

Sliding Windows \hookleftarrow This is new in Y.H's article

Matrix sketching \nwarrow Not that, that is old

$$\epsilon = \|A^T A - B^T B\|_2 / \|A\|_F^2 \quad \text{space complexity } \mathcal{O}\left(\frac{d}{\epsilon}\right)$$

Task: compress $A \in \mathbb{R}^{N \times d}$ to $B \in \mathbb{R}^{l \times d}$ where $l \ll N$

What's new:

\hookleftarrow need extra assumption, $A_{2(l)}$ is normalized

when A updates like a sliding window,

\hookleftarrow  can we also update B to track the new A .

Frequent Directions:

Initially, B are all zero.

If we have zero rows, then just put a into B

If not, we first append w to B , and then

apply SVD to shrink matrix B .

$$B = \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_{l-1} \end{pmatrix} V_{[1:l-1]}^T \quad \hookleftarrow \text{remove one row from } V$$

Why don't they just compute $A^T A$ and run SVD afterwards?

$$A \in \mathbb{R}^{N \times d} \quad A^T A$$

Useless shit X

Restate the problem (Vector codebook)

Given a function $f(\underset{\substack{\mathbf{w} \\ \mathbb{R}^{M \times d}}}{v}) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ (actually : $\{N: \mathbb{N}^+\} \rightarrow \mathbb{R}^{M \times d} \rightarrow (\mathbb{R}^d \rightarrow \mathbb{R}^d)$)

Goal : find another function $\hat{f}(\underset{\substack{\mathbf{w} \\ \mathbb{R}^{O(M \times d)}}}{w}) : \mathbb{R}^d \rightarrow \mathbb{R}^d$

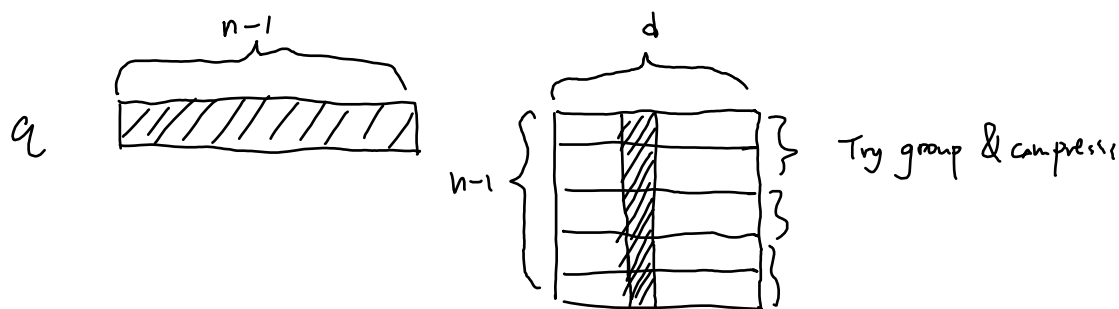
Expect some $p(N) = o(N)$ here

Heuristic : grab a machine learning model, update the model with v_n

$$g(w + v_n) \rightarrow w$$

such that $\hat{f}(w)$ behaves like $f(v)$

because $\hat{f}(w)$ approximates $f(w_{[1 \dots n-1]})$



$$f(w, q) = \text{softmax}(q \cdot v_{[1 \dots n-1][1 \dots d]}^T) v_{[1 \dots n-1][1 \dots d]}$$

$$f(w', q) = \text{softmax}(q \cdot v_{[1 \dots n][1 \dots d]}^T) v_{[1 \dots n][1 \dots d]}$$

$$f(w, q) \propto \exp(q \cdot v_{[1 \dots n-1][1 \dots d]}^T) v_{[1 \dots n-1][1 \dots d]}$$

$$f(w', q) \propto \exp(q \cdot v_{[1 \dots n][1 \dots d]}^T) v_{[1 \dots n][1 \dots d]}$$

Turn this into several constraints

$$f_v(q) = \text{softmax}(q v^T) v \quad v \in \mathbb{R}^{n \times d} \quad q \in \mathbb{R}^d$$

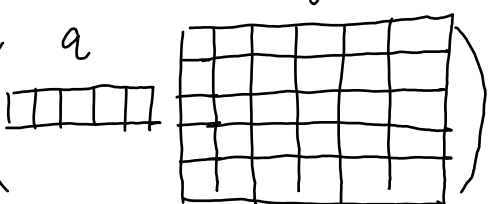
Don't view that as a direct computation problem, instead find a way to optimize it as some differential equation.

(Ask Jan for some ideas, or just read his paper)

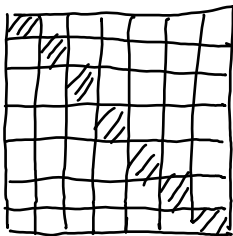
$v \rightarrow \phi(v)$ } embed v into some S
 • } a simple operation on S

$g(\cdot, \cdot)$ } a mapping from $\mathbb{R}^d \times S$ to \mathbb{R}^d
 can have some iterations in it

Doesn't feel very correct ...

$$\text{softmax}(q v^T) = \frac{1}{Z} \exp \left(\begin{array}{c} q \\ \hline \end{array} \begin{array}{c} v^T \\ \hline \end{array} \right)$$


$$= \exp(\text{diag}(q v^T)) \times \text{diag}(v_{\cdot j})$$



$$\exp(\text{diag}(q v^T)) \times V \times \mathbf{1}$$

$$\exp(\underbrace{A}_{\substack{\uparrow \\ n \times n}})$$

$$\frac{dx}{dt} = A x(t) \quad t \rightarrow \infty$$

Given q, v , we can run DE to obtain the result.

