

01 - Autoregressive

P₀₃ Motivation : AR model for sequence of observations X_1, X_2, \dots, X_T
 ↳ continuous observations at discrete time-step

P₀₄ AR model : AR model AR(p) : $X_t = c + \sum_{i=1}^p \varphi_i X_{t-i} + \varepsilon_t \Leftrightarrow \text{Lag}[x_{t-1}, \dots, x_{t-p}]_t \Rightarrow (x_t)_t$

$\left\{ \begin{array}{l} c: \text{constant} \\ \varphi_i: \text{parameters} \end{array} \right. / \begin{array}{l} X_{t-i}: \text{Lagged value at time } t \\ \varepsilon_t: \varepsilon_t \sim N(0, \sigma^2) \text{ white noise} \end{array}$
 shared through time

$\Rightarrow p(X_t | X_{t-1}, \dots, X_{t-p}) \sim N(c + \sum_{i=1}^p \varphi_i X_{t-i}, \sigma^2)$

P₀₆ AR model property : Mean function : $\mu(t) = E[X_t] = \frac{c}{1 - \sum_{i=1}^p \varphi_i}, \forall t$

Autocovariance : $\gamma(t, i) = \text{cov}(X_t, X_{t-i})$

Person Auto correlation function : $\rho(t, i) = \frac{\text{cov}(X_t, X_{t-i})}{\sqrt{\text{Var}(X_t)} \sqrt{\text{Var}(X_{t-i})}} \in [-1, 1]$

P₀₈ φ Parameter Learning : $\begin{bmatrix} \varphi_1 \\ \varphi_2 \\ \vdots \\ \varphi_p \end{bmatrix} = (X^T X)^{-1} X^T y$ / $X = \begin{bmatrix} X_{p-1} & \dots & X_0 \\ X_p & \dots & X_1 \\ \dots & \dots & \dots \end{bmatrix}$ $y = \begin{bmatrix} X_p \\ X_{p+1} \\ \dots \end{bmatrix}$

train | val | test → time

P₁₂ AR model - Stationarity : A process is stationary if : $-1 < \varphi_i < 1$

good modeling assumption $E[X_t] = E[X_{t-i}] = \mu, \forall t \forall i \Leftrightarrow \text{constant mean}$

$\Leftrightarrow \text{cov}(X_t, X_{t-i}) = \gamma_i, \forall t \forall i \Leftrightarrow \text{depend only on } i, \text{not } t$

$E[|X_t|^2] < \infty, \forall t$

$\Rightarrow \gamma_i = \text{cov}(X_t, X_{t-i}) = \text{cov}(X_{t-i}, X_t) = \gamma_{-i}$ 特殊值解法

$$X_t = (1 - 0.1L) X_{t-1} + \varepsilon_t$$

$$(1 - 0.1L) = 0 \quad L = 10 \quad |L| > 1$$

P₁₄ Parameter Learning under Stationary : Yule Walker

$$\gamma_0 = \sum_{j=1}^p \varphi_j \gamma_{-j} + \sigma^2$$

$$\gamma_1 = \sum_{j=1}^p \varphi_j \gamma_{1-j}$$

...

$$\gamma_p = \sum_{j=1}^p \varphi_j \gamma_{p-j}$$

$$\begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_{p-1} \\ \gamma_p \end{bmatrix} = \boxed{\begin{bmatrix} \gamma_0 & \gamma_{-1} & \dots & \gamma_{2-p} & \gamma_{1-p} \\ \gamma_1 & \gamma_{-2} & \dots & \gamma_{3-p} & \gamma_{2-p} \\ \vdots & & \ddots & & \vdots \\ \gamma_{p-2} & \gamma_{-p+1} & \dots & \gamma_{p-1} & \gamma_{-1} \\ \gamma_{p-1} & \gamma_{-p+2} & \dots & \gamma_1 & \gamma_0 \end{bmatrix}} \begin{bmatrix} \varphi_1 \\ \varphi_2 \\ \vdots \\ \varphi_{p-1} \\ \varphi_p \end{bmatrix}$$

Yule-Walker matrix

02 – Markov Chains

P03 Markov Chains : Markov Chains is a sequence of X_1, X_2, \dots, X_T which fulfills the **Markov property**

$$P(X_t | X_1, \dots, X_{t-1}) = P(X_t | X_{t-1}) \quad / \text{discrete time } t \text{ and discrete } X_t$$

Joint Distribution: $P(X_1=i_1, \dots, X_T=i_T) = P(X_1=i_1) \pi_{t=1}^{T-1} P(X_{t+1}=i_{t+1} | X_t=i_t)$

P04 General case : $P(X_1=i) = \pi_i$ and $P(X_{t+1}=j | X_t=i) = A_{ij}^{(t+1)}$

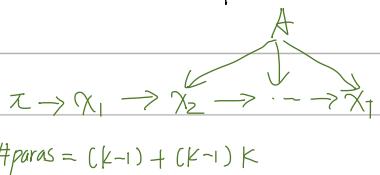
prior distribution: $\pi \in R^K$ / Transition matrix $A^{(t)} \in R^{K \times K}$

$$\begin{array}{c} A^{(2)} \\ \downarrow \\ \pi \rightarrow x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_T \\ A^{(1)} \\ \downarrow \end{array}$$

Joint Distribution: $P(X_1=i_1, \dots, X_T=i_T) = \pi_{i_1} \times A_{i_1, i_2}^{(2)} \times \dots \times A_{i_{T-1}, i_T}^{(T)}$ #params = $(K-1) + (K-1)K(T-1)$

P05 Stationary : $P(X_1=i) = \pi_i$ and $P(X_{t+1}=j | X_t=i) = A_{ij}$ / $A^{(t)} = A$, doesn't depend on t

Joint Distribution: $P(X_1=i_1, \dots, X_T=i_T) = \pi_{i_1} \times A_{i_1, i_2} \times \dots \times A_{i_{T-1}, i_T}$



#params = $(K-1) + (K-1)K$

P0 Parameter Learning : Given a set of $\{X_{1:T_n}\} \rightarrow$ learn π and A Maximum-likelihood

$$\max_{\pi, A} \log P(\alpha | \Pi) = \sum_{k=1}^K L(k) \log (\pi_k) + \sum_{i=1}^K \sum_{j=1}^K N(i,j) \log (A_{ij})$$

$$\text{subject to } \sum_k \pi_k = 1, \sum_j A_{ij} = 1 / A_{ij} = \frac{N(i,j)}{\sum_j N(i,j)}, \pi_k = \frac{L(k)}{\sum_k L(k)}$$

① compute $A_{ij}(n) = P(X_{t+n}=j | X_t=i) \Leftrightarrow$ probability n step from i to j

$$A(n) = P(X_{t+n}=j | X_t=i) = \sum_{k=1}^K A_{kj} A_{ik}(n-1) = A(n-1)A \Rightarrow \boxed{A(n) = A^n}$$

Chapman - Kolmogorov $\Rightarrow \boxed{A(m+n) = A(m)A(n)}$

(row vector)

② compute $\pi_j(t) = P(X_t=j) \Leftrightarrow$ probability t step reaching j

$$\pi_j(t) = P(X_t=j) = \sum_{i=1}^K A_{ij} \pi_i(t-1) = A \pi(t-1) \Rightarrow \boxed{\pi(t) = \pi A^{(t-1)}}$$

$\stackrel{\text{注}}{\rightarrow} A^\top \pi$

P03 Motivation : Probabilistic latent variable models for sequence observations: X_1, X_2, \dots, X_T
 discrete time step / discrete or continuous observations

problem ①: X_t can't capture all relevant information of $[x_1, \dots, x_t]$

②: State is unknown, can only be observed indirectly
 ↳ not observed / latent state Z_t

