

remove markup

detected & deletion of duplicate symbols

fit the parameters

task, data & algorithm

pre-processing

Train

1. data acquisition & preparation

Development

alter experimental condition

2. feature extraction

split 50:50

Validation

parameters tune

3. evaluation

evaluation

Tuning

hyperparameters : 2/ control process

different metrics & approaches

4. use

test holdout

unseen; once instead of looking

✓ Unbiased.

Supervised

Unsupervised

Semi-Supervised

Manually annotated (expert)

moderately

→ from (star review)

gradient descent
stochastic gradient descent

rm bias

algo robustness to unexpected data

consistency across datasets

No.

Task 1

Date

Jan 20

preannotated → sparsity
insufficient data

Review 1 ✓

~~x human labour~~
~~x what to include~~
~~x stance language changes~~

(specialized) fixed, predefined
standing going &

Task 2 y317

best.

well

satisfying

+

+

+

sentiment strength

Review 2 ✗

but

bland

lacking

not

like ✓

—

—

—

—

Review 3 ✗

no isn't ironic

Review 4 ✓

escapist

fun

capitalise

✓ tokeniser ~ split

preprocessing: normalization

tick.py

tokenize the review texts

utils / tokenizer.py.

definite: O

Classification.

records

positive / negative orientation (writer)

review - data.

lexicon Dict[str, int]

Evaluation

Output

Ground Truth

Accuracy

Raise accuracy

predict sentiment

List[str]

Dict[str, int]

punctuation, word tokenization (next ch)

Language modelling

LM

(next word)

Spam detection (spam/not-spam)

No. Tick'2

Text categorization / (sentiment analysis) Language, authorship, lib. subject Date category

Naive Bayes Classification: assign a category to each input

predefined classes.

sentiment
(labels)
 $C_1, C_2 \dots$

features: w_1, w_2, \dots, w_n

Reviews

features: O



classes:

$C_1, C_2 \dots$

labels

1. $\{w_1, w_2, \dots, w_n\}^D$: unordered set; position x $\{Positive; Negative\}$

2. features are independent (observable properties)

Bayes Theorem:

Multinomial naive Bayes classifier

$$C_{NB} = \underset{c \in C}{\operatorname{argmax}} p(c|O) = \frac{p(c) \cdot p(O|c)}{p(O)}$$

likelihood joint PDF

constant

part of $p(c|d)$

WLOG, represent document as a set of features

$$= p(w_1) p(w_2|w_1) \cdots p(w_n|w_1, w_2, \dots, w_{n-1}) = p(w_1, \dots, w_n | c)$$

marginal PDF

conditional PDF

for c

✓ training: collect $p(c)$

$p(w_i|c)$

✓ testing: $\rightarrow C_{NB}$ addition

mutually

MLB

$p(w_i|c) = \frac{\#(w_i|c)}{\# \sum_{w_j} (w_j|c)}$

No

New

assumption

Independence

$p(c) = \frac{\#(c)}{\# \sum_{w_j} (w_j|c)}$

□ avoid underflow & increase speed (+ than \times) count

$$2 \underset{c \in C}{\operatorname{argmax}} \log p(c) + \sum_{i=1}^n \log p(w_i|c)$$

binar classifier

concatenate

feature
separately

✓ data sparseness

sparsity

Add-one

Laplace

Smoothing

estimating parameters of statistical model

category 'c' test

$(w_i|c) + 1$

$(w_i|c) + 1$

If (w_i, c) doesn't appear before including zero counts MLE

$(w_i|c) + 1$

$(w_i|c) + 1$

✓ 0.1 for not related dummy log X

$(w_i|c) + 1$

$(w_i|c) + 1$

dict. [int, Dict [str, float]]

$(w_i|c) + 1$

$(w_i|c) + 1$

The unseen words in training set are dropped completely $\rightarrow p(w_i|c_1) = 1$

$(w_i|c) + 1$

$(w_i|c) + 1$

$p(w_i|c) = 1 - 0 = 1$

$(w_i|c) + 1$

$(w_i|c) + 1$

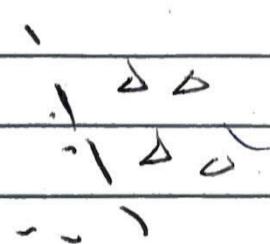
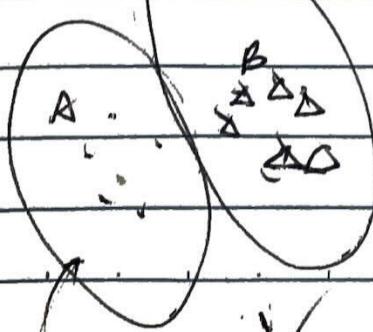
Merge dictionary

$(w_i|c) + 1$

$(w_i|c) + 1$

Generative vs Discriminative

+ weight.



