

Inference & Representation

Recitation #2

Today's Agenda:

1. Maximum Entropy Distribution
2. Intro to Graphs
3. Hidden Markov Models (HMM)

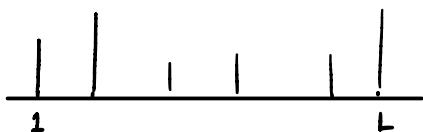
1.1 Variational Principle

\mathcal{X} : domain of the data

$P(\mathcal{X})$: space of probability distribution over \mathcal{X}

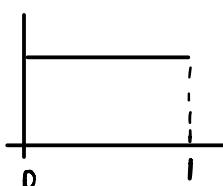
How to define a notion of smoothness for $p \in P(\mathcal{X})$?

Ex. \mathcal{X} discrete set $\mathcal{X} = \{1, \dots, L\}$



→ Uniform distribution maximizes "smoothness"

Ex. $\mathcal{X} = [0, 1]$



1.2 Entropy

We can quantify the smoothness of a probability distribution with its entropy:

$$H_p = - \sum_{i=1}^L p_i \log p_i$$

and under appropriate assumption

$$H_p = - \int_x p(x) \log p(x) dx \quad \text{for } \underline{x \in \mathbb{R}^d}$$

Entropy measures uncertainty

1.3 Lagrange Multiplier

We can use entropy to characterize distributions from a given set of constraints

$$(**) \quad \max_{p \in P(X)} H(p) \quad \text{subject to} \quad \begin{cases} \mathbb{E}_{x \sim p} X = \mu_0 \\ \mathbb{E}_{x \sim p} (X - \mu_0)(X - \mu_0)^T = \Sigma \end{cases}$$

Remark. This constrained optimization problem is not always well-defined.

It depends on the domain X as well as the constraints.

Ex. $X = \mathbb{R}$

$$\max_{p: \mathbb{R} \rightarrow \mathbb{R}} - \int p(x) \log p(x) dx \Leftrightarrow \min_p \int p(x) \log p(x) dx$$

$$\text{s.t. } \begin{cases} p(x) \geq 0 & \forall x \\ \int p(x) dx = 1 \\ \int x p(x) dx = \mu \\ \int (x - \mu)^2 p(x) dx = \sigma^2 \end{cases}$$

→ Constrained optimization problem : we introduce **Lagrange multipliers** :

$$\begin{aligned} L[p] &= \int p(x) \log p(x) dx + \lambda (\int p(x) dx - 1) \\ &\quad + \beta (\int x p(x) dx - \mu) \\ &\quad + \sigma (\int (x - \mu)^2 p(x) dx - \sigma^2) \end{aligned}$$

$$\frac{\partial L}{\partial p(x)} = 0 = \log p(x) + 1 + \lambda + \beta x + \sigma (x - \mu)^2$$

$$\log p(x) = a - bx - cx^2$$

$$p(x) \propto e^{-bx - cx^2} \quad \text{with } b, c \text{ adjusted s.t.}$$

$$\mathbb{E} p(x) = \mu, \mathbb{E} (x - \mu)^2 = \sigma^2$$

Same thing happens in d-dimensions :

the distribution of max entropy with $\mu \in \mathbb{R}^d$, $\Sigma \in \mathbb{R}^{d \times d}$ is $\mathcal{N}(\mu, \Sigma)$

1.4 Some Exercises :

Let $\mathcal{X} = \{x_1, \dots, x_N\}$ be a discrete set and let $P(\mathcal{X})$ be the space of probability distributions defined over \mathcal{X} . Define the entropy :

$$H(p) := - \sum_{i=1}^n p_i \log(p_i)$$

Ex1. We want to show its maximum is attained at uniform distribution, with maximum entropy $\log N$.

Lagrange Multipliers : ...