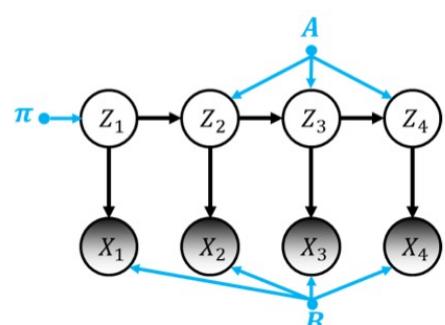
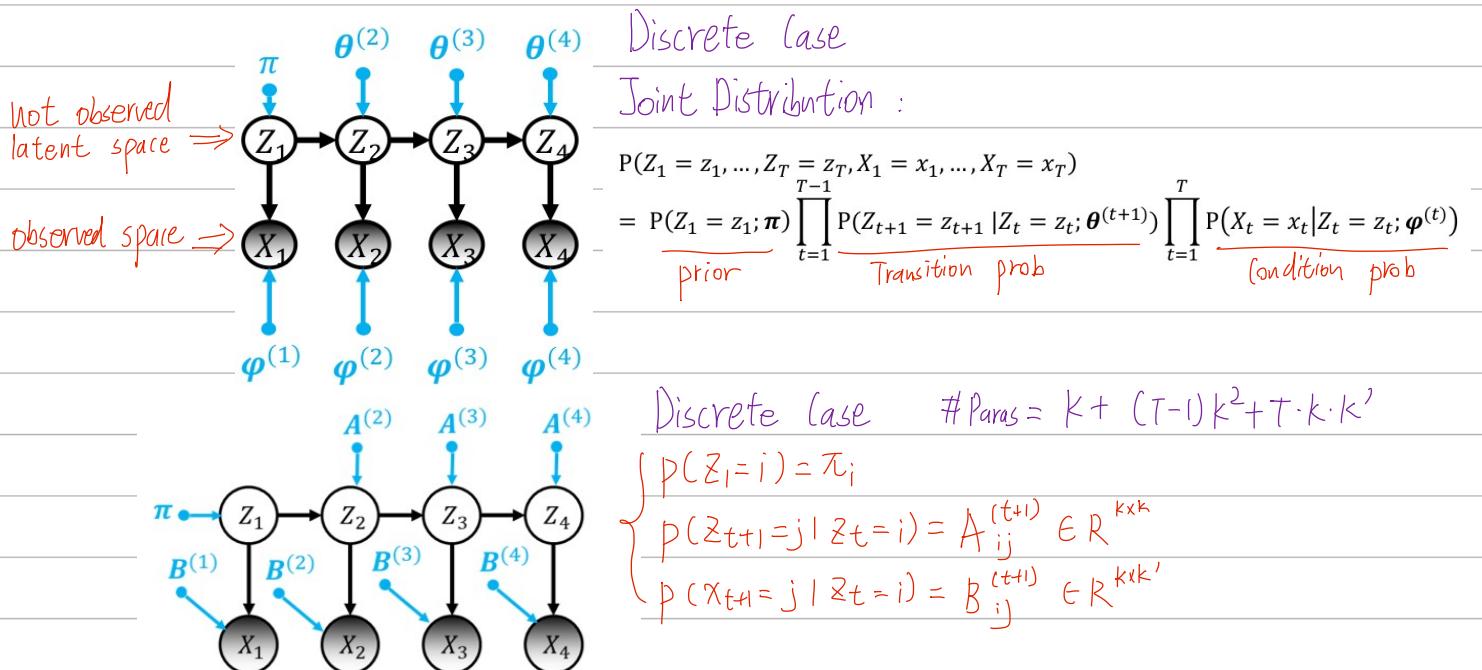
Markov - property: $[z_1, \dots, z_T]$

No Markov - property: $[x_1, \dots, x_T]$

P05 HMM: Composed of a sequence of hidden /latent variables $[z_1, \dots, z_T]$ and observed variables $[x_1, \dots, x_T]$

and $[z_1, \dots, z_T]$ fullfill Markov - property: $P(z_{t+1} | z_t, \dots, z_1) = \underbrace{P(z_{t+1} | z_t)}_{\text{Transition probabilities}}$

x_t depend only on z_t and not fullfill Markov - property: $P(x_{t+1} | z_1, \dots, z_t, x_1, \dots, x_T) = \underbrace{P(x_{t+1} | z_t)}_{\text{Emission probabilities}}$



Share Parameters: # para = $k + k^2 + k \cdot k'$

Joint Distribution :

$$P(Z_1 = z_1, \dots, Z_T = z_T, X_1 = x_1, \dots, X_T = x_T) = P(Z_1 = z_1; \pi) \prod_{t=1}^{T-1} A_{z_t z_{t+1}} \prod_{t=1}^T B_{z_t x_t}$$

Transition Emission

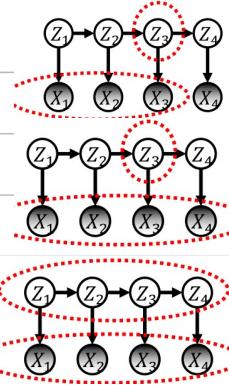
P12 Inference : ① Filtering $P(Z_t | X_{1:t})$

Online setting

② Smoothing $P(Z_t | X_{1:T})$

offline setting

③ MAP Inference $\arg \max_{Z_{1:T}} P(Z_{1:T} | X_{1:T})$



P13 Algorithm : ① Forward : Goal : $P(Z_t=k | X_{1:t}) = \frac{\alpha_t(k)}{\sum \alpha_t}$

$$\alpha_t(k) = \pi_k \odot B_{kx_1} / \alpha_{t-1}(k) = B_k(x_{t+1}) \sum_{j=1}^K A_{jk} \alpha_{t-1}(j)$$

$$\forall t \text{ for } \Rightarrow \alpha_{t+1}(k) = B_k(x_{t+1}) \odot (A' \alpha_t)$$

② Forward - Backward : Goal : $P(Z_t=k | X_{1:T}) \propto \alpha_t(k) \cdot \beta_t(k) = \frac{\alpha_t(k) \beta_t(k)}{\sum \alpha_t \cdot \beta_t}$

$$\beta_T(k) = 1 \quad / \quad \beta_t(j) = \sum_{k=1}^K A_{jk} B_{kx_{t+1}} \beta_{t+1}(k)$$

$$\beta_t^{(l)} = A C B_{k_l}(x_{t+1}) \odot \beta_{t+1}^{(l)}$$

③ MAP Inference : Goal : $\arg \max_Z P(Z_{1:T} | X_{1:T}) = \arg \max_Z \log [P(Z_1) p(X_1 | Z_1)] + \sum_{t=2}^T \log [P(Z_t | Z_{t-1}) p(X_t | Z_t)]$



- weights of edges connected to the Start node: $-\log[\Pr(Z_1=j) \Pr(X_1|Z_1=j)] = g(j)$

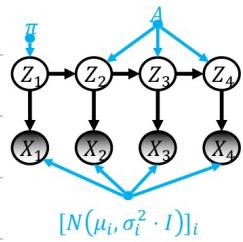
- weights of the intermediate layers: $-\log[\Pr(Z_t=j|Z_{t-1}=i) \Pr(X_t|Z_t=j)] = f(i,j)$

- weights of the edges connected to the End node: 0

Sum of edge weights = $-\log P(Z_{1:T}, X_{1:T})$

P25 Parameter Learning : model parameters : $\theta = \{\pi, A, B\}$ / Goal : solve $\max_{\theta} \log P(X|\theta)$ / EM-step

P30 Continuous Data :



Continuous case, e.g. time-series segmentation

$$\{ P(Z_1=i) = \pi_i$$

$$\{ P(Z_{t+1}=j | Z_t=i) = A_{ij}$$

$$\{ P(X_t=x_t | Z_t=i) = N(x_t | \mu_i, \sigma_i^2 \cdot I)$$

parameter Learning P33

64 - NN Approaches

Po3 Word representation : ① One-hot encoding (word are independent of each other)
 Distributional hypothesis : words in similar contexts should have similar meanings
 similar words should have vector near each other

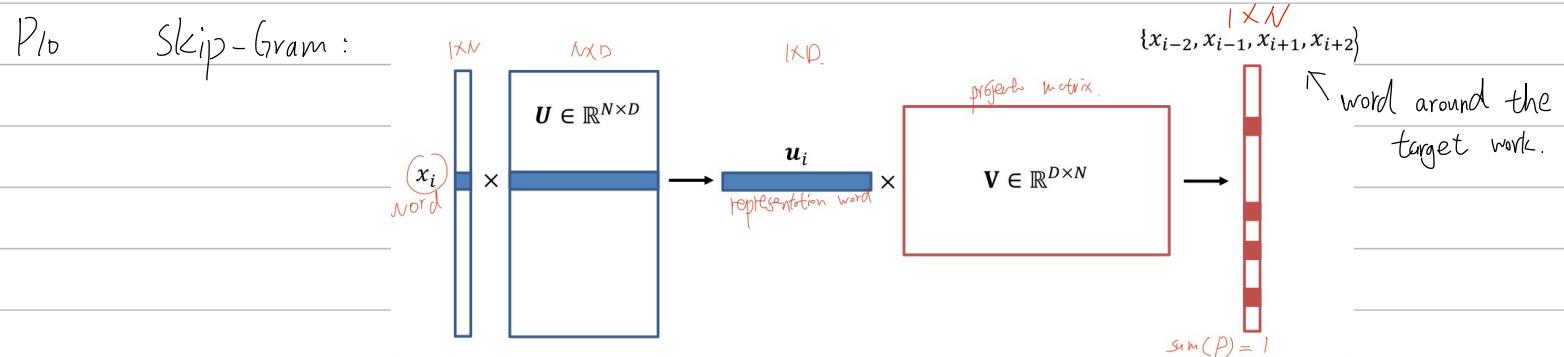
Po4 Co-occurrence Matrix M : aware of context \Rightarrow count word appearance



Words $\{x_1, \dots, x_N\}$, window size L , $\{x_{i-L}, \dots, x_{i-1}, x_{i+1}, \dots, x_{i+L}\}$

SVD : $M = U \Sigma V^T$ / first D columns of U = D -dim word vectors

Po4 Word2Vec : use NN to get word vectors / predict words based on context
 Two approaches: CBow / Skip-gram.