No. Topic 3

Date

Jan. 2

type ✓ ✓ ✓
taken the cat likes the dog ✓

(could remove) stop word

type: any unique word ~ dictionary

token: an instance of a type.

Zipf's Law:

$$fw = \frac{k}{rn} \sim \text{frequency rank}$$

n^{th} frequency of word $\propto \frac{1}{n}$

Heap's Law:

$$\text{type vocabulary size } V_n = kn^\beta$$

1. axis

2. dominant return

tokens
text

unique words vs # total

$$\log V_n \sim 2^0 \sim 2^n$$

(1). $|V| \uparrow$. (2) large

* unseen words \rightarrow smoothing
Estimate V .

$$fw = \frac{k}{(rn+\beta)^{\alpha}} \sim \text{shift}$$

2. $k \sim \text{language}$

frequency "a" "the"

sample a word
smoothing

include punctuation

rank

x-axis:

r_w

draw first 10,000 higher ranked words
chart - plot (data, title, xl, yl)

log - log \rightarrow linear

$$\log fw = -\alpha \log rn + \log k$$

linear least-square

$$y = mx + b$$

$$m = \frac{N \sum xy - \sum x \sum y}{N \sum (x^2) - (\sum x)^2}$$

Covariance

$$\text{Var}(x)$$

$$b = \frac{\sum y - m \sum x}{N} = \bar{y} - m \bar{x}$$

$$y = mx + b$$

didn't like this movie

NOT-like NOT-this NOT-movie
punctuation mark

$V \uparrow$ ~~variance~~ mean less

exists or not instead of frequency

word count \rightarrow 1

remove duplicate in a file

negation

didn't like this movie

NOT-like NOT-this NOT-movie
punctuation mark

Assumption: Individual tests on data are NOT linked
documents ~ isolation
ignores Ground Truth and Information

No.

Tick ✓

effects by chance. (natural variation) Date

Statistical Significance Testing (+ balanced by -) differences between values

BST

Null hypothesis: H_0 come from the same distribution baseline sceptical

1°: significance level 2. $\alpha = 0.01$ or 0.05.

Reject the null hypothesis with confidence 1- α

There is less than 5% we fooled ourselves, null hypothesis is true 0.99

0.95.

20 times on different sets of items random var X over all test datasets. p-value false (as good as)

Sign Test.

count

Type I error: declares a difference when it doesn't exist

Pos: $S_1 > S_2$ | +

test

Neg: $S_2 > S_1$ | -

ties

Tie: $S_1 \sim S_2$ | null

Binomial (N, q)

A negative outcome q

$$P_q(X=k|N) = \binom{N}{k} q^k (1-q)^{N-k}$$

$$P_q(X \leq k|N) = \sum_{i=0}^k \binom{N}{i} q^i (1-q)^{N-i}$$

symmetric

Type II error: declares no difference: β

power: $1-\beta$

use more data

change system (stronger)

powerful test.

permutation test

0 1 2 3 4 5 6 7 8 9 10

Intended

$B(N, 0.5)$

Two Tailed test

$2P(X \leq k) \leq 0.05$

Pos: Neg harder

top little down

→ normal approximation to binomial

One-Tailed Test <

$N!$

$i!(N-i)!$

P_{Neg}

$$P(X \leq k|N) = \sum_{i=0}^k \binom{N}{i} q^i (1-q)^{N-i}$$

$$K = \text{less sig [null]}^{0.5}$$

+ test statistic $\geq \min(+/-)$ len()

- evenly + & balanced

0 less. $n = 2 \lceil \frac{n_{\text{null}}}{2} \rceil + \text{plus + minus}$

common sign one more n odd

representative for the population

Magnitude classifiers.