Fix the # of samples we use to compute the DE. \hookrightarrow Let it be m

We just get to compute a telescopic sum for each v_{ij}

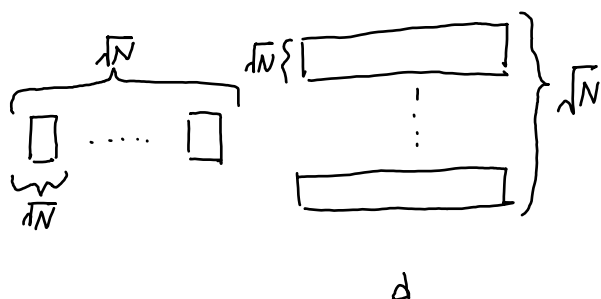
$$g_t((q v^T)_i, v_{ij}) = g_{t-1}((q v^T)_i, v_{ij}) \left(\frac{1}{m} + (q v^T)_i \right) \quad \text{不太有用}$$

q is incoming v

$$f(v, q) = \text{softmax}(q K v^T) v$$

$$\text{then, } v \leftarrow \begin{pmatrix} v \\ q \end{pmatrix}$$

不知道有啥用处



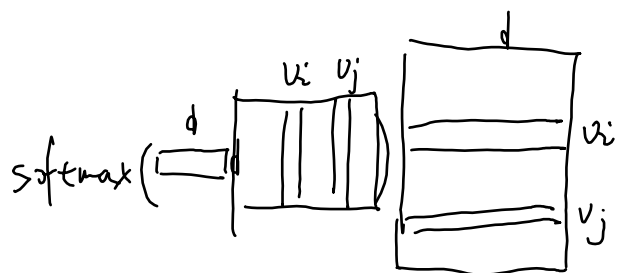
不知道引入什么误差能有什么收益

对新加入的 v , 找一个簇
然后找一个系数

Add more assumptions: v and q are both upperbounded by M

Simpler question: what happens when we merge v_i & v_j ?

what is optimal λ for i, j, v, q ?



$$\frac{\sum_{\ell} \exp(q v_{\ell}^T) v_{\ell}}{\sum_{\ell} \exp(q v_{\ell}^T)} = \frac{\sum_{\ell \in \{i, j\}} \exp(q v_{\ell}^T) v_{\ell} + v_i \exp(q v_i^T) + v_j \exp(q v_j^T)}{\sum_{\ell \in \{i, j\}} \exp(q v_{\ell}^T) + \exp(q v_i^T) + \exp(q v_j^T)}$$

$$\text{softmax}(qKv^T)v$$

q 和 $\text{softmax}(qv^T)v$ 有多像?

因为是算加权平均, 而和 q 最接近的 v_i 分到最大的权重

$$\hat{q} = qv^T(vv^T)^{\dagger}v \quad (\leftarrow N \text{ 不大的时候, 我猜是这种})$$

什么时候可以把 v_i 和 v_j 合并? 需要 v_i 和 v_j 各位都差不多

分开考虑 $v_{i,d}$ 可能更加合适

If $v_{i,d}$ is normalized, i.e. $\|v_{i,d}\|_2 = 1$

$$(\text{softmax}(qv^T))_i = \frac{\exp(-\frac{\|q - v_{i,d}\|_2^2}{2})}{\sum_j \exp(-\frac{\|q - v_{j,d}\|_2^2}{2})}$$

For each dimension $d \in [1..D]$

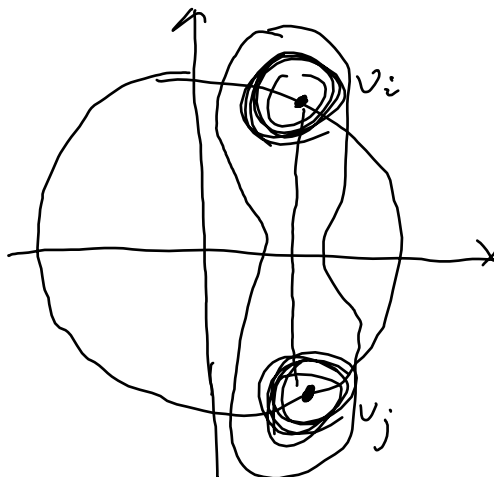
Compress or filter out Gaussians.

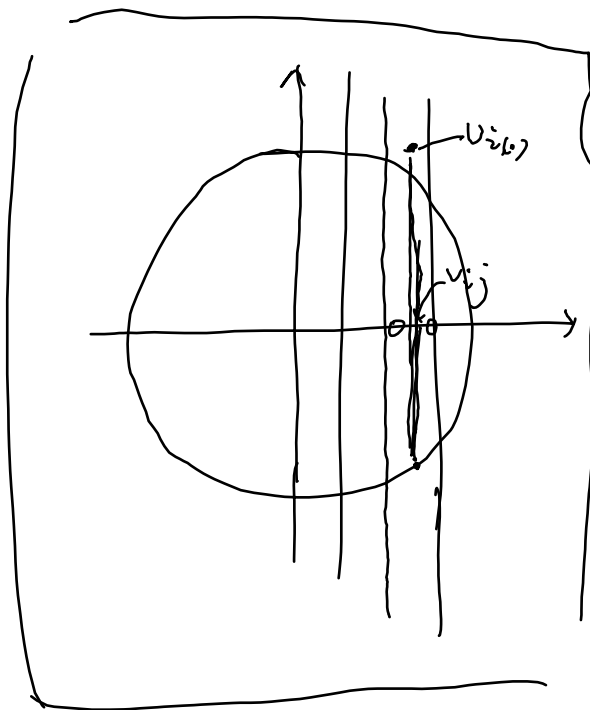
$$p_{\mathcal{I} \times \mathcal{X}}(i, x) = \frac{1}{N} \exp(-\frac{\|x - v_i\|^2}{2})$$

$$v_{i,1} = v_{j,1}$$

$$i \in [1..N]$$

$$\mathbb{E}[v_{\mathcal{I}} | X=q]$$





x_d

$$f_v(q)_j = \left[\begin{array}{c} \vdots \\ \vdots \end{array} \right] \cdot u_{(i,j)} \times d$$

$$j \in [1 \dots d]$$

$$\int_{x \in \mathbb{R}^d} \exp\left(\frac{x^T x}{2}\right) dx$$

$$= \int_{r \in \mathbb{R}^d} \exp\left(\frac{r^2}{2}\right) f(r) dr^2$$

p.d.f. (no dega...)

