \{\text{true}, \text{false}\}$   
 $P(X)$

- **Random Variable A** indicates an event that has intrinsic degree of **uncertainty** if A occurs or not
- A can be seen as a **function** that chooses a value from the **event space** according to probability distribution  $P(A)$  where an event is a subset of **sample space**
- **Discrete Boolean Random Variable**
  - Sample space {True, False}
- **Discrete Multivalued Random Variable (beyond PL)**
  - Sample space {Mon, ..., Sun}
- **Continuous Random Variable (beyond PL)**
  - Sample space  $\{x \in \mathbf{R}\}$

CSC872: PAMI – Kazunori Okada (C) 2025

5

5

# Foundation of Probability

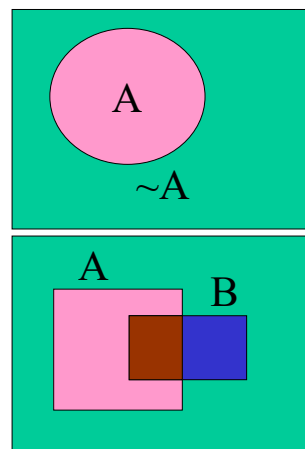
## • The axioms of probability

- $P(A) \geq 0$  and  $P(A) \in \mathbf{R}$  (non-negative real)
- $P(\Omega) = 1$  (unit measure)
- $P(A_1 \vee A_2 \vee \dots) = \sum_i P(A_i)$  (additivity)

- $P(\text{True}) = 1$
- $P(\text{False}) = 0$
- $0 \leq P(A) \leq 1$
- $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$

## • Total probability theorem

- $P(A) + P(\sim A) = 1$  given {true, false}
- $\sum_{A=v_i} P(A) = 1$  given  $\{v_1, \dots, v_K\}$



CSC872: PAMI – Kazunori Okada (C) 2025

6

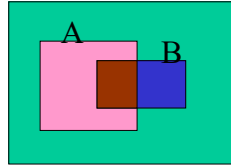
6

# Joint Distribution

*Probability*

- Probability of multiple events (A and B) in conjunction

-  $P(A \wedge B) := P(A, B)$

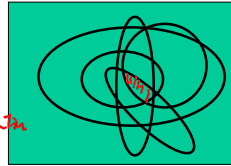


X = {CS, Math, History}  
Y = {Yes, No}

X = College Major  
Y = Likes "Games"

X	Y
Math	Yes
CS	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

-  $P(A \wedge B \wedge C \wedge D, \dots)$



P(X, Y)

$P(X=CS, Y=Yes) = P(CS, Yes)?$   $\frac{1}{4}$   
 $P(Math, No)?$   $\frac{1}{4}$   
 $P(History, Yes)?$  0

CSC872: PAMI – Kazunori Okada (C) 2025

7

7

# Marginal Distribution

- Probability of one event (A) regardless of other events (B)
- Derived from a joint distribution by integrating/summing out

-  $P(A) = P(A \wedge B) + P(A \wedge \sim B)$

X = College Major  
Y = Likes "Games"

- Sum Rule**

$P(x, y, z, \dots) = \sum_y P(x, y, z, \dots)$

$= P(A, B=7) + P(A, B=4) + \dots + P(A, B=2)$   
 $- P(A) = \sum_{B=v_i} P(A \wedge B)$   
 $- P(Yes)? = P(CS, Yes) + P(Math, Yes) + P(History, Yes) = \frac{1}{2}$   
 $- P(CS)? = P(CS, Yes) + P(CS, No) = \frac{3}{8}$

X	Y
Math	Yes
CS	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

CSC872: PAMI – Kazunori Okada (C) 2025

8

8

# Conditional Distribution

- Probability of one event (A) given the other event (B)
- Derived by combining a joint distribution and a marginal distribution

–  $P(A|B) = P(A \wedge B) / P(B)$

–  $P(\text{Yes}|\text{CS}) ? = \frac{P(\text{CS, Yes})}{P(\text{CS})} = \frac{1/4}{3/8} = \frac{2}{3}$