$S = \{x_{i-L}, \dots, x_{i-1}, x_{i+1}, x_{i+L}\}$, θ is the model parameters

$$\max_{\theta} \mathbb{E}[P(S|x_i, \theta)] = \min_{\theta} (-\mathbb{E}[P(S|x_i, \theta)])$$

where $P(S|x_i, \theta) = \prod_{x_k \in S} P(x_k|x_i, \theta)$ corresponding embedding

and $P(x_k|x_i, \theta) = \text{softmax}(u_i V)_k$ * choose $U = V \Rightarrow$ less expressiveness / parameters

Inefficient for large vocabularies

* Alternative: Negative Sampling \rightarrow BCE

P15 RNN : One to Many / Many to One / Many to Many

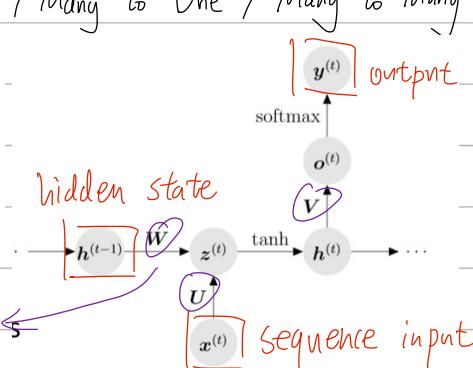
$$z^{(t)} = Wh^{(t-1)} + Ux^{(t)} + b$$

$$h^{(t)} = \tanh(z^{(t)})$$

$$o^{(t)} = Vh^{(t)}$$

$$\hat{y}^{(t)} = \text{softmax}(o^{(t)})$$

Weights share over time



The gradient $\frac{\partial L}{\partial h^{(t)}}$ depends on future times

$$\left(\frac{\partial L}{\partial h^{(t)}} \right) = V^T (\hat{y}^{(t)} - y^{(t)}) + W^T \text{diag} \left(1 - (h^{(t+1)})^2 \right) \frac{\partial L}{\partial h^{(t+1)}}$$

impact may vanish or explode

\rightarrow RNN can't retain informations for many steps

skipping / avoid gradient problem

P₂₁ GRU: Add Gate mechanism

$$\begin{aligned} \mathbf{z}^{(t)} &= \sigma(\mathbf{W}_z [\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}]) && \text{train parameter} \\ \mathbf{r}^{(t)} &= \sigma(\mathbf{W}_r [\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}]) && \text{determine gate which to learn / not to learn.} \\ \tilde{\mathbf{h}}^{(t)} &= \tanh(\mathbf{W} [\mathbf{r}^{(t)} \odot \mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}]) && \text{Simple RNN update - gives candidate state} \\ \mathbf{h}^{(t)} &= (1 - \mathbf{z}^{(t)}) \odot \mathbf{h}^{(t-1)} + \mathbf{z}^{(t)} \odot \tilde{\mathbf{h}}^{(t)} && \text{new information} \end{aligned}$$

How much to take from previous state vs. candidate state

P₂₂ LSTM: Add cell state $c^{(t)}$

$$\begin{aligned} \mathbf{f}^{(t)} &= \sigma(\mathbf{W}_f [\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}]) && \text{Forget gate} \\ \mathbf{i}^{(t)} &= \sigma(\mathbf{W}_i [\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}]) && \text{Input gate} \\ \mathbf{o}^{(t)} &= \sigma(\mathbf{W}_o [\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}]) && \text{Output gate} \\ \mathbf{c}^{(t)} &= \mathbf{f}^{(t)} \odot c^{(t-1)} + \mathbf{i}^{(t)} \odot \tanh(\mathbf{W} [\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}]) && \text{Simple RNN update - LSTM treats it as an input} \\ \mathbf{h}^{(t)} &= \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)}) && \end{aligned}$$

Update hidden state (now the output) using a cell state

P₂₅ Non-Recurrent Model: Weight in RNN rely on the whole history
use Convolutional Neural Networks (ConvNets)

P₂₈ Wave Net: use 1-D convNet

Causal convolutions: ensure the output only depends on the past

Dilated convolutions: skip some input to increase the receptive field

P₃₀ Transformers & Attention: Attention \leftrightarrow Weighting elements

positional Encoding: represent the token order

Training and Inference: During training, all masked embeddings are input

During inference, decode is step by step

GPT: \downarrow unsupervised pre training to predict next token
task-specific fine tuning

05 - Temporal Point Processes

P03 Event Data : Discrete events in continuous time

task: when next event happen / how many events happen

different time series : (synchronous) event data: regular intervals

(Asynchronous) event data: Irregular intervals

P05 Temporal Point Process: (TPP) \Rightarrow probabilistic model \Rightarrow describe distribution of discrete event
 \Rightarrow generative model for variable-length sequence $t = \{t_1, \dots, t_n\}$ in interval $[0, T]$
 \Rightarrow Likelihood Function $p(\{t_1, \dots, t_n\})$

P06 predict next event: ① conditional density as $p^*(t) := p(t | H(t))$, $H(t)$ is history: $H(t) = \{t_j < t\}$

② Cumulative distribution function (CDF)

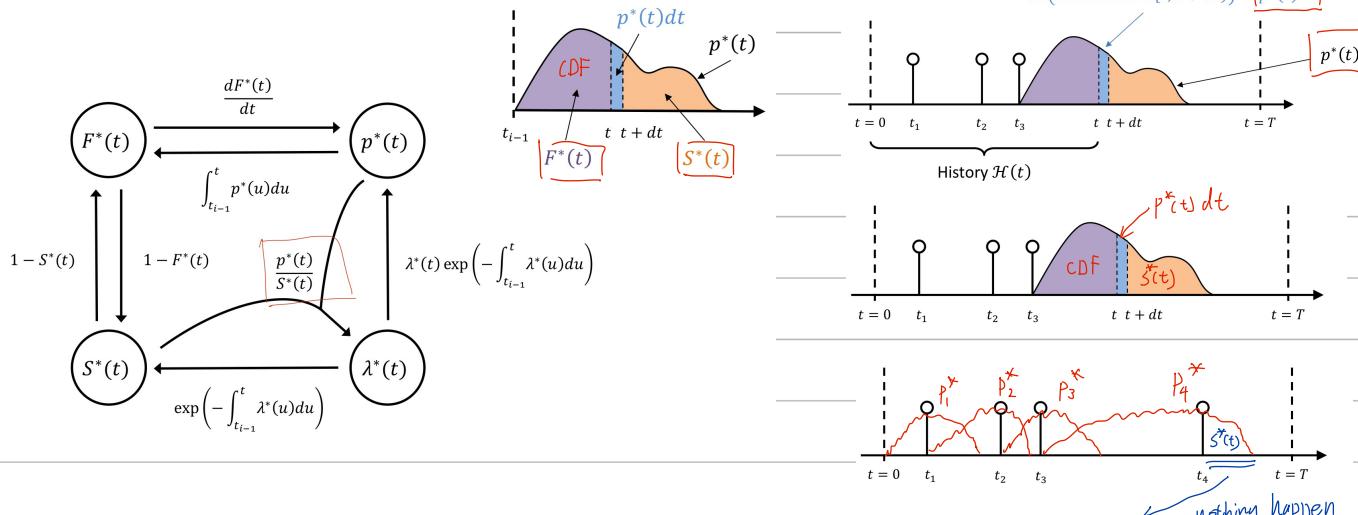
$F^*(t) = \int_{t_{i-1}}^t p^*(u) du \Leftrightarrow$ Probability that next event happen in $[t_{i-1}, t)$

Survival function

$S^*(t) = 1 - F^*(t) = \int_t^\infty p^*(u) du \Leftrightarrow$ Probability that next event happen after t

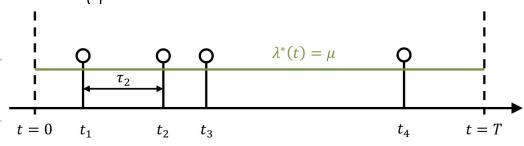
③ Conditional intensity

$\lambda^*(t) dt =$ probability of event in $[t, t+dt)$ given no event in $[t_{i-1}, t) = \frac{p^*(t) dt}{S^*(t)}$
 $=$ Expected # of events $\int_t^\infty \lambda^*(u) du$