$$f_1 \downarrow p(x) = \frac{1}{(\sqrt{2\pi})^d \det(\Sigma_p)} \exp\left(-\frac{(x-\mu_p)^T \Sigma_p^{-1} (x-\mu_p)}{2}\right)$$

$$f_0 \downarrow q(x) = \frac{1}{(\sqrt{2\pi})^d \det(\Sigma_q)} \exp\left(-\frac{(x-\mu_q)^T \Sigma_q^{-1} (x-\mu_q)}{2}\right)$$

$$(x, I) \sim (f_i, i)$$

$$\|\mu_p - \mu_q\|_2$$

V for Variance

$$V[X] = E[V[X|I]]$$

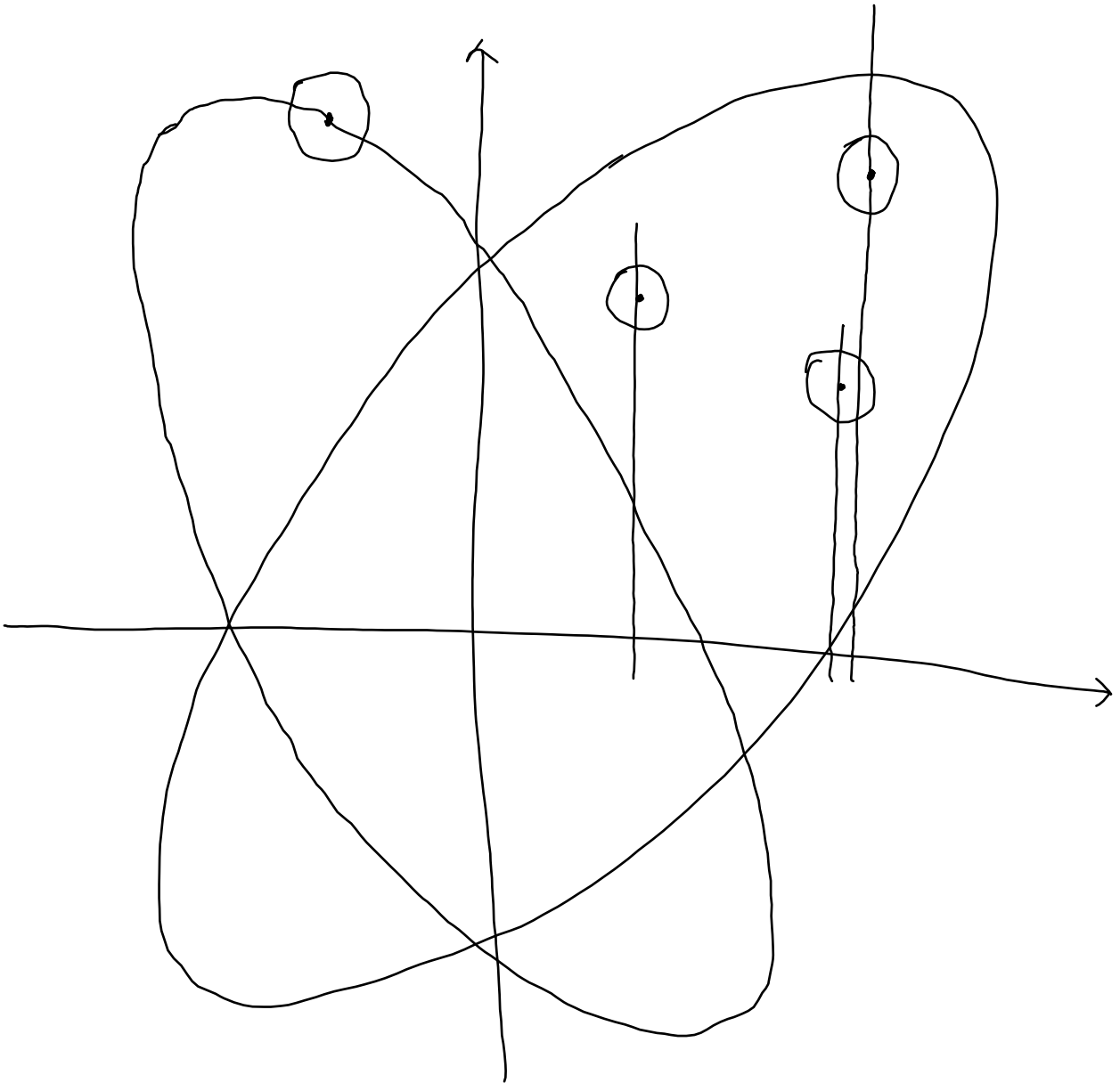
$$= \frac{1}{2} V[X|I=1] + \frac{1}{2} V[X|I=0]$$

$$= \frac{1}{2} (\Sigma_p + \Sigma_q)$$

$$E[X] = E[E[X|I]]$$

$$= \frac{1}{2} (\mu_p + \mu_q)$$

$$\sum_{\mathbf{q}, \mathbf{k} \in C} \exp(-\mathbf{q}^T \mathbf{q} / 2 + \mathbf{q}^T \mathbf{k} - \mathbf{k}^T \mathbf{k} / 2) \cdot (v \cdot \exp(\mathbf{k}^T \mathbf{k} / 2))$$



MCMC

Example : estimate π

$$X, Y \sim \text{Uniform}(-1, 1)$$

$$Z = [X^2 + Y^2 \leq 1]$$

$$Pr\{Z = 1\} = \frac{\pi}{4}$$

Theorem : Chernoff Bound

$$\{X_1, \dots, X_m\} \quad \mu = \frac{1}{m} \sum_{i=1}^m \mathbb{E}[X_i]$$

$$X_i(\omega) \in \{0, 1\} \quad Pr\left\{\left|\frac{1}{m} \sum_{i=1}^m X_i - \mu\right| \geq \epsilon \mu\right\} \leq \delta \Leftrightarrow m \geq \frac{3(\ln(2/\delta))}{\epsilon^2 \mu}$$