Times 2 0.65

4 0.663

8 0.665

10 0.665

No. Tick 5.

Date Feb. 3

Generalise

- Only those characteristics of data general enough.

Ignore specific training data. → test data (not training set)

Type III Errors
Overtraining
Testing h based on data

repeated on test data
Same small → properties of test data

"Wayne Rooney" time effects (change meaning)

Confusion Matrix

		Sys	TOTAL	
		Truth	pos	NB
Truth	pos	✓	W	90%
	NB	✗	Y	90%
Total		V+X	W+Y	180%

Ug

used in estimating it

"Leave-one-out" LOOW

Cross Validation

Random

Stratified

test the model's ability to predict new data that was not

in the training set

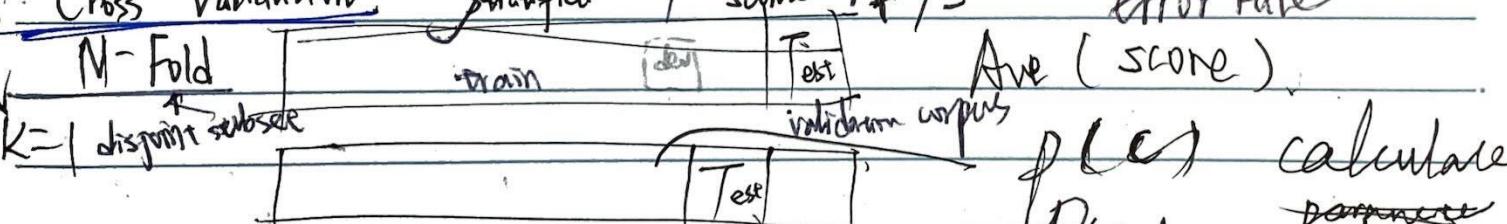
percentage

= same +/-

error rate

N-Fold

K=1 disjoint subsets



Ave (score)

Significance Testing

mean accuracy

Folds together. ⇒ one

mini-events

$$\text{Equally well for each fold. var} = \frac{1}{n} \sum_i (x_i - \mu)^2$$

p (es) calculate
p (N/C) probabilities

predict ← validation

Error bar

→ fewer significant differences

Validation Corpus

TxTx → parameters

new review set.

2004 → 2016

performance

p-value = 1

Human agreement

Positive, Negative, Neutral

Surjective

Luke-warm (words)
Pro-con

(threat)

Star rating: Inter-personal differences.

Surjective Component

Reader's / writer's perception

Kappa

Agreement metrics

$$\bar{P}_a = \frac{1}{N} \sum_{i=0}^{N-1} \frac{\text{observed rater}}{\text{possible rater}}$$

Above

By chance:

$$\bar{P}_e = p^2(\text{POS}) + p^2(\text{NEG})$$

✓ Majority Class

$$\frac{C_3^2 + C_4^2}{C_5^2}$$

✓ majority vote / min(,)

$$\kappa = \frac{\bar{P}_a - \bar{P}_e}{1 - \bar{P}_e}$$

← achieved reliability of agreement

← attainable

✓ D. 8 ✓

○ × small sample size

number of categories

$$\bar{P}_e = \sum_{j=1}^n \left(\frac{1}{Nk} \sum_{i=1}^N n_{ij} \right)^2$$

P for each category

N number of documents/reviews

K number of annotators

n_{ij} item i, class j

$$\bar{P}_a = \frac{1}{N} \sum_{i=1}^N \frac{1}{K(K-1)} \sum_{j=1}^n n_{ij} (n_j - 1)$$

annotator k

item	category
1	n _{ij}
2	n _{ij}
3	n _{ij}

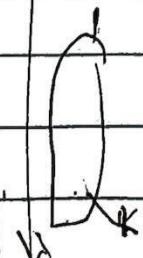
REVIEW

2 3 p

j / N

POS

NEG



AB

Initial State π_i start in certain state. $\sum \pi_i = 1$

No. 7

Date Feb. 9

more than
second order / twice in a row

special end state

Hidden Markov Models

Weather (Rainy & Cloudy)