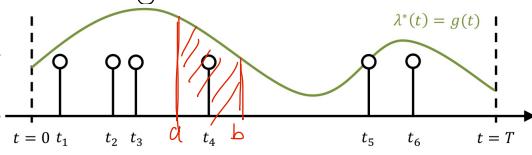
P10 Likelihood of Sequence: $p(\{t_1, t_2, t_3, t_4\}) = p^*(t_1) p^*(t_2) p^*(t_3) p^*(t_4) S^*(T)$
 $= \lambda^*(t_1) \lambda^*(t_2) \lambda^*(t_3) \lambda^*(t_4) \exp\left(-\int_0^T \lambda^*(u) du\right)$

P13 Homogenous Poisson Process : (HPP) Constant intensity : $\lambda^*(t) = \mu$



$$\text{Inter-event times : } p^*(t) = \mu \exp(-\int_{t_{i-1}}^t \mu du) = \mu \exp(-\mu(t - t_{i-1}))$$

P15 Inhomogeneous Poisson Process : (IPP) Intensity : $\lambda^*(t) = g(t) \geq 0$



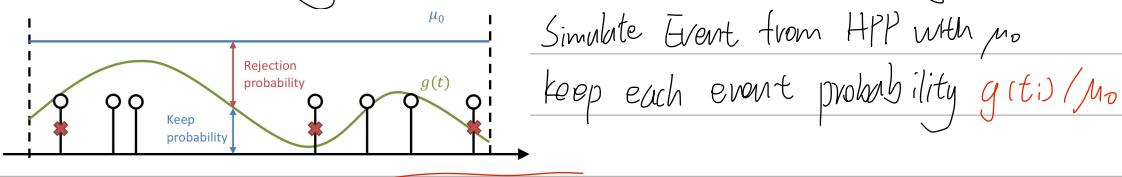
↳ independent of History / Capture global trend

$$\# \text{ Expected Events : } N([a, b]) \sim \text{Poisson} \left(\int_a^b g(t) dt \right)$$

Combination of two IPPs : $g(t) + h(t)$

$$p(x=k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

P18 Simulating an IPP : Find upper bound $M_0 \geq g(t)$



Simulate Event from HPP with μ_0

keep each event probability $g(t_i) / M_0$

P19 Hawkes Process : Self-exciting Process : $\lambda^*(t) = \mu + \alpha \sum_{t_j < t} k_w(t - t_j)$

↳ Triggering kernel = $\exp(-w(t - t_j))$

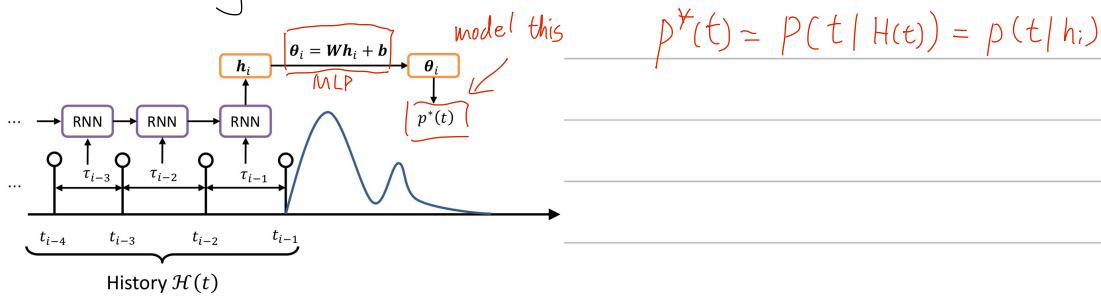


Intensity depend on history

P20 Parameters Training : MLE :

$$= \sum_{i=1}^n \log \lambda^*(t_i) - \int_0^T \lambda^*(u) du$$

P22 Modeling TPPs with RNN

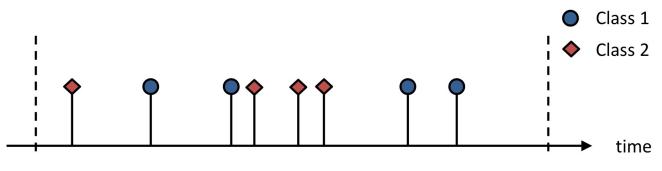


$$p^*(t) = p(t | H(t)) = p(t | h_i)$$

P23 Marked TPP : Categorical marks

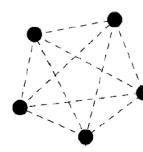
Events has associated class

Events may influence each other



01 - Graphs and Networks

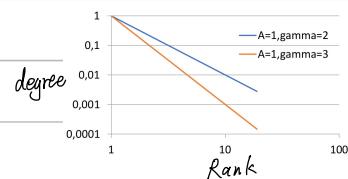
P05 Graphs : $G = (V, E)$ / V : sets of nodes, E : sets of edges $\Rightarrow E \subseteq V \times V$
 Adjacency matrix : $A \in \{0, 1\}^{|V| \times |V|}$



P08 Erdős-Renyi Random Graph Model : Start with N isolated Nodes

Each pair $v_1, v_2 \in V$ add edges with equal probability p

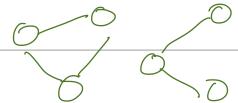
Real world graph is not random No obvious patterns



P10 Power Law Distribution : $y(x) = Ax^{-\gamma}$ / x and γ are variables, A and γ are positive constants
 Scale-free Law : $y(a \cdot x) = b \cdot y(x)$
 → A good descriptions / fit but not the data itself

P15 Patterns and Algorithms: Clustering Coefficient

$$c = \frac{3 \cdot \# \text{ triangles}}{\# \text{ connected triplets}}$$



Real world $\Leftrightarrow c \uparrow \Leftrightarrow$ Strong community // Expensive calculate

Efficient Triangle Counting : # triangles = $\frac{1}{6} \text{trace}(A^3) = \frac{1}{6} \sum_i x_i^3$

↪ Eigenvalues of adjacency matrix A

use power iteration to compute single eigen vector $v \leftarrow \frac{x \cdot v}{\|x \cdot v\|}$

P21 Real Graph : N nodes for N^2 edges

In Real, $E \ll N^2 \Rightarrow$ Sparse

P23 Characteristic Path Length : For starting Node V consider shortest path to each other

Take average length

Take average length for all starting Node and median
 median $\left\{ \frac{1}{|V|} \sum_{v_j \in V} \text{len } p_{\min}(v, v_j) \right\}$

Average Diameter : mean = $\frac{1}{|V|} \sum_{v \in V} \frac{1}{|V|} \sum_{v_j \in V} \text{len } p_{\min}(v, v_j)$

Eccentricity / Hop-plot : # reachable nodes via h hops : $N_h(u) = \{v \in V \mid \text{len } (p_{\min}(u, v)) < h\}$

Total neighborhood size : $N_h = \sum_{u \in V} |N_h(u)|$

Effective diameter : 90%

Q2 - Generative Models

P₆₂ Generative Model : Statistical model describe data distribution

Task: unsupervised learning $p(x)$ / generating data

GMM : prior probability for each cluster $\pi_k > 0 / \sum_k \pi_k = 1$
cluster indicator $z_i \sim \text{Cat}(\pi)$

Sample $x_i \sim N(\mu_{z_i}, \Sigma_{z_i})$

→ For Graph : Don't know the latent factors

a) What is the expected number of cliques of size m ?

$$\binom{n}{m} p^{\binom{m}{2}}$$

P₆₃ Erdős-Renyi : Given $p \in [0, 1] \rightarrow A_{ij} \sim \text{Bernoulli}(p)$ # 6 $\binom{N}{3} p^3$
Probability of vertex has degree k $P_k = \binom{N-1}{k} \cdot p^k \cdot (1-p)^{N-1-k} \approx \frac{ze^{-z}}{k!}$, $z = p(N-1)$
⇒ power-law \leftrightarrow real world is power-law ✗

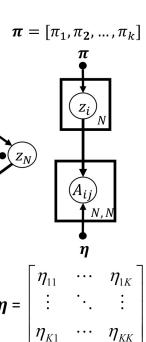
Diameter: $\log(N) / \log(z) \nrightarrow \text{grow slowly} \leftrightarrow \text{real world is constant or shrink} \times$
clustering coefficient: $p = z/(N-1) \nrightarrow \text{Real world look totally different} \times$

→ Generalization of Erdos-Renyi : node from same community connect with p
node from diff community connect with q , $p > q$
 Improvement

P₆₈ Planted Partition Model (PPM) : $V \rightarrow C_1, C_2$, assign node: $z_i \in \{1, -1\}$ ← latent variable
 $P(A_{ij}=1 | z_i, z_j) = \begin{cases} p & \text{if } z_i = z_j \\ q & \text{if } z_i \neq z_j \end{cases}$

disadvantage: two clusters / same p / ... "boring"

P₁₀ Stochastic Block Model (SBM) : advanced PPM 



Random variables $\Rightarrow z_i \in \{1, \dots, K\} / A \in \{0, 1\}^{N \times N}$