Definition : FPRAS (ϵ, δ) $\epsilon \ln \delta^{-1}$ size of input

Example : DNF Counting $f(x_1, \dots, x_n) = (x_1 \wedge x_2 \wedge \dots \wedge x_n) \vee (\dots)$
 \uparrow disjunction (or) \uparrow $\{0, 1\}$ \uparrow " x_i or " x_i " joined by " \wedge "

$$c(f) = \{(x_1, \dots, x_n) : f(x_1, \dots, x_n) = 1\}$$

$$\mathbb{E}[f(X_1, \dots, X_n)] = \frac{c(f)}{2^n}$$


$$X_i \sim \text{Uniform}(0, 1)$$

To give an (ϵ, δ) -estimation of $c(f)/2^n$, samples $m \geq \frac{3 \cdot 2^n \cdot \ln(2/\delta)}{\epsilon^2 c(f)}$

$$\exp \left(q \otimes \begin{array}{|c|} \hline \boxed{H} \\ \hline \end{array} Kx^T \right) xV$$

The type theory for notiz language

① I need the evaluation to be very efficient.

i.e. For a function term $x \rightarrow \underline{y}$. This syntax tree is traversed only once for subst x .  which means no reduce happen unless y has an "endpoint" type.

To keep track of values, we use a "Stack Variant" to store the value and removed abstraction layers.

↪ This is a bad idea. Because we have tuples.

② CPS to rescue

↪ Continuation Pass Style

Idea: for each term, find a way to represent "the next step"

Goal: for each term M , convert it to a term $\text{CPS}[M]$, such that $\text{CPS}[M, k]$ means when k is eventually called its argument $\equiv M$.

$$\text{CPS}[(M N), k] := \text{CPS}[N, \text{CPS}[M, k]]$$

$$\text{CPS}[\lambda x. M, k] := \lambda x. \text{CPS}[M, k]$$

$$\text{CPS}[N, \lambda x. \text{CPS}[M, k]] \equiv \text{CPS}[M[x/N], k]$$

This will cause an ever-growing stack.

So no X

③ De Bruijn with move

Define "last usage" of a variable, use movement

Case 1:

$\lambda x. \underbrace{CM}_{\substack{\uparrow \\ x}} \underbrace{N}_{\substack{\uparrow \\ x}}$ then "last usage" is in N
(CM is evaluated first)

Case 2:

$\lambda x. \underbrace{CM}_{\substack{\uparrow \\ x}} N$ then "last usage" is in M

Case 3:

$\lambda x. (\lambda y. \underbrace{M}_{\substack{\uparrow \\ x}})$ then "last usage" is in M

④ Composition

1. Modules are tree shaped (like rust)
2. Module-level parameters: different params results in different views when compiling standalone modules (code analysis).
3. Modules are imported and included with file path.
4. Module attributes are constructed using export statement (top level, "only once", same for mod-params)
5. "top level" is defined during lowering (before evaluation)

⑤ Defining Algebraic Data Types

Trouble: Consider evaluating function $f: A \times B \rightarrow C$

The problem is $\overset{\text{or } A}{\downarrow} B$ is not necessarily binded in for example
 $\lambda x. (f(a\ x))$, we shift x to stack,
but where to store a .

The evaluation should be designed against it.

Allow enum/tuple/struct types.

Markups will be encoded into these types.

Build core terms from face syntax.

During execution, values are somewhat byte encoded.

Divide terms to $\left\{ \begin{array}{l} \text{values} \\ \text{neutrals} \\ \text{reducibles} \end{array} \right.$

⑥ Clarify the misunderstanding of CPS

1. During CPS evaluation, it does not produce neutral terms.
2. Not free access is passed to continuation,
so no stack is involved.

⑦ Leave shallow captures and lift deep captures.

Neutrals should not be visited until they turn into values.

⑧ Modular design :

Each module is defined (applied) only once.

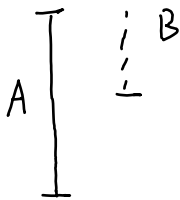
A module must be materialized before import.

In a file, dependency between modules can be determined by the appearance of import and module arguments.

If we import a module A before we define module B, then

B is dependent on A. Graph maintains start/end of a module.

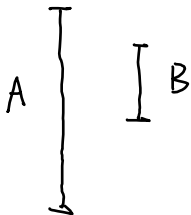
import (A; B) module A { import B; ... }



B must end before compile A ends.

Good: module A { };
module B { };

define (A; B) module A { module B(...); ... }



B starts/ends within A

after (A; B) ∃ module C { import A; module B(...); }



A compiles before B starts



⑨ Compiler Behaviour on Module Updates

Sps. we have a dep graph G , when updating a module, we first remove the original module from that graph, and mark the modules defined in that module and down stream modules as dirty.

Then, we add modified modules back and generate edges in graph G' .

During that process, check each module st. they have one and only one definition, and no dependency loop occurs.

Finally, run each module w.r.t. G' , update results when things are finished.

Dirty modules will not be re-evaluated if their input is equal.

⑩ Offloading : one observation is that exports are small but module output is large. We maintain an LRV-cached file system for that.

When evaluating include and import arguments, we always use the LRV-cached.

During evaluation, values are also offloaded to file system if they are big.

⑪ Module parameters on use v.s. on def

We may want module parameters to be different when we include modules.

* module parameters cannot appear in exported items

↑ This makes it so much easier

ref : reference to a codeblock in module

Allow generate state

⑫ How capturing works ?

⑬ With pointer to stack top



i means $\%rsp - 2$

so the "inner" translates as usual,

but the "outer" is captured in outer variables

⑭ extract all variables in one pass



Because we assume every nested function escapes



Make sure everything is captured



← During execution, this is everything we see.

When a symbol is discovered, we first find its definition level.

