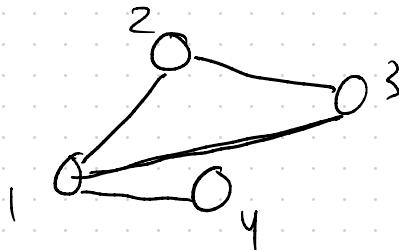


Some Probability Review

Network Analysis



$$\left. \begin{array}{l} V = \text{webpages} \\ E = \text{links in webpages} \end{array} \right\}$$

$$f: V \rightarrow \mathbb{R}$$

View f as a probability function. We can construct f as a centrality metric.

Most classical example page-Rank:

$$C_R(u) = \sum_{\substack{v \in V \\ \{u, v\} \in E}} \frac{C_R(v)}{\deg(v)}$$

How do we determine C_R ? Start with uniform

$$f(u) = \frac{1}{|V|} \quad \text{for all } u \in V, \text{ then}$$

look at Markov Process:

$\hookrightarrow X_0 = \text{first webpage}$

$X_1 = \text{second webpage chosen by clicking link.}$

- ~ Probability of clicking link independent of previous path
- ~ can only move to a new page if \exists link
- ~ Probability does not depend on how long we've been clicking.

Example

③ Metropolis-Hastings Algorithm and CTD^o

Define a matrix $P = (P_{uv}) \in \mathbb{R}^{n \times n}$ $n = \# \text{ vertices}$
on connected
graph

With

$$P_{uv} = \begin{cases} \frac{1}{d} \min\left\{1, \frac{\pi(u)}{\pi(v)}\right\} & \{u, v\} \in E \\ 0 & u \neq v \\ 1 - \sum_{i \neq v} P_{iv} & u = v \end{cases}$$

$\pi(u)$
 $\pi(v)$
only
ratio
→
i.e. know
more likely
link to click.

$\{u, v\} \in E$
 $u \neq v$
 $\pi: V \rightarrow \mathbb{R}$
probability
with
 $\pi(u) > 0$

Thm 2.43

P is transition matrix of

- aperiodic $\rightarrow \gcd(\{k \in \mathbb{N} \mid (P^k)_{uu} > 0\}) = 1$
- irreducible $\rightarrow \forall u, v \in V, \exists k \text{ with } (P^k)_{u,v} > 0$

Markov process

with unique stationary distribution π . $P\pi = \pi$

Why care about stationary distributions?

Theorem 2.40

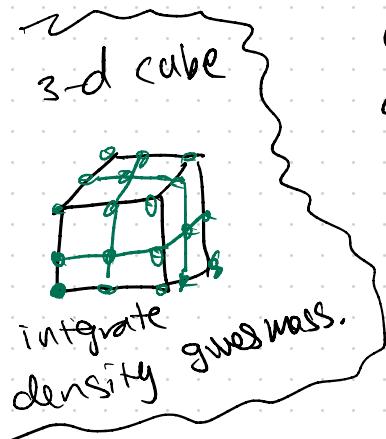
For any probability distribution $f: V \rightarrow \mathbb{R}$,

$$\lim_{k \rightarrow \infty} p^k f = \pi,$$

For example $f = \text{uniform}$, $\pi = CR$ page rank

Toy example Suppose we want to sample points

in the n -dimensional cube $\{-1, 1\}^n$,
and integrate a density / probability function
on this cube C



Idea $\int_C f = \sum_i \rho(x_i) \text{vol}(C_i)$

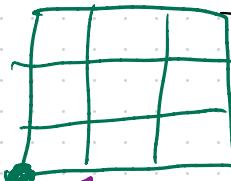
Split C into smaller cubes C_i

① Problem a uniform sampling x_i doesn't
See unit ball \rightarrow what if density is
concentrated inside the unit ball.

② Problem #sample points for uniform sampling grows exponentially

Solution given a probability density f

Sample x_i according to $f(i)$,
with f_i concentrated on unit



probability of going here chosen by f . Create.

Long term $P^k f \approx \pi$ = "correct" distribution of data.

Art of
Data
Science

- Need to choose initial f to match our problem
 - we don't know how many k to sample.
 - ↳ Let run for $k=100$ (burn-in period)
- Start sampling at $k=101 \rightarrow$

Why use Markov Chains?

- ↳ Limit curse of dimensionality
- ↳ Don't need to know exact distribution of data
- ↳ Only depend on current state, not all previous states [need less data to make future predictions]

Types of problems include

For example

+ Given a sample graph, describe transition matrix. + Given transition matrix + distribution, determine if aperiodic, stationary, irreducible.

Machine Learning + Probabilistic functions

Recall normal distribution, $x \in \mathbb{R}^n$

$$\Phi(x | \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n \det(\Sigma)}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

When $n=1$, $\Sigma = \sigma^2 = \text{variance}$
 $\mu = \text{mean}$

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

$f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}$, $\theta = \langle a, b \rangle \in \mathbb{R} \times \mathbb{R}^D$

$$f_\theta(x) = a + \langle b, x \rangle. \quad \text{Linear model}$$

Thm 3.12

$$P_\theta(y | x) := \Phi(y | f_\theta(x), \sigma^2)$$

normal distribution with
variance σ^2 and mean $f_\theta(x)$.

$$y = f_\theta(x) + \varepsilon \quad \text{for } \varepsilon \sim N(0, \sigma^2)$$

Thm 3.13

Nonlinear Compose with $\phi: \mathbb{R}^D \rightarrow \mathbb{R}^P$.

Now $\theta \in \mathbb{R}^P$

$$P_\theta(y|x) = \Phi(y|\theta^T\phi(x), \sigma^2)$$

Thm 3.14 We can make σ^2 variance random too!

$$\hat{\theta}_{ML} = \Sigma^+ Y, \quad \hat{\sigma}_{ML}^2 = \frac{1}{n} \|Y - \Sigma \hat{\theta}_{ML}\|^2$$

The Curse of Dimensionality: $n \ll D$

All theorems require $\text{rank } \approx D$ or P
nonlinear.

Add MAP Maximum A posteriori $\theta \in \mathbb{R}^P$

prior $\theta \sim N(\mu, \Sigma)$ $\mu \in \mathbb{R}^P, \Sigma \in \mathbb{R}^{P \times P}$

$\Sigma \in \mathbb{R}^{n \times P}$

MAP

Theorem 3.15

$$\hat{\theta}_{MAP} = (\underline{\Omega^T \Omega + \sigma^2 \Sigma^{-1}})^{-1} (\underline{\Omega^T y + \sigma^2 \Sigma^{-1} m})$$

from
idea of pseudoinverse

adds randomness
make full rank
matrix), i.e. invertible

adds randomness
to account for
random θ .

Recall Proof uses Bayes' Theorem + Def'n of $\hat{\theta}$.

$$\hat{\theta} = P(\theta | X, Y) = P(\theta) \frac{P(Y | X, \theta)}{P(Y | X)}$$

$$P(\theta | Y) = \frac{P_\theta(\theta) P(Y | \theta)}{P_Y(Y)}$$

independent of X

both depend on X .

Moral of Story Probabilistic approach is

a key to obtaining unique solutions
to regression problems

↳ as in verified
computer output is
correct/only one
possible output.