Markovianity Q

special start state

ϵ_{B1g}

q_0 $2N$ Z

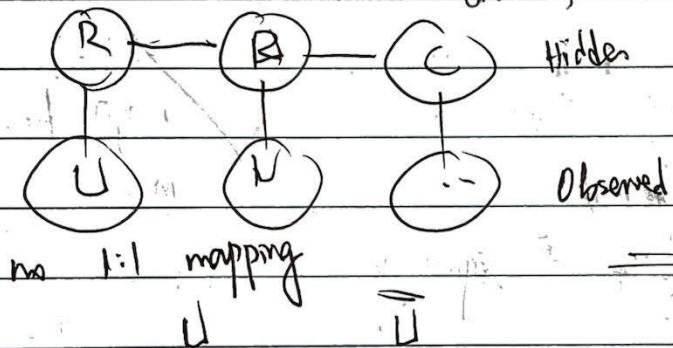
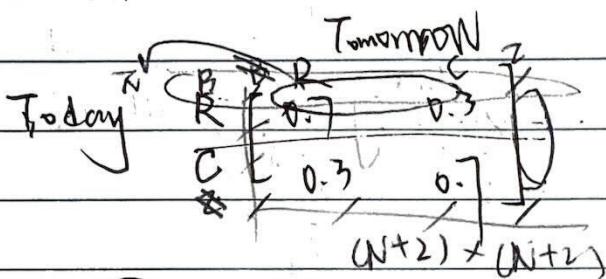
g_f

→ Markov Assumption (first order):

$$P(W_t | W_{t-1}, \dots, W_1) = P(W_t | W_{t-1})$$

JOINT $P(W_1, \dots, W_n) \propto \prod_{t=1}^n P(W_t | W_{t-1})$

Markov Chains



$$\begin{matrix} R & 0.9 & 0.1 \\ C & 0.2 & 0.8 \end{matrix}$$

N emitting hidden states $S: \{S_0, S_1, \dots, S_N\}$, S_N {spelled end}

hidden states $X = B X_1 \dots X_T Z$
observations $O = O_1 \dots O_T$
hidden states O end

start M output alphabet.

$$K: \{k_1, \dots, k_m\}$$

K alphabet. Vocabulary

k_f

observations
independent of

Output Independence: only current state, not other states or any other observation