[...] $\lambda(x, y, z) \dots$
 [[...] $\lambda(w, v, u) \dots$
 (... y ...)

Then, we put the symbol along each level of current expressions' definition.



[...] $\lambda(x, y, z) \dots$
 [y ...] $\lambda(w, v, u) \dots$
 (... y ...)
 ↑
 add it here

The captured variables should link to their immediate upper level.

Problem: if we have multiple levels for insertion, we may have to shift previous de Bruijn indices.

Exact counting of Perfect Matching on Planar Graph

Theorem: Exactly counting PM in bipartite graph is
#P-hard

Ex: #SAT/count $f: \{0,1\}^N \rightarrow \{0,1\}$

$$\{x: f(x) = 0\}$$

Definition: A Fully Poly-time Randomized Approx Scheme (FPRAS)

for estimating z , with ϵ, δ error:

$$\Pr((1-\epsilon)z \leq \hat{z} \leq (1+\epsilon)z) \geq 1-\delta$$

in time $\text{poly}(n, \frac{1}{\epsilon}, \log(\frac{1}{\delta}))$

Definition: A $(\mathbb{F}) \mathbb{P} \mathbb{A}$ Sampler for target distribution μ ,
a randomized algo, with ϵ error:

$$d_{TV}(\nu, \mu) \leq \epsilon \text{ in } \text{poly}(n, \log(1/\epsilon))$$

(Thm) \exists FPRAS & FPAS for MP in bipartite graphs.

(Lemma) FPRAS for $|M|$ on all graphs (M : matchings)

\Leftrightarrow FPAS for μ on all graphs

\downarrow

uniform distribution over M

Proof. $\frac{1}{|M|} = \Pr_{M \sim \mu}(M = \emptyset) \quad E = \{e_1, \dots, e_m\}$

$$= \Pr_{M \sim \mu}(e_1 \notin M \wedge \dots \wedge e_m \notin M) \quad \begin{array}{l} \text{decompose} \\ \Pr(e_i \notin M | \dots) \end{array}$$

$$= \prod_{i=1}^m \Pr_{M \sim \mu / \{e_1, \dots, e_{i-1}\}}(e_i \in M)$$

estimate it with error ϵ/m

then est the count failure probability δ/m

\Rightarrow For $i=1$ to m : treat $M \in \bar{M}$ as $\sigma \in \{0,1\}^m$

$$p_i = \Pr_{\sigma \sim \mu}(\sigma_i = 1 \mid \sigma_1, \dots, \sigma_{i-1})$$

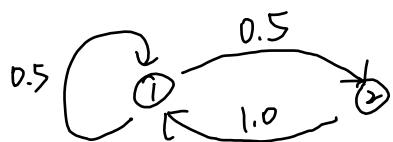
$$= \frac{|\bar{M}_{G' \setminus \{u_i, v_i\}}|}{|\bar{M}_{G'}|} \quad e_i = \{u_i, v_i\}$$

Running time will be $\text{poly}(m, \frac{1}{\epsilon})$

Boost to $\text{poly}(m, \log(1/\epsilon))$ with rejection sampling

(Def) Finite Markov Chain

(= random walk on directed graph)



(Def) A discrete-time & time-homogeneous Markov chain of a finite state space Ω is a sequence of R.V. $(X_t)_{t=0}^{\infty}$ valued in Ω satisfying:

$$\Pr(X_{t+1}=y \mid X_t=x, X_0) = \Pr(X_{t+1}=y \mid X_t=x) = p(x,y)$$

Remark: row sum $\sum_y p(x,y) = 1$

Basic Properties:

$$\textcircled{1} \Pr(X_{t+2}=y \mid X_t=x)$$

$$= \sum_z \Pr(X_{t+1}=z, X_{t+2}=y \mid X_t=x)$$

$$= \sum_z \Pr(X_{t+2}=y \mid X_{t+1}=z, X_t=x) \cdot \Pr(X_{t+1}=z \mid X_t=x)$$

$$\textcircled{2} \Pr(X_{t+l}=y \mid X_t=x) = p^l(x,y)$$

$$\textcircled{3} \text{ if } X_0 \sim \mu_0 \text{ then } \Pr(X_t=y) = (\mu_0 P^t)(y)$$

(Def) stationary distribution of Markov chain M

(P : transition matrix)

$$\pi P = \pi$$

Every Markov chain has at least one stationary distribution.

($P \mathbf{1} = \mathbf{1}$; take the row eigenvector of P with eigenvalue 1;
Perron-Frobenius Theorem)

Prf: Take row vector $vP = v$

$$\text{Consider } |v| := \begin{pmatrix} |v_1| \\ |v_2| \\ \vdots \\ |v_n| \end{pmatrix}$$

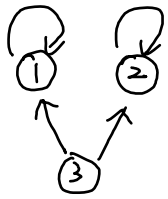
$$\text{Triangle Ineq: } |v| = |vP| \leq |v|P \quad \swarrow \text{entries}$$

$$\text{Also, } |v|P \mathbf{1} = |v| \mathbf{1}$$

$$\Rightarrow |v| = |v|P$$

Q: When can we find a unique stationary distribution?

Bad examples: The graph must be strongly connected.



$$\forall x, y \in \Omega : \exists t \in \mathbb{N} \text{ s.t. } P^t(x, y) > 0$$

Bad examples: The graph must be aperiodic.



$$\forall x \in \Omega : \text{period of } x := \gcd \{ t \in \mathbb{N} : P^t(x, x) > 0 \} = 1$$

$$\forall x \in \Omega \quad P(x, x) > 0 \Rightarrow \text{aperiodic}$$

Therefore: $P' = \frac{1}{2}(I + P)$ is aperiodic.

