

Lecture 1. Coupling and Markov chain, basics

▷ Self-intro, email

▷ X : prob space, typically discrete & finite

1.1. Coupling

Motivating example:

Two players A, B each tosses a coin. A's coin has head prob = 0.75, B's coin has head prob = 0.5, each tosses 100 times.
Prove: $P(A \text{ gets } \geq k \text{ heads}) \geq P(B \text{ gets } \geq k \text{ heads})$

[Formal: $P(\text{Binom}(100, 0.75) \geq k) \geq P(\text{Binom}(100, 0.5) \geq k)$]

At each time, let $X_A \in \{0, 1\}$ be the outcome of A.

$X_B \in \{0, 1\}$ — — — of B

We can couple X_A and X_B s.t. $X_A = X_B$

Therefore the inequality becomes trivial.

Observation: LHS only depends on A's distribution

RHS — — — B's distribution

Therefore we can "design" their correlation in whatever way.

Def: A coupling of two prob distributions $\mu, \nu \in \mathcal{P}(X)$
is a joint distribution on $X \times X$.

if $(X, Y) \sim \pi$ then

$$P(X=x) = \mu(x) \quad P(Y=y) = \nu(y) \quad \forall x, y$$

P. Joint distribution that maintains marginal

Given μ, ν , a coupling can be described as

a matrix $q \in \mathbb{R}^{|\mathcal{M}| \times |\mathcal{N}|}$

with $q(x, y) = P(X=x, Y=y)$ satisfying

$$\sum_y q(x, y) = \mu(x), \quad \sum_x q(x, y) = \nu(y)$$

Example :

①. Couple $\text{Ber}(0.5)$ and $\text{Ber}(0.75)$

②. Couple $\text{Ber}(0.5)$ and $\text{Ber}(0.5)$

③. Couple $\text{Ber}(p)$ and $\text{Ber}(q)$ (2-d bond percolation)

1.2.

Coupling and total variation distance

Def: Given μ, ν on X , their TV distance is

$$\|\mu - \nu\|_{TV} := \sup_A |\mu(A) - \nu(A)|$$

$$\text{Lemma: } \|\mu - \nu\|_{TV} = \sup_A |\mu(A) - \nu(A)|$$

$$= \sup_B \mu(A) - \nu(A)$$

$$= \sup_B \nu(B) - \mu(B)$$

$$= \sum_{\{x: \mu(x) > \nu(x)\}} \mu(x) - \nu(x)$$

$$= \frac{1}{2} \sum_{x \in X} |\mu(x) - \nu(x)|$$

[TV distance is L^1 -distance]

Coupling lemma: Fix μ, ν on X , let $\Gamma(\mu, \nu)$ be the set of all the couplings between μ and ν . Then

$$\|\mu - \nu\|_{TV} = \inf_{(X, Y) \sim \Gamma(\mu, \nu)} P(X \neq Y)$$

($(X, Y) \sim \Gamma(\mu, \nu)$ means joint dist. of (X, Y))

$$= 1 - \sup_{(X, Y) \sim \Gamma(\mu, \nu)} P(X = Y)$$

D. Any coupling satisfying $\|\mu - \nu\|_{TV} = 1 - P(X = Y)$ is called "maximal coupling".

Pf:

$$\|M - V\|_{TV} = \frac{1}{2} \sum_x |M(x) - V(x)|$$

$$= \frac{1}{2} \sum_{\{x : M(x) \geq V(x)\}} |M(x) - V(x)| + \frac{1}{2} \sum_{\{x : V(x) > M(x)\}} |V(x) - M(x)|$$

$$= \frac{1}{2} \sum_x \max \{M(x), V(x)\} - \min \{M(x), V(x)\}$$

$$= 1 - \sum_x \min \{M(x), V(x)\}$$

For any coupling $q \in \Gamma(M, V)$

$$P(X=Y=x) \leq \min \{M(x), V(x)\}$$

$$\Rightarrow P(X=Y) \leq \sum_x \min \{M(x), V(x)\}. \quad \square$$

Lemma: Maximal coupling always exists

Def $\tilde{M}(x) \propto M(x) - \min \{M(x), V(x)\}$

$\tilde{V}(x) \propto V(x) - \min \{M(x), V(x)\}$

Consider the following way of generating (X, Y)