Model parameters $\Rightarrow \pi = [\pi_1, \dots, \pi_K]$: community proportions / η_{uv} : edge probability

Condition distributions $\Rightarrow P(z_i=k) = \pi_k / P(A_{ij}|z_i, z_j) = \text{Ber}(\eta_{z_i, z_j})$ in 2 cluster

$$p/q$$

Not capture all patterns ✗

$$B_{ii} = \frac{\# \text{edges within group } i}{\binom{\#\text{nodes in group } i}{2}}$$

$$B_{ij} = \frac{\# \text{edges between groups } i, j}{\#\text{nodes in group } i \cdot \#\text{nodes in group } j} \quad \text{if } i \neq j$$

$$\eta \begin{pmatrix} a & b \\ b & a \end{pmatrix} \# \sigma = \binom{\frac{N}{2}}{3} a^3 \text{ in } \mathbb{R}$$

$$+ x_2$$

$$\# \sigma = \binom{\frac{N}{2}}{2} \binom{\frac{N}{2}}{1} b^2 \cdot a \text{ in } \mathbb{R}$$

$$2 \text{ in } c_1 \quad 1 \text{ in } c_2$$

P₁₅ Preferential Attachment Models : ER / PPM / SBM assume edge generated independently

↳ Generate Network in two steps

↳ Growth : start with a set node, add new nodes and edges

Preferential attachment : "rich get richer", # connecting nodes \propto degree

P₁₇ Initial Attractiveness : start with m_0 nodes \rightarrow add a new node $w \rightarrow$ insert m edges

→ Probability that endpoint of edge (u, v) corresponds to v

is proportional to $A_v = A + \text{indeg}(v)/A$: initial attractiveness for all nodes

→ currently # incoming edges

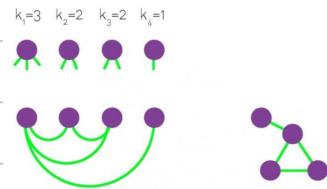
parameters : m (new edge per step) / A (initial attractiveness)

Degree distribution : $\gamma = 2 + \frac{A}{m} \Leftarrow$ power law ✓

Diameter : increase \Leftarrow not same as the real world X

Average degree : constant \Leftarrow increase in real world X

P₁₉ Configuration Model : assign degree to each node at stub
stub \Leftarrow half link \rightarrow random graph



B₂ Deep Generative Model : Difficult to hand-craft a graph fit for all real world data

Solution : VAE / GAN / ... / Net GAN

03 - Ranking

P04 Page Rank : Ranking of Nodes \Leftrightarrow pages not equally important



\hookrightarrow Page is important if many important pages point to it

$$\text{Rank of page } j: r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i} / d_i: \text{out-degree of node } i$$

\checkmark $i = \text{column}$

Gaussian elimination \rightarrow Adjacency Matrix $A \xrightarrow{M \rightarrow \text{sum(column } M) = 1 \Leftrightarrow \sum r_i = 1}$

\hookrightarrow if $i \rightarrow j$, then $M_{ji} = \frac{1}{d_i}$, else = 0

Eigen vector: $r = Mr \Rightarrow I \cdot r = Mr \Rightarrow r$ is eigen vector

Finding r = Finding eigenvector of M with largest eigenvalue

\hookrightarrow power iteration: $r \leftarrow \frac{Mr}{\|Mr\|}$ M is sparse

P10 Random Walk Interpretation: Surfer path \Leftrightarrow Markov chain

\downarrow Start from random pages: $p(X_1=i) = \pi_i$

\downarrow Transition probabilities: $p(X_{t+1}=j | X_t=i) = M_{ji} \Leftrightarrow B = M^T$

$$p(X_t=i) = \pi_i(t) \Leftrightarrow \pi(t) = \pi B^{(t-1)}$$

Stationary distribution $\Leftrightarrow \pi^\infty = \pi^\infty B \xrightarrow{\pi^\infty} r^\top$

limiting distribution: $\lim_{t \rightarrow \infty} \pi(t)$

$\hookrightarrow r = \lim_{t \rightarrow \infty} \pi(t) // \text{rank score of page } i = r_i = \lim_{t \rightarrow \infty} p(X_t=i)$

$$r = \pi B, (\pi B)B, [(\pi B) \cdot B] \cdot B \dots$$

super efficient \Leftrightarrow Vector \cdot Matrix

\hookrightarrow Probability of reaching a node is not depend on start point

\hookrightarrow assume $t \rightarrow \infty$ B^t has same row

$$\lim \pi(t) = \lim \pi B^{(t-1)} = [\pi_1, \pi_2, \pi_3] \begin{bmatrix} a & b & c \\ a & b & c \\ a & b & c \end{bmatrix} = \begin{bmatrix} ka & kb & kc \end{bmatrix}$$

\uparrow
 $k = \pi_1 + \pi_2 + \pi_3$

Above limiting distribution must be Irreducible and Aperiodic

任意状态可以到达任意状态 \rightarrow Irreducible: possible to get any state from any state

从任意状态出发可以回到当前状态 \rightarrow Aperiodic: exists $n \rightarrow$ for all $n' \geq n \Rightarrow p(X_{n'}=i | X_i=i) > 0$

$n' \geq n$ 且能到

\hookrightarrow Markov chain is Aperiodic if all states are Aperiodic

\hookrightarrow Irreducible markov chain: one Aperiodic \rightarrow all Aperiodic

P₁₅ Problems and solutions : Dead ends (not irreducible)
 Spider traps (not irreducible)
 Periodic (not aperiodic)

详细

Random Teleports β : follow a random link

$1-\beta$: jump to random page

$$\Rightarrow r_j = \sum_{i \in j} \beta \frac{r_i}{d_i} + (1-\beta) \frac{1}{N}$$

$$\boxed{A = \beta M + (1-\beta) \left[\frac{1}{N} \right]_{N \times N} \Rightarrow r = Ar} \Rightarrow B = M^T \Rightarrow B = A^T$$

problem: A is dense Matrix

$$\text{Instead: } r = \beta M r + (1-\beta) \left[\frac{1}{N} \right]_N \leftarrow \text{Sparse vector}$$

Page Rank Problems: Measure generic popularity / Link Spam / Single measure of importance

P₂₂ Topic-Sensitive PageRank : based on topic popularity

Core idea: Bias the random walk

Standard PageRank : S = all pages with equal probability

Topic-Sensitive PageRank: S = set of "relevant" pages / diff rs

$$r = \beta M r + (1-\beta) \pi \quad \text{where } \pi_i = \begin{cases} \frac{1}{|S|} & \text{if } i \in S \\ 0 & \text{otherwise} \end{cases}$$

$$r = \beta M r + (1-\beta) \pi \quad \text{where } \sum_i \pi_i = 1$$

↓

$$r = (1-\beta) (I - \beta M)^{-1} \pi$$

* important: when convergence, set $S = \{i\}$, $r_i \geq 1 - \beta$

P₂₆ Variants: Normal $\pi_i = c$

Personalized page

Spam \leftrightarrow Trust Rank

随机走跳

$$\frac{a_0 (1-q)}{1-q}$$

随机走跳

$$c(n,k) \frac{n!}{k! \times (n-k)!}$$

04 — Clustering

P₀₃ Clustering : Given $G = (V, E)$, find clusters of vertices

Goal : vertices in same clusters has more edges

vertices in diff clusters has less edges

Two categories : Partitioning approaches : Each vertex belong to one cluster

Non-Partitioning approaches : Cluster can overlap . Vertex \in None

P₁₀ Partitioning approaches : Constrained optimization problem

Given $G = (V, E)$, assign vertex V to a set of clusters $C = \{C_1, \dots, C_K\}$

such that $Q(C)$ is optimized / s.t. $\begin{cases} C_1 \cup C_2 \cup \dots \cup C_K = V \\ \forall 1 \leq i < j \leq K : C_i \cap C_j = \emptyset \end{cases}$

P₁₂ Global Minimum Cut : Goal: Minimize the number of edges/weights between clusters

$$\text{Cut}(C_1, C_2) = \sum_{v_i \in C_1, v_j \in C_2} w(v_i, v_j) \text{ is minimized}$$

Problem : Tend to cut small vertex sets / consider only inter-cluster edges

P₁₃ Normalized Cut :

Ratio Cut : Minimize $\frac{\text{cut}(C_1, G)}{|C_1|} + \frac{\text{cut}(C_2, G)}{|C_2|}$

Normalized Cut : Minimize $\frac{\text{cut}(C_1, G)}{\text{vol}(C_1)} + \frac{\text{cut}(C_2, G)}{\text{vol}(C_2)}$

volume of a set of nodes
$\text{vol}(C_i) =$
$\text{assoc}(C_i, V) =$
$\text{cut}(C_i, V) =$
$\sum_{v_i \in C_i, v_j \in V} w(v_i, v_j) =$
$\sum_{v_i \in C_i} \deg(v_i) =$