(Def) A finite MC is ergodic if it's aperiodic and irreducible.

$$\Leftrightarrow \exists t \in \mathbb{N} \text{ s.t. } \forall x, y \in \Omega : P^t(x, y) > 0$$

Fundamental Theorem of Markov Chains:

\forall ergodic finite MC P , \exists unique $\pi P = \pi$

$$\wedge \lim_{t \rightarrow \infty} P^t = \mathbf{1} \pi = \begin{pmatrix} \pi \\ \pi \\ \vdots \\ \pi \end{pmatrix}$$



(Def) Total variation distance

\forall distributions μ & π on finite Ω

$$d_{TV}(\mu, \pi) = \sum_{x \in \Omega} \mu(x) - \pi(x)$$

Property: $0 \leq d_{TV}(\mu, \pi) \leq 1$ \wedge $d_{TV}(\mu, \pi) \leq d_{TV}(\mu, \nu) + d_{TV}(\nu, \pi)$

(Def) A distribution w is a coupling of (μ, π) s.t.

$$\forall x \in \Omega \quad \sum_{y \in \Omega} w(x, y) = \mu(x)$$

$$\forall y \in \Omega \quad \sum_{x \in \Omega} w(x, y) = \pi(y)$$

(Def) Identity coupling for μ and μ :

$$w(x, y) = \begin{cases} \mu(x) & \text{if } x=y \\ 0 & \text{otherwise} \end{cases}$$

(Def) Product coupling for μ and π :

$$w(x, y) = \mu(x) \cdot \pi(y) \quad \leftarrow \text{independent case}$$

(Lemma) Coupling lemma :

$$(a) \quad d_{TV}(\mu, \pi) \leq \sum_{x, y} w(x, y) \quad \text{if } w \text{ is a coupling of } (\mu, \pi)$$

$$(b) \quad \exists w \text{ s.t. it takes "="}$$

(Proof of (b)) Let $\theta = \sum_x \min \{ \mu(x), \pi(x) \}$

$$\text{Set } w(x, x) = \min \{ \mu(x), \pi(x) \}$$

(Def) Consider a MC on Ω with transition matrix P
and two instances of (X_t, Y_t) ... joint distribution (X_t, Y_t)

(Proof) Fundamental theorem of MC (uniqueness)

$\forall x, y \in \Omega$ Ergodic means $P^t(x, y)$

Consider two MC (X_t) & (Y_t)

$$X_0 \sim \mu \quad Y_0 \sim \pi$$

$$\forall t \in \mathbb{N}, X_t \sim \mu, Y_t \sim \pi$$

Construct a coupling of X_t, Y_t

If $X_{t-1} = Y_{t-1} = x$ then

$$\text{Let } X_t \sim Q(x, \cdot) \text{ and } Y_t = X_t$$

If $X_{t-1} = x \neq y = Y_{t-1}$

$$\text{Let } (X_t, Y_t) \sim Q(x, \cdot) \otimes Q(y, \cdot)$$

$$Pr(X_t = Y_t | X_{t-1} = x, Y_{t-1} = y') \geq \epsilon$$

$$d_{TV}(\mu, \pi) \leq Pr(X_t \neq Y_t)$$

(Def) Pfaffian of a skew-symmetric matrix

$$\text{pf}(A) = \frac{1}{2^n n!} \sum_{\sigma \in S_n} \text{sgn}(\sigma) \prod_{i=1}^n a_{\sigma(2i-1)\sigma(2i)}$$

$$\text{pf}^2(A) = \det(A) \quad \left(\quad \right)$$

Motivation: to avoid sum over all possible permutations
find

(Lemma) If P is a symmetric matrix, then $\Omega = \text{uniform}(\Omega)$

(Proof) $P1 = 1 \Rightarrow P^T 1 = 1$

so $1^T P = 1^T$, $(\frac{1}{|\Omega|}, \dots, \frac{1}{|\Omega|})$ is the stationary distribution

(Definition) P is reversible w.r.t. π if

$$\forall x, y \in \Omega: \pi(x) \cdot P(x, y) = \pi(y) \cdot P(y, x)$$

(Lemma) If P is reversible w.r.t. π then π is stationary

$$(\pi P)(x) = \sum_y \pi(y) P(y, x) = \sum_y \pi(x) P(x, y) = \pi(x)$$

(Example) Random walk on undirected graph $G = (V, E)$

$$P(u, v) = \begin{cases} \frac{1}{\deg(u)} & \text{if } \{u, v\} \in E \\ 0 & \text{otherwise} \end{cases} \quad \begin{array}{l} \text{Ergodic means connected} \\ \& \text{not bipartite in this case} \end{array}$$

$$\pi(u) = \frac{\deg(u)}{2|E|}$$

(Definition) For $\epsilon \in (0, 1)$, $T_{\text{mix}}(\epsilon, P) = \max_{x \in \Omega} \min \{t \in \mathbb{N} : d_{TV}(P^t(x, \cdot), \pi) \leq \epsilon\}$

(Lemma) $T_{\text{mix}}(\epsilon) \leq T_{\text{mix}}(\frac{1}{4}) \lceil \log_2(1/\epsilon) \rceil$

(Example) Random walk on hypercube $H_n = \{0, 1\}^n$

From $x_t \in \Omega$, uniformly pick a bit and mutate it with $p = \frac{1}{2}$

$$\forall x: P(x, x) \geq \frac{1}{2}$$

$P(x, y) = P(y, x)$ so it is symmetric

(Theorem) Mixing Time of hypercube H_n

Pick two chains with H_n distribution (X_t, Y_t)

Sps. at t , $X_t(\omega) = Y_t(\omega)$ If two bits are equal, they stay equal.

$$\begin{aligned} d_{TV}(P_t^x(\cdot), \pi) &\leq P_r(X_t \neq Y_t) \\ &\leq E[d_H(X_t, Y_t)] \\ &\leq (1 - \frac{1}{n})^t n \leq \frac{1}{4} \end{aligned}$$

$$T_{mix}(\frac{1}{4}) = O(n \log n)$$

(Example) Proper Coloring

$G = (V, E)$ graph with max degree Δ

$\Omega \subseteq [q]^V$: a set of q -colorings of G , $q \geq \Delta$

s.t. $\forall \sigma \in \Omega$: $\sigma(v) \neq \sigma(u)$ for $uv \in E$

Goal: sample a q -coloring u.a.r. from G .