$$P(O_t | X_1, \dots, X_T, O_1, \dots, O_T) = P(O_t | X_t)$$

sequence

label

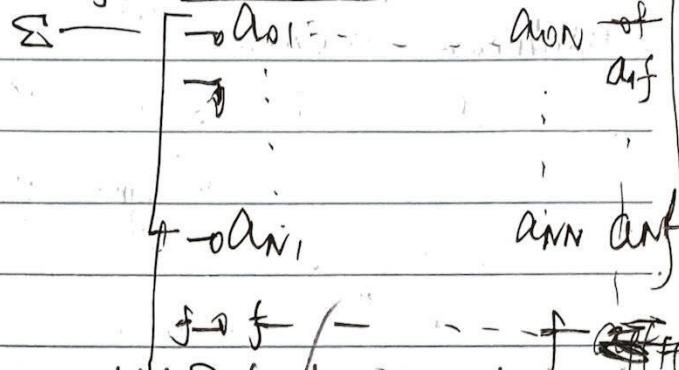
H

only option state: unlabeled

State Transition Probability

$$\text{State } i \rightarrow j \quad a_{ij} = P(X_t = s_j | X_{t-1} = s_i)$$

$$\forall i \sum_{j=0}^{N+1} a_{ij} = 1$$



Abi jaif \vee ; Abi x ; abi X

B

Emission Probability

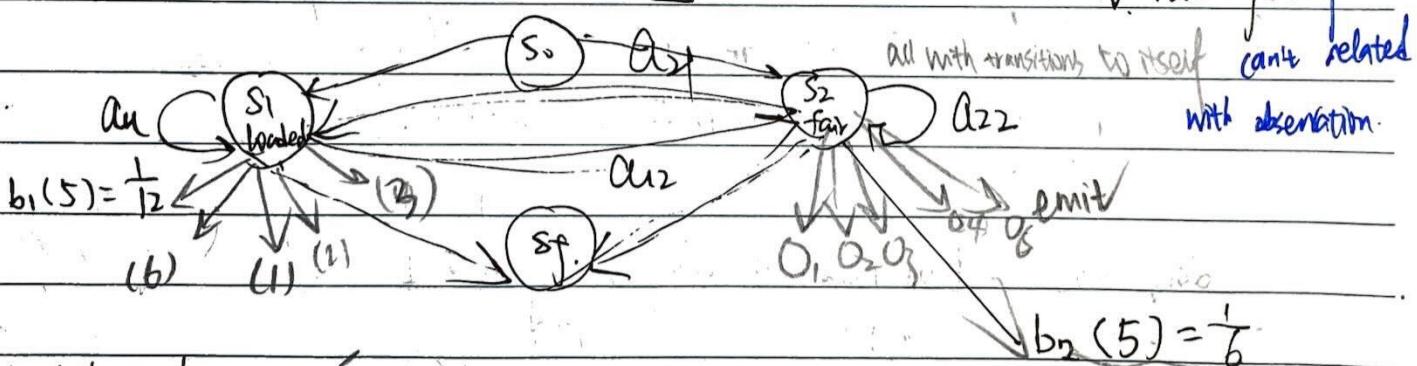
$$B = \begin{bmatrix} 1 & 2 & \times & 5 & 6 \\ b_0(k_0) & - & - & - & - \\ F & - & b_1(k_1) & \dots & - & b_N(k_N) \\ L & - & - & \vdots & - & - \\ \text{state} & - & b_1(k_M) & \dots & - & b_N(k_M) \end{bmatrix} \quad (N+2) \times (M+2)$$

B is an output alphabet of $k = \{k_1, \dots, k_M\}$
 k_0 start symbol, k_f special end.

$$b_i(k_j) = P(O_t = k_j | X_t = s_i)$$

v. item from ↑

all with transitions to itself can't relate with observation



Labelled learning

Dataset:

observation: J n b

$$1. P(F \rightarrow L) = \frac{\#(F \rightarrow L)}{\#(F \rightarrow K)}$$

hidden states: F W

$$2. \frac{\#(O_t \mid s_i)}{\#S_i} = z_i^p$$

ns B Z 4x4

$B \rightarrow Z : Z \rightarrow B$

$$\begin{array}{c} Z \rightarrow \\ \rightarrow B \end{array} \quad \begin{array}{c} 4 \\ 4 \end{array} \quad \begin{array}{c} -2 \\ -2 \end{array}$$

Viterbi Algorithm. hidden state sequence X explains observation O

Decoder

$$\hat{x} = \underset{X}{\operatorname{argmax}} P(X, O | \mu)$$

Conditional

HMM parameter: $\mu = (A, B)$

Basis & \wedge drop denominator

Likelihood

$$= \underset{X}{\operatorname{argmax}} P(O | X, B) P(X | A)$$

$$= \underset{X_1 \dots X_T}{\operatorname{argmax}} \prod_{t=1}^T P(D_t | X_t, B) P(X_t | X_{t-1}, A)$$

$O(N^T)$

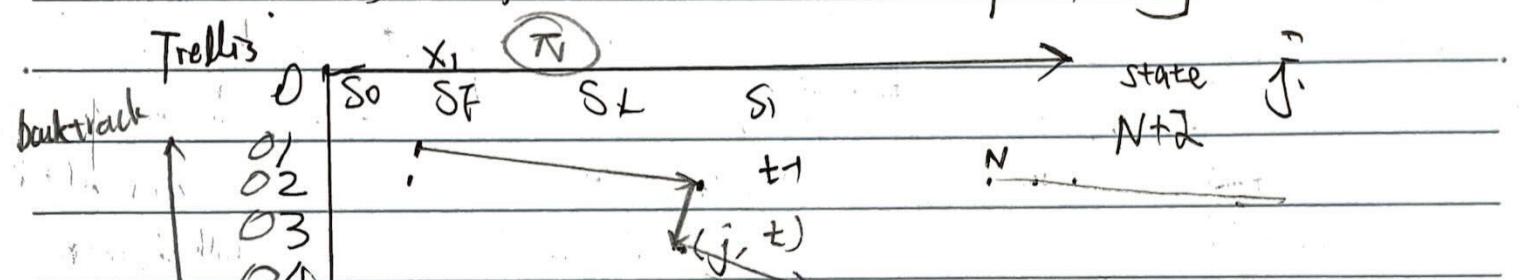
DP: optimal substructure ✓ Overlapping subproblems ✓