P₁₄ Multi-way Cut : $K \geq 2$

- Cut: $\min_{C_1, \dots, C_K} \sum_{i=1}^k \text{cut}(C_i, V \setminus C_i)$

- Ratio Cut: $\min_{C_1, \dots, C_K} \sum_{i=1}^k \frac{\text{cut}(C_i, V \setminus C_i)}{|C_i|}$

- Normalized Cut: $\min_{C_1, \dots, C_K} \sum_{i=1}^k \frac{\text{cut}(C_i, V \setminus C_i)}{\text{vol}(C_i)}$

P₁₅ Graph Laplacian : Effective Computation

Laplacian matrix : $L = D - W$

D : degree matrix 对角矩阵 , 该节点所有权重

W : weighted Adjacency matrix 未连接的 = 0

$$L_{uv} = \begin{cases} -w_{uv}, & (u, v) \in E \Rightarrow u \neq v, D=0 \\ \deg(v), & u=v \Rightarrow u=v, D \neq 0, W=0 \\ 0, & \text{else} \Rightarrow u \neq v, D=0, (u, v) \notin E, W=0 \end{cases}$$

$$\text{For any vector } f \Rightarrow f^\top \cdot L \cdot f = \frac{1}{2} \sum_{(u, v) \in E} W_{uv} (f_u - f_v)^2$$

Assign "number" to each node : $f: V \rightarrow \mathbb{R}$

Laplacian transform f to another function g : $\delta f = g$

P17 Properties : ① L is symmetric and PSD $\Rightarrow \mathbf{x}^T L \mathbf{x} \geq 0$

$\hookrightarrow D$ and W are symmetric

② Smallest eigenvalue of L is 0, vector is $\mathbf{x}_1 = \mathbf{0}$

$$L\mathbf{x} = \lambda \mathbf{x} \Rightarrow \mathbf{x}^T L \mathbf{x} = \lambda = \frac{1}{2} \sum_{(u,v) \in E} W_{uv} (x_u - x_v)^2$$

Trivial solution set $\mathbf{x} = c \cdot \mathbf{i}$ and $\lambda = 0$

③ L has n non-negative eigenvalues

④ Additive : $L_{GH} = L_G + L_H$

⑤ Graph Laplacian : $L_G = \sum_{(u,v) \in E} L_{(u,v)}$

⑥ # Eigenvector of L with eigenvalue 0 = # connected components = $\# k$

C_k : the set of nodes of the k -th connected component

let $f_{C_k}[i] = 1$ if $v_i \in C_k$, else = 0

$\Rightarrow f_{C_k}^T L f_{C_k} = 0, \forall k \Rightarrow f_{C_k} \Leftarrow$ smallest eigenvector = eigenvalues = 0

$$\left\lceil f^T \cdot L \cdot f = \frac{1}{2} \sum_{(u,v) \in E} W_{uv} (f_u - f_v)^2 = \underline{\lambda} \right\rceil$$

If graph is connected, $\lambda_2 > 0 \Leftarrow (\lambda_0 = \lambda_1 < \lambda_2 < \dots < \lambda_n) \Leftarrow$ ~~not~~ fully graph

$\hookrightarrow \lambda_2(L) \Leftarrow$ Fiedler eigenvalue

\rightarrow reflect how well the connect is

$\lambda_2(k_n) = n$, if k_n is a fully connected graph with n nodes

⑦ Isoperimetric ratio of S : $\theta(S) = \frac{\text{cut}(S, \bar{S})}{|S|} \geq \lambda_2 (1 - \frac{|S|}{|V|})$

Cheeger constant of graph $\theta_G = \min_{|S| \leq \frac{|V|}{2}} \theta(S) \geq \frac{\lambda_2}{2}$

\nwarrow cluster smaller than half graph

λ_2 is big = graph is well connected

⑧ Exist co-spectral graphs, not isomorphic but share same spectrum

P22 Graph Laplacian and Minimum Cut : For $k=2$, $f \in \{-1, 1\}^n$, $f_{C_1}[i] = \begin{cases} +1, & v_i \in C_1 \\ -1, & v_i \in V \setminus C_1 = C_2 = \bar{C} \end{cases}$

$$f_{C_1}^T L f_{C_1} = 4 \cdot \text{cut}(C_1, C_2) / \text{Minimize } f_{C_1}^T L f_{C_1} \rightarrow \text{Minimize the cut}$$

P₂₃ Spectral Clustering : 2 clusters : Ratio Cut for $k=2 \Rightarrow \min_{C_1 \in V} \frac{\text{Cut}(C_1, C_2)}{|C_1|} + \frac{\text{Cut}(C_2, C_1)}{|C_2|}$

Indicator vector f_{C_1} for cluster C_1 / $f_{C_1}[i] = \begin{cases} +\sqrt{\frac{|C_1|}{|V|}} & \text{if } v_i \in C_1 \\ -\sqrt{\frac{|C_1|}{|V|}} & \text{else} \end{cases}$

\downarrow

① $\sum_i f_{C_1}[i] = 0$ / ② $f_{C_1}^T f_{C_1} = \|f_{C_1}\|_2^2 = |V|$ / ③ $f_{C_1}^T L f_{C_1} = |V| \left[\frac{\text{Cut}(C_1, C_2)}{|C_1|} + \frac{\text{Cut}(C_2, C_1)}{|C_2|} \right]$
constant length

Finally : Minimize $\{f_{C_1}^T L f_{C_1}\}$, s.t. $f_{C_1} \perp \vec{1}$ and $\|f_{C_1}\|_2 = \sqrt{|V|}$

f_{C_1} is the second smallest eigenvector of L
 $\Rightarrow L \cdot f_{C_1} = f_{C_1} \cdot \lambda_2 \Leftrightarrow f_{C_1}^T L f_{C_1} = |V| \cdot \lambda_2$ Relax iso constraint

P₂₇ General case : $k > 2$ / $h_C[i] = \begin{cases} \frac{1}{\sqrt{|C_i|}} & \text{if } v_i \in C_i \\ 0 & \text{else} \end{cases} \Rightarrow H = [h_{C_1}; h_{C_2}; \dots; h_{C_k}]$

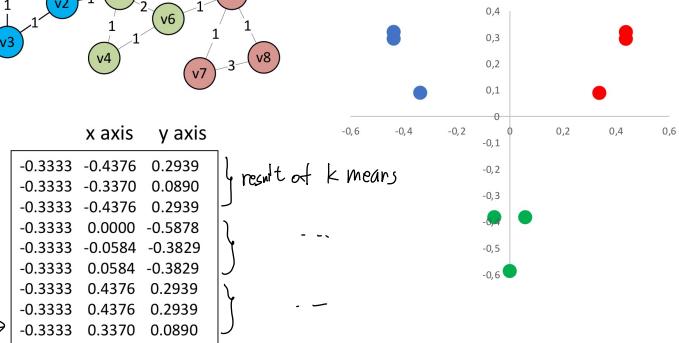
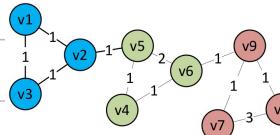
$\Rightarrow H^T H = Id$ / $h_{C_i}^T L h_{C_i} = \frac{\text{Cut}(C_i, V \setminus C_i)}{|C_i|}$ and $h_{C_i}^T L h_{C_i} = (H^T L H)_{ii}$

Ratio cut = $\sum_{i=1}^k \frac{\text{Cut}(C_i, V \setminus C_i)}{|C_i|} = \sum_{i=1}^k (H^T L H)_{ii} = \text{Tr}(H^T L H)$

Minimize ratio cut = minimize $\text{Tr}(H^T L H)$ s.t. $H^T H = Id$

\Rightarrow optimal H = first k smallest eigenvectors of L

Spectral Layout:



P₃₇ Inference in PPM : Likelihood of community assignment $\mathbf{z} \in \{-1, 1\}^N$ / observed symmetric $A \in \{0, 1\}^{M \times N}$

$$\Pr(A|\mathbf{z}) = \prod_{i < j} [p^{A_{ij}}(1-p)^{1-A_{ij}}]^{\mathbb{I}(z_i=z_j)} [q^{A_{ij}}(1-q)^{1-A_{ij}}]^{\mathbb{I}(z_i \neq z_j)}$$

$$\text{MLE} : \mathbf{z}^* = \underset{\mathbf{z} \in \{-1, 1\}^N}{\operatorname{argmax}} \Pr(A|\mathbf{z})$$

P38 PPM and Spectral Clustering : assume balanced cluster $|C_1| = |C_2| = \frac{N}{2} \Leftrightarrow \sum_i z_i = 0$

$$E_{cut}(\mathbf{z}) = |\{(i, j) \in E \text{ s.t. } z_i \neq z_j\}| = \sum_{(i, j) \in E} \mathbb{I}(z_i \neq z_j)$$

$$\Pr(\mathbf{A}|\mathbf{z}) \propto \left(\frac{(1-p)q}{(1-q)p} \right)^{E_{cut}(\mathbf{z})}$$