X = College Major  
Y = Likes “Games”

- **Product Rule**

–  $P(A \wedge B) = P(A|B) P(B)$

$P(B \wedge A) = P(B|A) P(A)$

- **Chain Rule (factorization)**

–  $P(A \wedge B \wedge C) = P(A|B \wedge C) P(B|C) P(C)$

$P(A \wedge B \wedge C \wedge D) = P(A|B \wedge C \wedge D) P(B|C \wedge D) P(C|D) P(D)$

X	Y
Math	Yes
CS	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

CSC872: PAMI – Kazunori Okada (C) 2025

9

# Bayes Rule

- **Product rule is symmetric**  $P(A, B) = P(B, A)$

–  $P(A \wedge B) = P(A|B) P(B)$

–  $P(A \wedge B) = P(B|A) P(A)$

- This leads to the following form known as Bayes Rule

–  $P(A|B) = P(B|A) P(A) / P(B)$

CSC872: PAMI – Kazunori Okada (C) 2025

10

10

# How to use Bayes Rule

- $h$ : hypothesis (A)  $\alpha$   $\{B\} \rightarrow \alpha$
- $d$ : evidence (B)  $\{K.B., data\}$
- Compute  $P(h|d)$  given  $P(d|h)$  and  $P(h)$  known

$$p(h|d) = \frac{p(d|h)p(h)}{\sum_{h' \in H} p(d|h')p(h')} = p(d)$$

Posterior probability  $\alpha$   
 Likelihood  $\downarrow$   
 Prior probability  $\leftarrow$   
 Sum over space of hypotheses  $\leftarrow$   
 Decision Making!  $\checkmark$

11

# Bayesian Inference

$$P(h) \xrightarrow{d} P(h|d)$$

Best we can do without 'd'  
 Evidence Perception

- I've got this evidence, what's the chance that this hypothesis is true?
  - a hypothesis 'h'
  - prior probability  $P(h)$ : chance of  $h$  being true without any evidence
  - a new evidence (data) 'd': more information about 'h'
  - posterior probability  $P(h|d)$
- Inference is to compute the posterior!!!
  - I've got a stomach pain, how likely that I have appendicitis?
  - Given NY Dow Jones index goes down by 5%, what is the chance of the amount of foreclosures increasing this month?
  - I see the traffic signal goes green and no cars are visible approaching from right and left, what is the chance of me getting involved in an accident by crossing the intersection?

12

## More on Bayesian Inference

- What this Bayes rule tell us?
- You get high posterior when:
  - Hypothesis is plausible: high  $P(h)$
  - Hypothesis strongly predicts the observed evidence/data: high  $P(d|h)$
  - Evidence/data is very surprising: low  $P(d)$

$$P(h|d) \propto \frac{P(d|h)P(h)}{P(d)}$$

posterior  $\propto$  likelihood  $\times$  prior

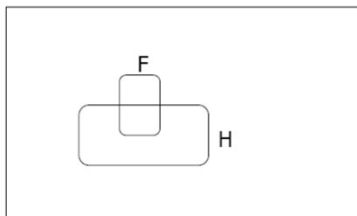
*Handwritten notes:  $\alpha$  (pointing to the proportionality symbol), Likelihood (under  $P(d|h)$ ), prior (under  $P(h)$ ),  $P(d)$  (under  $P(d)$ ),  $KB_1$  (top right).*

CSC872: PAMI – Kazunori Okada (C) 2025

13

13

## Example



H = "Have a headache"  
F = "Coming down with Flu"

$P(H) = 1/10$   
 $P(F) = 1/40$   
 $P(H|F) = 1/2$


- One day you wake up with a headache: You think "Nooo! 50% of flues are associated with headaches so I must have 50-50 chance on coming down with flu.."
- What is prior?  $P(F) = 1/40$
- What is likelihood?  $P(H|F) = 1/2$
- What is posterior?  $P(F|H) = \frac{1/2 \cdot 1/40}{1/10} = 1/8$
- Inference: what is the prob. of flu given a headache evidence

