

LECTURE 1 · Intro - Course Overview

Molecular Biology Review

DNA is a polymer of : adenine (A), thymine (T), guanine (G), and cytosine (C).

DNA is replicated and recombined.

5' (upstream) 3' (downstream)

RNA polymerase reads the anti-sense strand (template) → produces coding (sense) strand.

mRNA → codes for proteins

rRNA → translation

tRNA → ribosome structure

long intergenic noncoding RNAs (lincRNAs) → regulation? (e.g. XIST lincRNA silences X chromosomes in women).

miRNAs → bind to RNA (post-transcriptional regulation)

snRNAs → help guide chemical modification of other RNAs.

Genes begin at TSS. The coding sequences: 5' UTR → 3' UTR.
 Regulatory regions { upstream TSS → promoter (proximal reg. regions)
 distal reg. → enhancers → conformational DNA change
 ↳ insulators
 exons
 introns
 alternative splice isoforms

Epigenetics → { DNA methylation → CpG → repression
 histone modification → acetylation

Probability Review

Discrete random variable X can take multiple values from a sample space S .

$P(X=x) \rightarrow$ Probability mass function: $P(X)$

$$| 0 \leq P(X=x) \leq 1 |$$

$$\sum_{x \in S} P(X=x) = \sum_x P(X) = 1$$

Conditional Probability: $P(A|B) = \frac{P(A, B)}{P(B)}$

$$P(A, B) = P(A \cap B)$$

or $P(B|A) \rightarrow$ given A

Chain Rule: $P(A, B, C, \dots) = P(A) \cdot P(B|A) \cdot P(C|A, B) \dots$

or $P(A|B) \rightarrow$ given B what is the prob that it comes from A
 engineering it could also come from B .

Bayes's Rule: $P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$

Independence: $P(A|B) = P(A)$

Conditional independence: $P(A|B, C) = P(A|C)$

A & B are only conditionally independent given C if, given knowledge C occurs, knowledge of whether A occurs provides no info of whether B occurs and vice versa.

Inference problems:

LEARNING: updating posterior probabilities, we want to compute $P(X|Y)$.

From Bayes', we'll need $P(Y|X)$. From chain rule we can factor joint distribution of X, Y . Then use conditional independence to simplify.

EST. MOST LIKELY: maximum a posteriori estimate $\rightarrow \arg \max_x (P(x|Y))$.

Probability distributions

Discrete: take specific countable values (e.g. rolling a dice $\{1, 2, 3, 4, 5, 6\}$).

Continuous: take any value within a range

For any: $P(x) \geq 0$

Discrete:

$$\sum P(x) = 1$$

Continuous:

$$\int P(x) dx = 1$$

Gaussian distribution (Continuous)

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

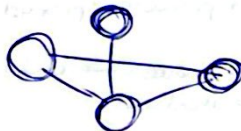
μ = mean
 σ^2 = variance

Graphical Probabilistic Models (GPMs)

Nodes: random variable

Edge: dependencies / relationship

Structure: directed or undirected



Bayesian Networks

\rightarrow edges represent conditional dependencies (Directed acyclic graphs or DAGs)

\rightarrow encodes the joint probability

$$P(x_1, x_2, x_3, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{Parents}(x_i))$$

E.g.: $R(\text{rain}) \rightarrow W(\text{wet})$
 $S(\text{sprinkler}) \rightarrow W(\text{wet})$

$$P(R, S, W) = P(R) \cdot P(S) \cdot P(W | S, R)$$

Markov Networks

\rightarrow undirected; edges represent pairwise relationships

\rightarrow encode joint probability distributions as products of potential functions over cliques (fully connected subsets of nodes)

$$P(x_1, x_2, \dots, x_n) \propto \prod_{\text{cliques}} \phi(c)$$

where $\phi(c)$ is the potential function for clique c .

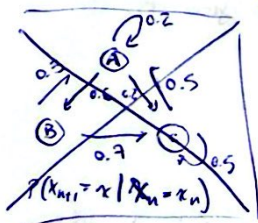
E.g.: social networks

$$A \leftrightarrow B, B \leftrightarrow C, A \leftrightarrow C$$

$$P(A, B, C) \propto \phi(A, B) \cdot \phi(B, C) \cdot \phi(A, C)$$

Normalization:

$$P(A, B, C) = \frac{1}{Z} \phi(A, B) \cdot \phi(B, C) \cdot \phi(A, C)$$



A & B tend to agree

B & C "

A & C are directly influenced

A	B	$\phi(A, B)$
0	0	0.8
1	0	0.2
0	1	0.2
1	1	0.8

Bayes' Rule

$$P(H|E) = \frac{P(E|H) P(H)}{P(E)}$$

- $P(H|E)$: probability of hypothesis after evidence **POSTERIOR**
- $P(H)$: p. of hypothesis **PRIOR**
- $P(E|H)$: how likely E is if hypothesis is true. **LIKELIHOOD**
- $P(E)$: evidence itself probability **NORMALIZATION**

Ex:

$$P(D) = 0.01$$

$$P(P|D) = 0.95 \text{ (true positive)}$$

$$P(P|N) = 0.05 \text{ (false positive)}$$

$$P(N) = 0.99$$

$$P(D|D) = \frac{P(P|D) \cdot P(D)}{P(P|D) \cdot P(D) + P(P|N) \cdot P(N)} = \frac{0.95 \cdot 0.01}{0.95 \cdot 0.01 + 0.05 \cdot 0.99} = \underline{0.6}$$

Markov Chains

Probability of the next state depends only on the current state:

$$P(X_{t+1} | X_t, X_{t-1}, \dots) = P(X_{t+1} | X_t)$$

components

$$\text{States: } S = \{s_1, s_2, s_3, \dots, s_n\}$$

$$\text{Transition prob (T): } T_{i,j} = P(X_{t+1} = s_j | X_t = s_i), T_{i,j} \geq 0 \text{ and } \sum_j T_{i,j} = 1.$$

Ex:

$$s_1 = \text{sunny}$$

$$s_2 = \text{rainy}$$

$$T = \begin{bmatrix} 0.8 & 0.2 \\ 0.4 & 0.6 \end{bmatrix} \quad \left| \begin{array}{l} T_{11} = \text{if today sunny} \Rightarrow 80\% \text{ also sunny} \\ T_{12} = \text{if today sunny} \Rightarrow 20\% \text{ rainy} \\ T_{21} = \text{if today rainy} \Rightarrow 40\% \text{ sunny} \\ T_{22} = \text{if today rainy} \Rightarrow 60\% \text{ rain} \end{array} \right.$$

Hidden Markov Models (HMMs)

1. hidden states (S): true underlying states we aim to observe/infer.
2. observations (O): symbols emitted by S according to a probability dist.
3. Emission probabilities (E): likelihood of each O given S.

Assuming that the Markov property holds and observations are conditionally independent given S.

Example: detecting DNA motifs in a sequence

$$S = \begin{cases} \text{motif (M)} \\ \text{background (B)} \end{cases}, O = \{A, C, T, G\}, \quad \begin{array}{l} T_{M \rightarrow M} = 0.9 \\ T_{M \rightarrow B} = 0.1 \end{array}, \quad \begin{array}{l} E_M(A) = 0.4 \\ E_M(C) = 0.3 \\ E_M(G) = 0.2 \\ E_M(T) = 0.1 \end{array}$$

To infer hidden states we can use:

Forward-Backward Algorithm

Viterbi Algorithm

Maximum Likelihood and Max A post Estimators

MLE

Find θ that max. likelihood

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(\text{data} | \theta)$$

MAP

$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(\text{data} | \theta) P(\theta)$$