$\rightarrow p > q \Rightarrow$ Find a minimum balanced cut
 $p < q \Rightarrow p > 1 \Rightarrow$ Find maximum cut

P41 Inference in SBM : Assume η and π are known $\Pr(\mathbf{z}|\mathbf{A}, \boldsymbol{\eta}, \boldsymbol{\pi}) = \frac{\Pr(\mathbf{A}|\mathbf{z}, \boldsymbol{\eta}, \boldsymbol{\pi}) \Pr(\mathbf{z}|\boldsymbol{\pi})}{\Pr(\mathbf{A}|\boldsymbol{\eta}, \boldsymbol{\pi})}$



use variational inference to find approximate solution

$$\Pr(\mathbf{z}|\mathbf{A}, \boldsymbol{\eta}, \boldsymbol{\pi}) \approx q(\mathbf{z}|\boldsymbol{\Psi}) = \prod_{i=1}^N q(\mathbf{z}_i|\boldsymbol{\psi}_i)$$

Learning in SBM : If both \mathbf{z} and \mathbf{A} are observed, then MLE for η and π

$$\pi_k^{MLE} = \frac{\# \text{ nodes in cluster } k}{N} =: \frac{N_k}{N}$$

$$\eta_{uv}^{MLE} = \frac{\text{observed # edges between } u \text{ and } v}{\text{possible #edges between } u \text{ and } v} = \frac{\sum_{(i,j) \in E} \mathbb{I}(z_i = u) \mathbb{I}(z_j = v)}{P_{uv}}$$

where $P_{uv} = \begin{cases} \binom{N_u}{2} & \text{if } u = v \\ N_u \cdot N_v & \text{if } u \neq v \end{cases}$ is the number of possible edges between clusters u and v

If only \mathbf{A} are observed, then use variation inference

05 - Classification

P03 Introduction : Supervised learning problems

P04 collective classification : known \rightarrow labeled nodes / classify \rightarrow unlabeled nodes
 ↳ based on structure / edges / relationship

P05 Label Propagation : Consider binary case

- $V = S \cup U$ / S : Labeled instances S (seeds) ; U : Unlabeled instances
- $W \in \mathbb{R}^{V \times V}$: Symmetric weighted adjacency matrix / $w_{ij} \geq 0$ = similarity
- $\hat{y}_i \in \{0, 1\}$ for $i \in S$: Given class label
- $y_i \in \{0, 1\}$ for $i \in U$: Unknown class label

Idea: Energy Minimization : $E(y) = \frac{1}{2} \sum_{ij} w_{ij} (y_i - y_j)^2 \Leftrightarrow \min E(y) = \text{min cut}$

- ↓
- smoothness : adjacent nodes should have same class label
 - matching the seed labels

$$\min_{y \in \{0, 1\}^V} \frac{1}{2} \sum_{ij} w_{ij} (y_i - y_j)^2, \text{ s.t. } y_i = \hat{y}_i \text{ for all } i \in S$$

$$\min_{y \in \{0, 1\}^V} y^T L y, \text{ s.t. } y_i = \hat{y}_i \text{ for all } i \in S$$

$\uparrow R^V$ $\uparrow L = D - W$

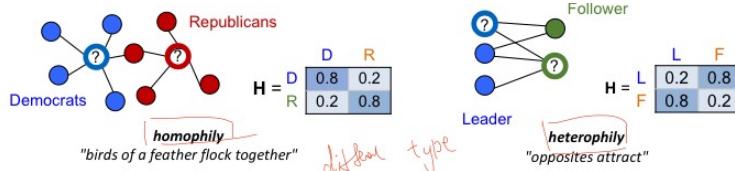
solution : assume $L = \begin{bmatrix} L_{ss} & L_{su} \\ L_{us} & L_{uu} \end{bmatrix} \quad \min y^T L y$

$$y = \begin{bmatrix} y_s \\ y_u \end{bmatrix} = \begin{bmatrix} \hat{y}_s \\ y_u \end{bmatrix} \quad \Rightarrow y_u = -L_{uu}^{-1} L_{us} \hat{y}_s$$

k labels : one-hot $y_{ik} = \begin{cases} 1 & \text{if node } i \text{ is of class } k \\ 0 & \text{else} \end{cases}$

$$E(Y) = \sum_{i,j} w_{ij} (y_i - y_j)^T (y_i - y_j)$$

Network effect H :



estimate pars.
 ↓

SBM vs LP \Rightarrow predict label

$$E(Y) = \sum_{i,j} w_{ij} (y_i - y_j)^T H (y_i - y_j)$$

P₁₀ Transductive Learning \Rightarrow Label Propagation (LP)

Given: label training instances T / unlabeled test instances V
other prior knowledge

Goal: predict label only for unlabeled instances V (no new data!)

P₁₂ Semi-supervised learning (SSL) \Rightarrow Given T, V , Goal: learn a mapping f
inductive ($f: X \rightarrow Y$) / transductive ($f: V \rightarrow Y$)

Unlabeled data help to model the data distribution

Motivation: Large unlabeled data + small labeled data
label is expensive

P₁₄ Attributed Graph: Each node with a feature vector
 \Rightarrow Markov random fields

P₁₅ Graph construction: construct a graph connecting similar data points
apply LP

P₁₆ Graph classification: $G = (V, E) \Rightarrow$ classify V

Given: a set of training graphs $\{G_i = (V_i, E_i)\}_{i=1, \dots, n}$
label y_i for each G_i

Goal: Learn a function $f: g \rightarrow Y$, map new graph with label

Method: Graph kernel function: $k: g \times g \rightarrow \mathbb{R}$ / $k(G_1, G_2) = \phi(G_1)^T \phi(G_2)$
SVM / GNN ...

06 - Node Embeddings

P₀₃ Motivation: Want to apply traditional ML methods to graph
→ have no concept of graphs

P₀₄ Learning Node Representation: Transform graph \Rightarrow vertex is represented by vector
Node embedding function: $\phi: V \rightarrow \mathbb{R}^d$
close nodes are similar in embedding space

P₀₅ Application: Clustering / Semi-supervised classification / Link prediction

P₀₆ Type of Embedding:
Spectral Embedding \rightarrow nodes are similar, if belong to cluster
Deep Embedding \rightarrow depend on loss function and architecture
first k eigenvector of L

P₀₇ Spectral Embedding: based on the eigenvectors of the graph Laplacian L
 $|V| \times |V|$ Adjacent matrix $\rightarrow |V| \times k$ matrix H

P₁₂ Deep Embedding: Capture more complex structure

Deep Walk: Graph \Rightarrow Random Walk / learn word2vec model

Every node sample multiple random walks

Train Word2vec

close node has similar vector representation

attributed features

Graph2Gauss: Learning f_θ to map node v_i to Gaussian distrib. $/ v_i \rightarrow N(\mu_\theta(x_i), \Sigma_\theta(x_i))$

1st node closer / 2nd node far away

Loss Function: $\sum_{(u,v)} E_{uv}^2 + e^{-E_{uv}} / E_{uv} = KL(f_\theta(u) || f_\theta(v))$

P₁₅ Graph Embedding: Aggregate all node embedding into one graph embedding
→ mean pooling

07 - GNN

P₀₃ Motivation : Label propagation consider only graph structure
 Real-world graph have attributes
 GNN: ANN operates on graph-structured data

Adjacent matrix is not image fed into CNN

multiple GCN

Two versions: { Spatial GNNs
 Spectral GNNs

P₀₇ Spatial GNNs: Differentiable Message Passing Network Framework

Message passing = local aggregation

Given: $G = \{V, E\}$ and node $v = x_v \in \mathbb{R}^d \Rightarrow G = (A, x)$

Even with edge features E_{vu}

$h_v^{(k-1)}$: hidden representation of node v at previous $k-1$ layer

Solve: Message gather: $m_v^{(t)} = \sum_{u \in N(v)} M(h_v^{(k-1)}, h_u^{(k-1)}, E_{vu}) = \sum \frac{1}{d_v} (w^{(k)} h_u^{(k-1)} + b^{(k)})$

Hidden representation update: $h_v^{(k)} = V(h_v^{(k-1)}, m_v^{(k)}) = \text{relu}(Q^{(k)} h_v^{(k-1)} + p^{(k)} + m_v^{(k)})$

V and M are trainable NN



P₁₁ GCN: $X_i^{(t+1)} = X_i^{(t)} + \sum_{j \in N(i)} X_j^{(t)} / X^{(t+1)} = X^{(t)} + A X^{(t)}$