$O(N^2 T)$

first order

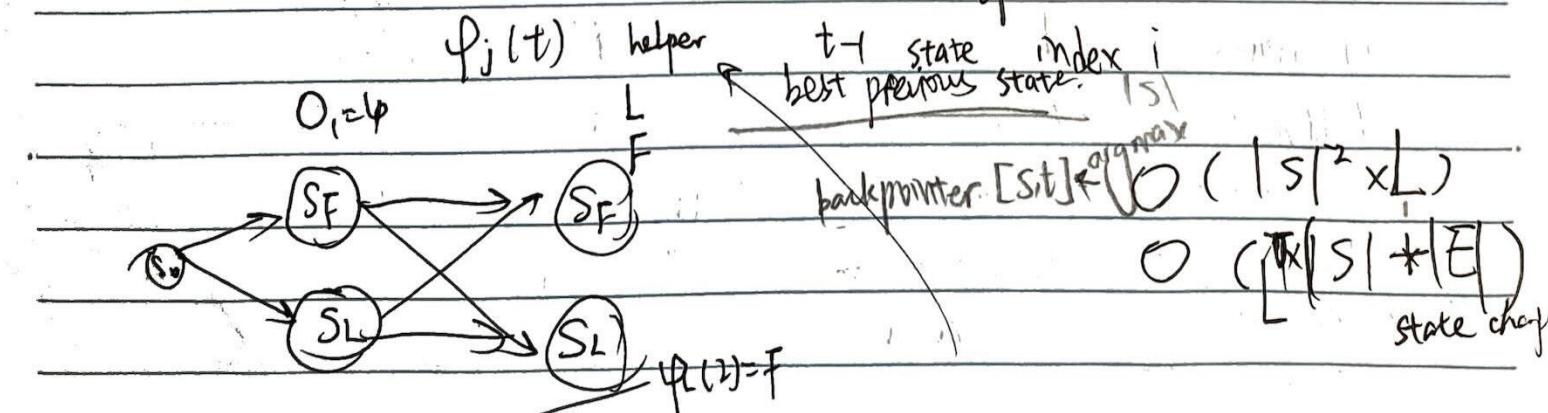
memorize the pr / only cal once

Forward Algo: Greedy



Best route path ending in S_j $O(1 \dots O(T))$ $N(N \cdot T)$

$$S_j^{(t+1)} = \max_{1 \leq i \leq N} [S_i^{(t-1)} a_{ij} b_j(D_t)]$$



$$\text{backpointer } [S_t] \leftarrow \underset{i}{\operatorname{argmax}} (|S|^2 \times L)$$

$$O(|T| \cdot |S| + |E|)$$

state change

$x_1 = F$

$\phi_L(t) = F$

A Q: in a state : sequence count

most likely sequence

backtrack
starte individually

✓ For each class / category!

sys / Pred

To detect spam

Difference

L

F

Total

TP

TN

✓

Gold

standard

L

TP

FN

Recall

FP

FN

x

standard

F

FP

FN

precision

more weight $\uparrow x$
FP \downarrow (unspam \rightarrow spam) medicare

system \rightarrow in favor Precision of L: $P_L = \frac{TP}{TP+FP}$

true/total Recall of L: $R_L = \frac{TP}{TP+FN}$

F-measure (β^2+1) $F_L = \frac{2P_L R_L}{P_L + R_L}$

Fn \downarrow (spam \rightarrow unspam)

correct \uparrow total

accuracy

TP+TN
FP+FN+TP+TN

1: equal importance
2: unbalanced to FP, FN
baseline: majority class prob
conservative metric
lower

L_{sys}

L_{truth}

SYS

* precisely cover I

coverage II

harmonic mean

mean

monotonically increase

when $t = \infty$

$$\text{table I: } S_{\text{Sj}}(t) = \log(\frac{p_{\text{Sj}}(C|t)}{p_{\text{Sj}}(U|t)})$$

Iterations: semi-supervised.

hidden: B Z i o M

self-training

O: HMM training data.

1. predict all

unlabeled data

2. Merge

{ observed: o, for i in O }

3. Train on merged

hidden: prediction

observed references as unlabeled data: larger dataset: performance

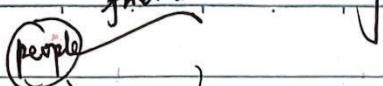
dev training data: meet before \propto use of

seed = 1 ↓ different rule.