CSC872: PAMI – Kazunori Okada (C) 2025

14

14

## Statistical Modeling

- Nice formulae ! but it **requires** you to provide numerical values for all the distributions: prior and likelihood all that...
- So you **collect data** and **estimate these distributions from the data.**
  - Interview appendicitis patients then ask how many had headache
  - Call randomly 1000 people and ask how many had appendicitis
  - Count times that DowJones went down by more than 5% given the times when foreclosure rate was coming down
- This is called **Statistical Modeling** 

CSC872: PAMI – Kazunori Okada (C) 2025

15

15

## More on Statistical Modeling

- One caveat: **population can be huge!!!**
- You cannot get all data!!!
- Strategy
  - Data = Statistical Samples drawn from a distribution
  - Sampling process assumption:
    - **IID, independent and identically distributed**
    - Estimate the distribution from the samples or
    - Estimate population statistics from sample statistics

CSC872: PAMI – Kazunori Okada (C) 2025

16

16



## Two Approach in Statistical Modeling

- Do you have statistical insights of the underlying distribution? *Shape?*
  - Counting flip-coins ===== Binomial distribution
  - Weight of SF residents ===== Normal distribution
- **No: Non-Parametric** Modeling
  - Histogramming and more...
- **Yes: Parametric** Modeling (2 steps)
  - Find a **parameterize function** for the known distribution
  - Solve a **parameter estimation problem** by fitting the function to data

CSC872: PAMI – Kazunori Okada (C) 2025

17

17

## Joint Distribution is GREAT!!!



- Hassle to model every different distributions, right?
- You have a problem with **N random variables**
- Statistically model a **N-variate joint distribution  $P(X_1, \dots, X_N)$**
- **You can derive the probability of any logical expression from this joint distribution !!!**

- Use various inferential rules

- Sum Rule:  $P(A, B, C, \dots)$  →  $P(A, B), P(A), P(B)$
- Product Rule:  $P(A, B)$  →  $P(A|B)P(B)$
- Chain Rule:  $P(A, B, C, \dots, Z)$  →  $P(A|B, C, \dots, Z) P(B|C, \dots, Z) \dots P(Z)$
- Bayes Rule:  $P(A|B)$  →  $P(B|A)$

$$P(h|d) = \frac{P(d|h)P(h)}{P(d)}$$

- E.g., you can derive  **$P(A), P(B), P(A|B), P(B|A)$**  from  **$P(A, B)$**

CSC872: PAMI – Kazunori Okada (C) 2025

18

18

## Joint Dist. is TERRIBLE!!!



- Well here is a catch... *Curse of dimensionality*
  - **Impossible to create it for more than small number of variables**  $P(x_1, x_2, x_3, \dots, x_n)$ 
    - M = 10 variables =  $2^{10}$  entries for Boolean RV
    - M = 10 variables =  $N^{10}$  entries for N-valued RV (8-bit pixel)
- M: number of variables
- M can be a number of all stocks in the market
  - M can be all possible pixel locations for a video tracker
  - M can indicate # of all stars in the universe
  - M can represent # of all neurons in our brain
  - M can be all users in Tiktok

CSC872: PAMI – Kazunori Okada (C) 2025

19

19

## Statistical Independence

- Some domain knowledge make things better
- If A, B are **statistically independent** then:
  - $P(A \wedge B) = P(A, B) = P(A)P(B)$
  - $P(A|B) = P(A)$
  - $P(B|A) = P(B)$
  - A: scores of riverpool in the next game
  - B: scores of my exam and I am not a riverpool fan ( $A \perp B$ )
- If A, ..., Z are **mutually independent** then:
  - $P(A \wedge B \wedge \dots \wedge Z) = P(A)P(B)\dots P(Z)$
  - “**Naïve Assumption**”  $2^n \leftrightarrow 2n$
  - $2^{26}$  entries --- 52 entries (for Boolean)