Choose $v \in V$ u.a.r. & $c \in [q] \setminus \underbrace{X_t(N(v))}_{\text{color of neighbors}}$

$$\text{Set } X_{t+1}(\omega) = \begin{cases} c & \text{if } w=v \\ X_t(\omega) & \text{if } w \neq v \end{cases}$$

(Lemma) If $q \geq \Delta + 2$, then G^D is ergodic

& $\pi = \text{unif}(\Omega)$ is stationary

(Theorem) If $q \geq 2\Delta + 1$ then $T_{mix} = O(n \log n)$

(Proof) Path Coupling

1. Construct $X_t = Z_0 - Z_1 - \dots - Z_n = Y_t$

$$Z_i(v_j) = \begin{cases} Y_t(v_j) & \forall j \leq i \\ X_t(v_j) & \text{if } j > i \end{cases}$$

2. Construct a coupling

3. "Glue"

(Theorem) Path Coupling

P : MC on finite Ω

Suppose $T = (\Omega, S)$ is a connected graph

$w: S \rightarrow [1, +\infty)$

$$d(x, y) = \min_{\substack{P: \text{path connecting } x \text{ and } y \\ \text{in } T}} \sum_{e \in P} w(e)$$

If $\forall x_0, y_0 \in S, \exists$ coupling of transition $(x_0, y_0) \xrightarrow{P} (x, y)$

$$\text{s.t. } \mathbb{E}[d(x, y)] \leq (1-\gamma) d(x_0, y_0) \quad X \sim P(x_0, \cdot)$$

$$Y \sim P(y_0, \cdot)$$

$$\text{then } T_{\text{mix}} = O\left(\frac{1}{\gamma} \log(\text{diam}(T, w))\right)$$

$$\text{where } \text{diam}(T, w) = \max_{x, y \in \Omega} d(x, y)$$

(Proof) $X_0 = x \quad X_t \sim P^t(x, \cdot)$

$Y_0 \sim \pi \quad Y_t \sim \pi$

$$d_{TV}(P^t(x, \cdot), \pi) \leq P_r(X_t \neq Y_t) \leq \mathbb{E}[d(X_t, Y_t)]$$

Coupling $(X_t, Y_t) \xrightarrow{P} (X_{t+1}, Y_{t+1})$

1. Take the shortest path in (T, w)

$$X_t = z_0 \cdots \cdots z_k = Y_t \quad z_{i-1}, z_i \in S$$

$$\text{s.t. } d(X_t, Y_t) = \sum_{i=1}^k w(z_{i-1}, z_i)$$

2. For $i=1$ to k :

Take a coupling of $(z_{i-1}, z_i) \xrightarrow{P} (w_{i-1}, w_i)$

$$\text{s.t. } \mathbb{E}[d(w_{i-1}, w_i) | (z_{i-1}, z_i)] \leq (1-\gamma) d(z_{i-1}, z_i)$$

3. Compose these coupling to get a ground coupling,

$$(z_0, z_1, \dots, z_k) \xrightarrow{P} (w_0, w_1, \dots, w_k)$$

where $w_i \sim P(z_i, \cdot)$

Set $X_{t+1} = w_0 \quad Y_{t+1} = w_k$

$$\begin{aligned}
\mathbb{E}[d(X_{t+1}, Y_{t+1}) | (X_t, Y_t)] &\leq \mathbb{E}[d(w_0, w_k) | (z_0, z_k)] \\
&\leq \sum_{i=1}^k \mathbb{E}[d(w_{i-1}, w_i) | (z_0, z_k)] \\
&\leq (1-\gamma) \sum_{i=1}^k d(z_{i-1}, z_i) \\
&= (1-\gamma) d(X_t, Y_t)
\end{aligned}$$

$$d_{TV}(P_t^X, \pi) \leq (1-\gamma)^t d(X_0, Y_0) \leq \frac{1}{4}$$

$$\text{when } t \geq \frac{1}{\gamma} \log(4 \text{diam}(T, w))$$

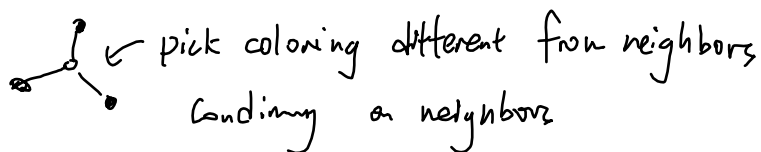
Product Space : $\Omega = [q]^n \quad n \in \mathbb{N}^+$

μ : distribution on Ω

$$X_t \xrightarrow{GD} X_{t+1}$$

1. Pick $i \in [n]$ u.a.r.

2. Sample $c \sim \mu_i^0 := \mu_i(\cdot | \mathcal{G})$ where $\mathcal{G} = (X_t)_{[n] \setminus \{i\}}$



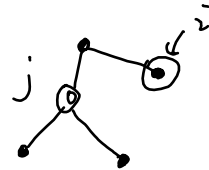
$$X_{t+1}(j) = \begin{cases} c & \text{if } j=i \\ X_t(j) & \text{if } j \neq i \end{cases}$$

Lemma: GD is reversible w.r.t. μ , so μ is stationary

(Def) Dobrushin Influence Matrix

$$R \in \mathbb{R}^{n \times n} \quad R(i, i) = 0 \quad \forall i \in [n]$$

$$R(i, j) = \max_{(b, \pi): \text{pair of} \\ \text{contig on } [n] \setminus \{j\} \\ \text{that differ only at } i} d_{TV}(\mu_j^b, \mu_j^\pi) \quad (i \neq j)$$



sample j conditioned on b & π respectively

(Theorem) Dob...-Shlosman uniqueness \Rightarrow mixing

$$\text{If } \|R\|_\infty \leq 1 - \delta \quad (\text{i.e., } \forall i \in [n], \sum_{j=1}^n R(i, j) \leq 1 - \delta)$$

$$T_{\text{mix}} = O\left(\frac{1}{\delta} n \log n\right)$$

You can replace $\|\cdot\|_\infty$ with $\|\cdot\|_p$

(Proof) $T = (\Omega, S)$ Hamming Graph $(x, y \in S \Leftrightarrow d_H(x, y) = 1)$
 \wedge
 unweighted