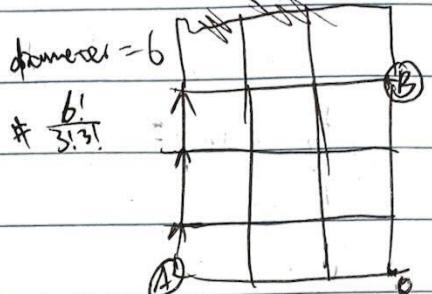
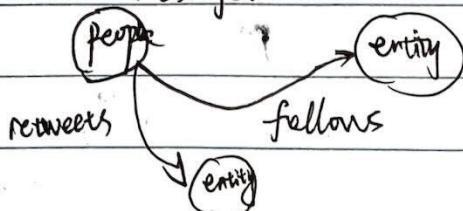
$$G = \langle V, E \rangle$$

Social Network friend clustering

Facebook



Twitter



Undirected, unweighted graphs.

the path that connects two nodes with the fewest edges

Distance: Length of shortest path between two nodes.

Diameter: Maximum distance between any pair of nodes.

$2^{n/4}$: Degree of a node: # neighbours / links of the node.

central & highly connected

$2^{n/4}$ / $\frac{2^{n/4}}{\text{periphery}}$ \downarrow $\frac{2^{n/4}}{\text{subcenter}}$

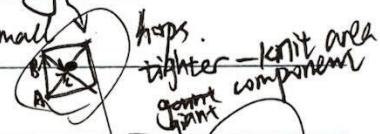
closely clustered regions, few links

most nodes are not neighbour of each other.

can be reached via small steps.

but nodes are likely to be neighbours of each other
small diameter & high clustering coefficient

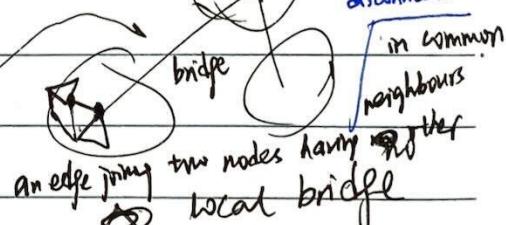
'Chain-links'!



giant component: connected; most of the graph

weak ties: socially distant ties

VS infrequent interaction.



Strong links: close friends and family

Cut \rightarrow increase len(path)

Triadic closure: ✓ $A \rightarrow B$, $B \rightarrow C$ strong ties, $A \rightarrow C$

closed triangles

↑ Global Clustering Coefficient

(closed + open) triads

number of edges between pairs of neighbours of v

/ # pairs of neighbours of node v .

Random links: $\frac{1}{d^2}$ nodes

edges between pairs of neighbours

/ # pairs of neighbours of

C_n^2

Clustering coefficient

Gracekeeper

"hub"; center of tightly-knit network

$y \rightarrow z$

$$(B(V)) = \sum_{S \in V} S(S|V)$$

$$S(s|V) = \sum_{t \in V} S(s,t|V)$$

No. Date ... 11.

"connectors" \nwarrow \nearrow \rightarrow disperse
only one

bridge functionality \rightarrow tight-knit area

node important in information flow

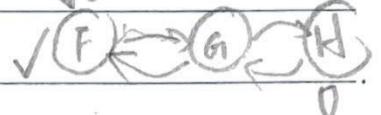
Betweenness centrality \sum relative # shortest path ($s \rightarrow t$) rely on the node gate-keeping

$$\cancel{S(s,t|V)} - CB(V) = \sum_{\substack{s \in V \\ t \in V}} \frac{\sigma(s,t|V)}{\sigma(s,t)}$$

path starting or ending from the $\cancel{S(s,t|V)}$ $\cancel{S(t,t)}$ total amount of flow it carries.

$\forall v \in V, CB(V) = 0$

$\forall s \in V$. $\cancel{- BFS \rightarrow \sigma(s,t)}$ storage efficiency.
count # v . $\Rightarrow \sigma(s,t|V)$ and count



Initialization: for each node w :

1. Mark w as unvisited by setting $dist[w]$ (the dist from s to w) to ∞ .
2. Set $Pred[w]$ precede w on a shortest path to $[]$.
3. Set $Paths[w]$ (the list of all shortest path from s to w) to $[]$.
4. While Q $\neq [s] \quad dist[s] = 0$: not empty.