CSC872: PAMI – Kazunori Okada (C) 2025

20

20

## More on Statistical Independence

- Given three events  $A, B, C$
- $A$  and  $B$  can be independent *once you know  $C$  is true/false*

- **Conditional Independence**

- $P(A \wedge B|C) = P(A|C)P(B|C)$

- **Applications**

- Naïve Bayesian Classifier
- Bayesian Network

$$P(x_1, x_2, \dots, x_n | d) = P(x_1|d) \cdot P(x_2|d) \cdot \dots \cdot P(x_n|d)$$

$$= \prod_{i=1}^n P(x_i|d)$$

CSC872: PAMI – Kazunori Okada (C) 2025

21

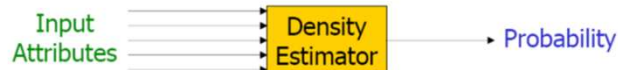
21

## What for?

*Bayesian Reasoning*  
 Posterior Prob  $\leftarrow$  d Prior Prob

- Providing probabilistically sound algorithms for realizing useful functions (PF)

- Classification
- Regression
- Density estimation = *Statistical Modeling*



CSC872: PAMI – Kazunori Okada (C) 2025

22

22

# PF: Bayesian Classification

$H = \{h_1, \dots, h_k, \dots, h_K\}$

- Hypothesis space  $H$  consists of  $N$  classes  $h_i$
- Get data ' $d$ ' as an input feature
- Which class does this input belong to???

- 1) Compute posterior  $P(h_i | d)$  for all  $h_i \in H$
- 2) Take the class ' $i$ ' with the highest posterior of  $d$

$P(h_1|d), P(h_2|d), \dots, P(h_K|d)$

- Formally

$$k^* = \operatorname{argmax}_i P(h_k \in H | d)$$

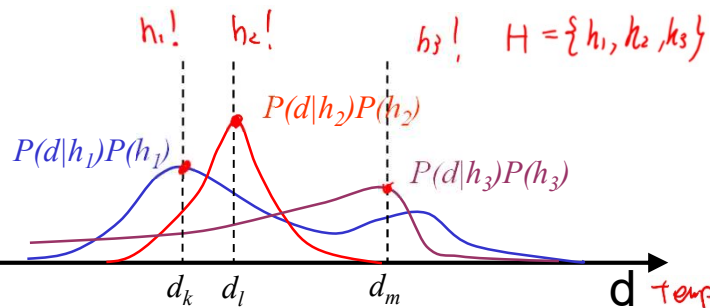
result

# PF: Bayesian Classification

$$k^* = \operatorname{argmax}_k P(h_k | d)$$

$$k^* = \operatorname{argmax}_k \frac{P(d|h_k)P(h_k)}{P(d)} \text{ Bayes Rule}$$

$$k^* = \operatorname{argmax}_k P(d|h_k)P(h_k) \text{ Simplify}$$



## PF: Naïve Bayesian Classifier

- Typically we work with multiple <sup>inputs</sup> features

$$k^* = \operatorname{argmax}_k P(h_k | d_1, \dots, d_M)$$

$$k^* = \operatorname{argmax}_k \frac{P(d_1, \dots, d_M | h_k) P(h_k)}{P(d_1, \dots, d_M)}$$

$$k^* = \operatorname{argmax}_k P(d_1, \dots, d_M | h_k) P(h_k)$$

$$P(d_1 | h_k) P(d_2 | h_k) \dots P(d_M | h_k) P(h_k)$$

- But cannot easily build the joint dist so...

$$k^* = \operatorname{argmax}_k \prod_m P(d_m | h_k) P(h_k)$$

- This is much more **practical** (Data Mining) !