□

Now let us move from a discrete to continuous domain, by setting $\mathcal{X} = [0, R]$. Let $P(\mathcal{X})$ denote the space of probability distribution admitting a density $p(x)$, $x \in \mathcal{X}$, and define the Shannon entropy as

$$H(p) = - \int_{\mathcal{X}} p(x) \log p(x) dx$$

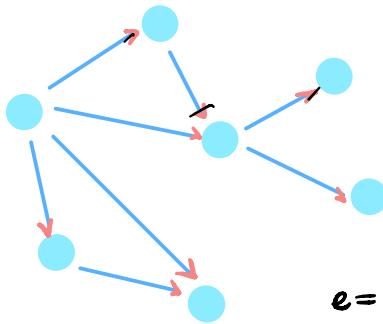
Ex2. Show that the maximum entropy distribution in $P(\mathcal{X})$ is the uniform measure in \mathcal{X} , with entropy $\log R$.

Similarly to discrete domain,

Lagrange multipliers , ...

□

2. Graph Theory



$$G = (V, E)$$

$V = \{v_1, \dots, v_n\}$ set of vertices

$E = \{e_1, \dots, e_m\}$ set of edges

$e = (v_i, v_j) \in E$

$v_i \sim v_j$

Two types of graphs:

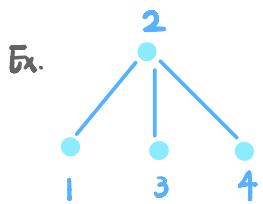
- generic : directed

$$E: \{(i, j), i, j \in V\}$$

- symmetric/ undirected : $(v_i, v_j) \in E \Leftrightarrow (v_j, v_i) \in E$

$i \rightarrow j$ ordered

Adjacency matrix : $A_{ij} = \begin{cases} 1 & \text{if } v_i \rightarrow v_j, A \in \{0, 1\}^{n \times n} \\ 0 & \text{o.w.} \end{cases}$



$$A = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

Weighted graph : graph with weights on edges

$$W = (w_{ij}) \in \mathbb{R}^{n \times n}$$

Some concepts in graph theory

* path $v_1, v_2, \dots, v_k \in V : v_i \sim v_{i+1}, v_i \neq v_j$

* cycle path with $v_k = v_1$

* connected A subset of vertices is connected
if \exists path between each pair of vertices.

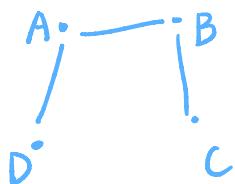
* connected component maximal connected subset

* tree symmetric + connected + acyclic

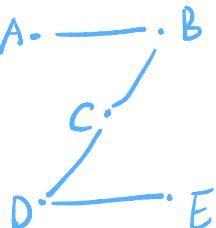
(*) if G is a tree. $|E| = N - 1$

* spanning tree subgraph $T \subseteq G : V(T) = V(G)$ + T is a tree.

Ex. Tree



Spanning Tree :



* Clique $H \subseteq G$: there is an edge between each pair of vertices in H

- maximal clique There is no other clique containing it.

- clique number of G $c(G) = \max_{\text{clique } C \subseteq G} |C|$

* Complementary notion

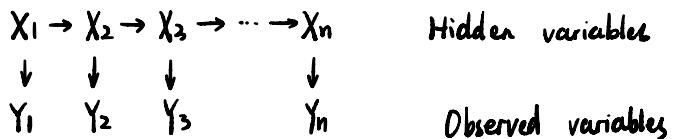
- independent set : no two nodes share an edge

\Leftrightarrow a clique of G^c

- independence number size of maximal independent set

The independence number of G^c = clique # of G

3. Hidden Markov Model (HMM)



Ex1. Hidden : events in the atmosphere
Observation : daily weather

Ex2. Part-of-Speech tagging

Y_i : words

X_i : tags (n./verb./adj)

Ex3. Speech recognition

Y_i : sounds, recordings, audio

X_i : words, letters

Densities of HMM

$$P(x, y) = P(x_1) \prod_{i=2}^n P(x_i | x_{i-1}) \prod_{i=1}^n P(y_i | x_i)$$

\uparrow \uparrow
parent of x_i parent of y_i

Conditional Independence implied by HMM

$$x_i \perp x_j | x_{i-1} \quad \forall j < i-1 \quad \leftarrow \text{Markov property}$$

$$y_i \perp \{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n, y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n\} | x_i$$

Why "Markov"?

Future \perp Past | Present

Ex. Compute Probabilities in HMM

TBC...