$H^{(L)} = [\quad]$ encode the hidden representations after layer L

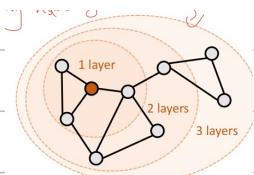
$$H^{(L+1)} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(L)} W^{(L)})$$

Message passing
 Feature Transformation
 non linearity

$\Rightarrow \begin{cases} \tilde{A} = A + I_n \\ \tilde{D}_{ij} = \sum_j \tilde{A}_{ij}, \text{ degree matrix} \end{cases}$

At graph step k , node v aggregates information from k -hop neighbours
 ↳ all nodes in its k -hop

Transformation depth and propagation depth are orthogonal
 ⇒ more layers / more node information



P15 Graph Signal Processing:

Given: Normalized Laplacian matrix $\hat{L} = I - D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$

Eigen decomposition: $\hat{L} = U \Lambda U^\top$, $0 = \lambda_1 \leq \lambda_2 \dots \leq \lambda_n \leq 2$

$$U^\top U = I$$

Graph Signal: A function of nodes $x: V \rightarrow \mathbb{R}$, $x \in \mathbb{R}^N$

Graph Fourier transformer: $\tilde{x} = U^\top x / x = U^\top \tilde{x}$

Spectral filtering: $y = U g(\Lambda) U^\top x / g: \text{spectral filter}$

P16 Spectral Graph NN: $H^{(l+1)} = \phi(U g(\Lambda) U^\top \varphi(H^{(l)}))$ $g_\theta(\Lambda) = \sum_{k=0}^K \theta_k \Lambda^k \in \text{polynomial filter}$

$$= \phi(g_\theta(\hat{L}^k) \varphi(H^{(l)}))$$

Computational complexity: avoid eigendecomposition

Receptive field size: whole graph not the neighbour

Transferability:

Application of GNN

P19 Semi-Supervised Node Classification: probability of node v belong to class c

$$\Rightarrow p_v = \text{softmax}(h_v^{(k)})$$

$$\Rightarrow \text{use CE} \quad \min - \sum_{v \in \mathcal{V}} \sum_{c \in \mathcal{C}} y_{vc} \log p_{vc}$$

P20

Graph Classification: $H_{G_i} = [h_1^{(k)}, \dots, h_{|V_i|}^{(k)}]$ of $G_i = \{V_i, E_i\}$

$$\Rightarrow \text{Aggregation function} \quad h_{G_i} = R(H_{G_i}) = \frac{1}{|V_i|} \sum_{v \in V_i} h_v^{(k)}$$

$$\Rightarrow \text{use CE} \quad \min L(y_{G_i}, \text{softmax}(h_{G_i}))$$

P21 Equivariant Machine Learning: Given $G_1 = (A_1, x_1)$, $G_2 = (A_2, x_2)$

isomorphic \Rightarrow exist a permutation matrix P , let $PA_1 P^\top = A_2$, $Px_1 = x_2$

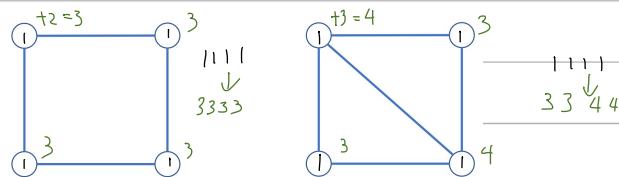
Graph isomorphism \Rightarrow invariance $\Rightarrow f(Px, PA P^\top) = f(x, A)$

equivariance $\Rightarrow f(Px, PA P^\top) = Pf(x, A)$

$E(n)$ equivariant: $E(n) = \{Rx + T : R \in O(n)\}$

P₂₄ Limitations of GNN: Expressiveness / Reliability (Non-Robust) / Scalability

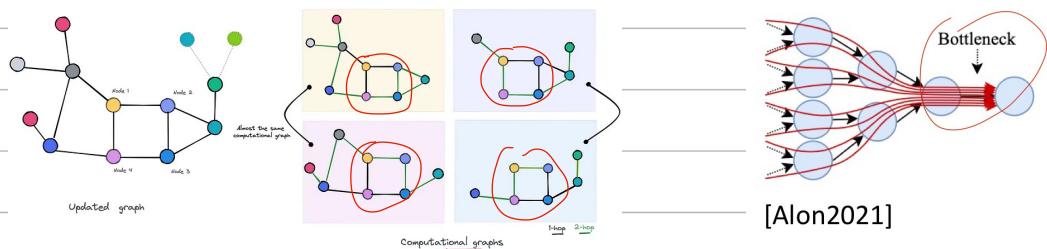
P₂₆ Weisfeiler-Lehman Test:



\Rightarrow can distinguish many non-isomorphic graph but not all of them

P₂₈ Oversmoothing: Infinite layers \Rightarrow stationary distribution

FCN Graph



P₃₀ Overquashing: long distance \Rightarrow in one node \Rightarrow bottleneck
tree-based Graph

$$\begin{aligned} P_{31} \text{ Page Rank : } & \text{Page Rank : } r^{(t+1)} = \boxed{AD^{-1}} \boxed{r^{(t)}} \\ G_{CN} : & H^{(t+1)} = \boxed{\tilde{D}^{-\frac{1}{2}} \tilde{A} D^{-\frac{1}{2}}} \boxed{G(H^{(t)}) W^{(t)}} \\ & \text{Message Passing} \quad \text{Features} \end{aligned}$$

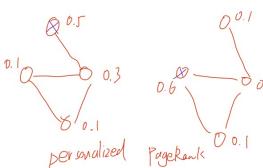
P₃₃ Personalized Propagation of Neural Prediction (PPNP) : Personalized teleport $\pi = e$

\hookrightarrow Separate Transformation and Propagation

$$① H^{(0)} = X \Rightarrow H_i^{(0)} = f_\theta(x_{i,:})$$

$$② r^{(t+1)} = (1-\alpha) A D^{-1} r^{(t)} + \alpha \pi \quad / \quad H^{(t+1)} = (1-\alpha) \tilde{D}^{-\frac{1}{2}} \tilde{A} D^{-\frac{1}{2}} H^{(t)} + \alpha H^{(0)}$$

$$③ \text{PPNP for Node Classification} \Rightarrow Z = \text{softmax} \left[\alpha (I_N - (1-\alpha) \tilde{A} \tilde{D}^{-\frac{1}{2}})^{-1} H^{(0)} \right]$$



prevent oversmoothing / Separation of Transformation and Propagation / Scalability

$$P_{37} \text{ APPNP : } Z = \text{softmax} \left[\alpha (I_N - (1-\alpha) \tilde{A} \tilde{D}^{-\frac{1}{2}})^{-1} H^{(0)} \right] \Rightarrow Z = H^{(k)}$$

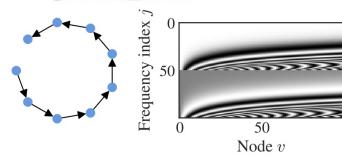
P₃₈ Graph Transformers \Rightarrow Positional Encodings

P₃₉ Sequence problem : Sinusoidal Encodings $PE \in \mathbb{R}^{n \times d}$

positional encodings \Leftarrow Discrete Fourier Transformation

$$\hat{x}_j = \sum_{v=0}^{n-1} x_v \left[\cos\left(v \cdot \frac{2\pi}{n} j\right) - i \cdot \sin\left(v \cdot \frac{2\pi}{n} j\right) \right]$$

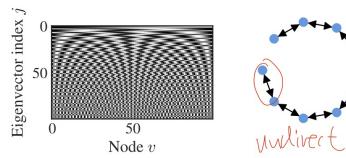
▪ Sinusoidal Encodings



$$PE_{v,j+d/2}^{(\sin)} = \sin(v \cdot 1/10,000^{2j/d})$$

$$PE_{v,j}^{(\sin)} = \cos(v \cdot 1/10,000^{2j/d})$$

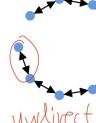
\Leftrightarrow Eigenvectors of Laplacian



Closed form for $L = D - A = U \Lambda U^T$ of a sequence:

$$U = PE_{v,j}^{(\text{lap})} = \boxed{\pm \cos\left(v \cdot \frac{j\pi}{n} + \frac{j\pi}{2n}\right)}$$

Eigenvalues



P₄₁ Direct Graph : $L = D - W / \hat{L} = I - D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$ are real symmetric
if A is symmetric and graph is undirected

\Rightarrow Direction-aware Magnetic Laplacian $L^{(q)} \Leftrightarrow$ encode edge directions

$$L^{(q)} := D_S - A_S \odot \exp(i\Theta^{(q)}) \quad \hat{L}^{(q)} := I - D_S^{-\frac{1}{2}} A_S D_S^{-\frac{1}{2}} \odot \exp(i\Theta^{(q)})$$

$$\Theta_{u,v}^{(q)} := 2\pi q (A_{u,v} - A_{v,u}) \quad \text{direction}$$