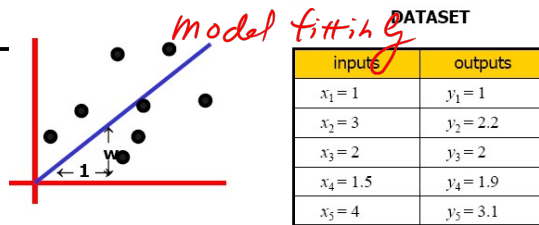
CSC872: PAMI – Kazunori Okada (C) 2025

25

25

## PF: Bayesian Regression

- A function  $f: X \mapsto Y$
- A parameterized model  $y = f(x, w) = wx + \text{noise}$
- Noise is independent, 0 mean, variance  $\sigma^2$
- $P(y | w, x)$  models this function probabilistically
- Get a set of IID data samples  $\{x_m, y_m\}, m=1, \dots, M$
- Find parameters 'w' that make the model fit to the data most**

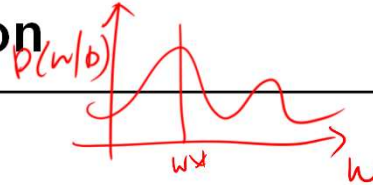


CSC872: PAMI – Kazunori Okada (C) 2025

26

26

# Bayesian Regression



- Bayesian Regression

$$w^* = \operatorname{argmax}_w P(w | (x_1, y_1), \dots, (x_M, y_M))$$

- **Maximum Likelihood Regression**

$$w^* = \operatorname{argmax}_w P((y_1, \dots, y_M) | w, (x_1, \dots, x_M))$$

$$w^* = \operatorname{argmax}_w \prod_m P(y_m | w, x_m) \quad \text{IID naive}$$

$$w^* = \operatorname{argmax}_w \prod_m \exp\left(-\frac{1}{2} \left(\frac{y_m - wx_m}{\sigma}\right)^2\right) \quad \text{Normal}$$

$$w^* = \operatorname{argmax}_w \sum_m -\frac{1}{2} \left(\frac{y_m - wx_m}{\sigma}\right)^2 \quad \text{Log} \quad \text{when } g \text{ is monotonic}$$

$$w^* = \operatorname{argmin}_w \sum_m (y_m - wx_m)^2 \quad \text{Simplify}$$

CSC872: PAMI – Kazunori Okada (C) 2025

27

27

# Bayesian Regression

- So the maximum likelihood regression

$$w^* = \operatorname{argmax}_w P((y_1, \dots, y_M) | w, (x_1, \dots, x_M))$$

- Is the same as this...

$$w^* = \operatorname{argmin}_w \sum_m (y_m - wx_m)^2 = \sum_m \text{error}_m^2$$

- This is known as **least squares problem**

- **Least Squares (Linear) Regression is Maximum Likelihood Regression when we model the PDF function by Gaussian/Normal distribution !!!**

CSC872: PAMI – Kazunori Okada (C) 2025

28

28

# Bayesian Learning

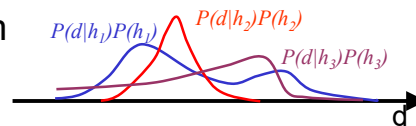
Bayesian KB,

- Bayesian learning is an estimation of **probability distribution** given a data set (as evidence)

- E.g., counting frequency from

X	Y
Math	Yes
CS	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

- E.g., making histogram



- E.g., estimating a function by least squares

CSC872: PAMI – Kazunori Okada (C) 2025

29

29

## Any Other Reasoning with Uncertainty?

- Oh Yes!
  - Fuzzy Logic
  - Dempster-Shafer
  - Three-valued Logic
  - Non-monotonic Logic
- But “probability theory” is the one of the methods that is most sound mathematically
  - If you are gambling using them, you can't be unfairly exploited by an opponent using some other system [di Finetti 1931]
  - [http://en.wikipedia.org/wiki/Bruno\\_de\\_Finetti](http://en.wikipedia.org/wiki/Bruno_de_Finetti)

CSC872: PAMI – Kazunori Okada (C) 2025

30

30

## Summary

- Probability Theory for Reasoning with Uncertainty
  - Bayesian Inference
  - Statistical Modeling
  - Bayesian Classification
  - Bayesian Regression
  - Bayesian Learning
- Next
  - Non-Parametric Statistical Modeling
  - Kernel Density Estimation
  - Mean Shift
  - **FP#1 on Eigenface (2<sup>nd</sup> Session)**