$v = Q$. pop dequeue.

for w in v .neighbours:

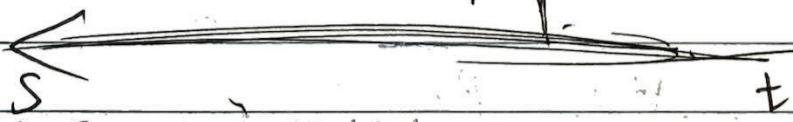
- if $dist[w] = \infty$:

$$dist[w] = dist[v] + 1$$

enqueue w .

- if $dist[w] = dist[v] + 1$:

append v to $Pred[w]$



$\forall s \in V$: $\underbrace{v \in s, t}_{\sigma(s,t|V)=0}$

1. set $S(v, t) = 0$,

2. BFS. $s \rightarrow t$

3. backward, increment $S(v, t)$

4. /

No.

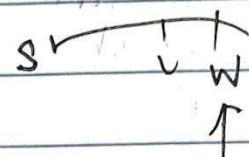
Date

$$\delta(S_t, t | V) = \text{unique}$$

1. $\delta(t) = 0$ if t is a terminal node

2. Increment $\delta(v)$ based on $\delta(w)$ $v \leftarrow w$.

3. $C_B(v)$



1. $\forall v \in V, C_B(v) = 0$

2. $\forall s \in V:$

1. set $\delta(v) = 0$.

$$\delta(v) = \sum_w (1 + \delta(w))$$

2. While Q is not empty:

1. $v = Q.pop()$

$S.push(v)$

2. $w \in v.\text{neighbours}$

if $\text{dist}[w] = \infty$:

$\text{dist}[w] = \text{dist}[v] + 1$

$Q.push(w)$

if $\text{dist}[w] = \text{dist}[v] + 1$:

$$\sigma(s, w) = \sigma(s, v) + \sigma(s, w)$$

append v to $\text{Pred}(w)$

3. while S not empty

~~$w = S.pop()$~~

correct order.

$\forall v \in \text{Pred}(w)$

$$\delta(v) \neq \text{Magic}(\delta(w))$$

$$[(v, w)] + \infty$$

if $w \neq s$:

$$C_B(w) \neq \delta(w)$$

/2

\forall pair of nodes, there is a path between them.
 Strongly Connected Component

dfs(u):
 not visited:
 visited [v] = true
 dfs($)$

for each node:

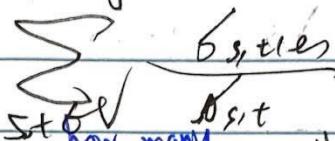
if u is not visited:
 visited [u] = true

connected-component $+1$

dfs(u)

union set

Edge betweenness



between all nodes
 ratio of the # shortest path from s to t
 crossing the edge and the total # shortest paths

Girvan-Newman : partition of network

while number of connected component < clusters & #edges > 0:

1. Calculate edge betweenness

2. Remove highest (with ties)

e^{-6}

3. Connected Components

closed

nodes \rightarrow similar

triadic closure

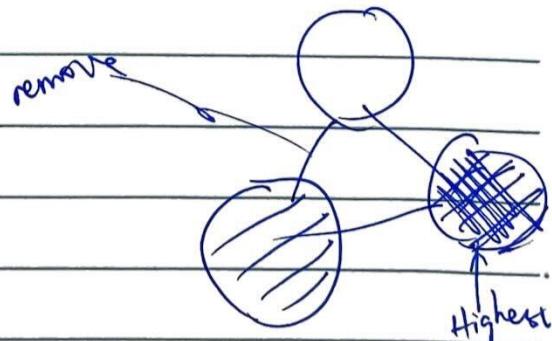
$\langle A, B, C \rangle$ Ratio

= 1 fully connected ✓



std. ()

Variance of betweenness
distribution



cluster criteria

✓ connectivity is uniform

No.

Date

General Case of S : directed
 X_S $\dashv \vdash$ can't compare.

T : undirected.
 X_T .

Bridge: an edge that connects two nodes that would otherwise be in disconnected components of the graph.

~~remove~~

A local bridge: an edge joining two nodes that have other neighbours in common.