For $x, y \in \Omega$ with $d_H(x, y) = 1$, say $x_i \neq y_i$ & $x_{-i} = y_{-i}$

Construct a coupling $(x, y) \xrightarrow{GD} (X, Y)$

1. Pick $j \in [n]$ u.a.r.

2. Pick $(c_1, c_2) \in [q] \times [q]$ from an optimal coupling of $\mu_j(\cdot | X_{-j})$ & $\mu_j(\cdot | Y_{-j})$

$$3. \quad X(k) = \begin{cases} c_1 & \text{if } k=j \\ x(k) & \text{if } k \neq j \end{cases} \quad Y(k) = \begin{cases} c_2 & \text{if } k=j \\ y(k) & \text{if } k \neq j \end{cases}$$

$$\begin{aligned} \mathbb{E}[d_H(X, \gamma)] - 1 &\leq \frac{1}{n}(-1) + \sum_{j \neq i} \frac{1}{n} R(i, j) \\ &\leq -\frac{\delta}{n} \end{aligned}$$

Applications

Hardcore Model

$G = (V, E)$ graph of max deg Δ

$\mathcal{I}(G)$ independent sets of $G \subseteq \{0, 1\}^V$

$$\mu(\sigma) = \frac{1}{Z} \lambda^{|\sigma|} \quad \forall \sigma \in \mathcal{I}(G) \quad \text{where } |\sigma| = \sum_{v \in V} \sigma_v$$

GD for hardcore

1. $v \in V$ u.a.r.

2. If $(\exists w \in N(v) : X_t(w) = 1)$ set $X_{t+1}(v) = 0$

otherwise $(\forall w \in N(v), X_t(w) = 0)$ set $X_{t+1}(v) = \begin{cases} 1 & \text{w.p. } \frac{\lambda}{1+\lambda} \\ 0 & \text{w.p. } \frac{1}{1+\lambda} \end{cases}$

$$R(u, v) = \begin{cases} 0 & \text{if } v \notin N(u) \\ \frac{\lambda}{1+\lambda} & \text{if } v \in N(u) \end{cases}$$

So $R = \frac{\lambda}{1+\lambda} A$ A adj matrix

$$\|R\|_{\infty} = \frac{\lambda}{1+\lambda} \|A\|_{\infty} \leq \frac{\lambda}{1+\lambda} \Delta \leq 1 - \delta$$

$$\Leftrightarrow \lambda \leq \frac{1-\delta}{\Delta-1+\delta} \leq \frac{1-\delta'}{\Delta-1} \quad T_{\text{mix}} = O\left(\frac{n}{\delta} \log n\right)$$

Ising Model

$G=(V, E)$ graph of max deg Δ

$$u(\sigma) = \frac{1}{Z} \exp\left(\beta \sum_{(u,v) \in E} \sigma_u \sigma_v\right) \quad \sigma \in \{-1, +1\}^V$$

where $\beta \in \mathbb{R}$

More generally, $\mu(\sigma) = \frac{1}{Z} \exp\left(\frac{1}{2} \sigma^T J \sigma\right) \quad J \in \mathbb{R}^{n \times n} \quad \sigma \in \{-1, +1\}^n$

$$R(i, j) \leq \tanh |J_{ij}| = |J_{ij}|$$

$$\Leftrightarrow \|R\|_\infty \leq \|J\|_\infty$$

If $J = \beta A$ then $R \leq (\tanh |\beta|) A$

$$\Rightarrow \|R\|_\infty \leq \tanh |\beta| \cdot \Delta$$

$$\|J\|_\infty \leq 1 - \delta \Rightarrow T_{\text{mix}} = O\left(\frac{n}{\delta} \log n\right)$$

Random Matching

$G = (V, E)$ a graph

$\bar{\mu} = \bar{\mu}(G)$ set of matchings of G

Goal: sample $\mu \in \bar{\mu}$ u.a.r.

Construction of MC:

1. Choose $e = uv \in E$ u.a.r.

2. Consider four cases:

(i) (Add) u, v are unmatched



(ii) (Remove) u, v already matched

(iii) (Slide) If u is unmatched & v is matched by $e' = vw \in E$

then $X' = X_t \cup e \setminus e'$

(iv) Otherwise $X' = X_t$

$$3. X_{t+1} = \begin{cases} X' & \text{with prob } \frac{1}{2} \\ X_t & \end{cases}$$

Claim: This MC is ergodic & symmetric.
stationary distribution is $\text{unif}(\bar{\mu})$

For all $\sigma, \tau \in \Omega$, define a path γ in (Ω, P)

Observation: $\sigma \oplus \tau$ consists of alternating paths & alternating even cycles

Define the path $\sigma \sim \tau$

① Order components in $\sigma \oplus \tau$ by min vertex index

Given a transition:

$$P(\mu, \mu') = \frac{1}{\pi(\mu) \cdot P(\mu, \mu')} \sum_{(\sigma, \tau) \in \mathcal{P}_{\mu\mu'}} \pi(\sigma) \pi(\tau)$$

$$\left. P(\mu, \mu') = \frac{1}{2m} \right\} = 2m \cdot \frac{|\mathcal{P}_{\mu\mu'}|}{|\bar{\mu}|} \sim \text{try to prove } |\mathcal{P}_{\mu\mu'}| < |\bar{\mu}|$$

$$T_{\text{mix}} = O(m n^2 \log n)$$

construct an injection

$$\pi(\mu) = \frac{1}{2} \chi^{|\mu|} \quad \forall \mu \in \bar{\mu} \quad (\text{hardcore on line graph})$$