①. Flip coin with head prob $\sum_i \min\{Mx_i, \gamma_i x_i\}$

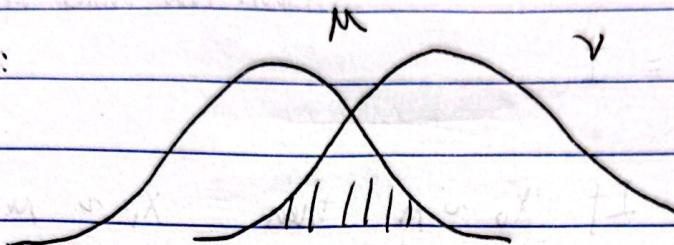
② If heads, generate $X \sim \min\{M, \gamma\}$, set $Y = X$

If tails, generate $X \sim \tilde{M}$, $Y \sim \tilde{\gamma}$ independently

Therefore $P(X=Y) \geq \sum_i \min\{Mx_i, \gamma_i x_i\}$

D

Geometric:



TV distance = 1 - shaded area

1.3. Markov chain

A Markov chain on \mathcal{X} is a stochastic process (X_0, X_1, \dots)

such that for any (x, y) and any x_0, \dots, x_{t-1} , and any t
such that $P(X_0 = x_0, X_1 = x_1, \dots, X_{t-1} = x_{t-1}, X_t = x) \neq 0$

we have

$$\begin{aligned} & P(X_{t+1} = y \mid X_0 = x_0, \dots, X_{t-1} = x_{t-1}, X_t = x) \\ & \cong P(X_{t+1} = y \mid X_t = x) \quad [\text{Markov property}] \end{aligned}$$

⇒ Future depends only on present, but not past

If $P(X_{t+1} = y \mid X_t = x)$ does not depend on t , then Markov chain is homogeneous.

We say $P \in \mathbb{R}_+^{(X \times X)}$ is the transition matrix of M.C.

$$\text{if } P(X_{t+1} = y | X_t = x) = P(x, y)$$

↑ ↑
y-th col
x-th row

- ▷ Transition matrix: non-zero non-negative entry
each row sums up to 1.

$$P \vec{1} = \vec{1}$$

Exercise: If $X_0 \sim \mu$ then $X_1 \sim \mu P$

Exercise: $P^t(X_{s+t} = y | X_s = x) = P^t(x, y)$

(Transition matrix after t steps is just P^t).

Def: μ is called stationary distribution of IP

$$\text{if } \mu P = \mu.$$

Def: IP is called irreducible if for any (x, y)

there is t s.t. $P^t(x, y) > 0$.

IP is called aperiodic if $\text{gcd } T(x)$

$$:= \text{gcd} \{n, P^n(x, x) > 0\} =$$

IP is reversible with respect to π

$$\text{if } \pi(x) P(x, y) = \pi(y) P(y, x).$$

Prop: IP reversible w.r.t. $\pi \Rightarrow \pi P = \pi$.

Example : Random walk on \mathbb{Z} ($* x = \mathbb{Z}$)

$$P(x, x+1) = P(x, x-1) = \frac{1}{2}$$

Example : Random walk on graphs.

A graph $G = (V, E)$

$\begin{matrix} \uparrow & \uparrow \\ \text{vertex} & \text{edge} \end{matrix}$

$$E \subseteq \{(x, y) \in V \times V, x \neq y\}$$

When $\{x, y\} \in E$, we say y is a neighbor of x ,
 $\deg(x) = \# \text{ of neighbors of } x$.

Fix G , define $P(x, y) = \begin{cases} \frac{1}{\deg(x)} & y \sim x \\ 0 & \text{otherwise} \end{cases}$

▷ $G = \star$

▷ check : P reversible w.r.t. $\pi(x) \propto \deg(x)$.

▷ application : search, pagerank, spectral graph theory.

Example : Metropolis-Hastings algorithm.

Target : sample from π

Input : Initialization $x \in \mathcal{X}$, proposal distribution $q(\cdot | \cdot)$

For t in $1:T$:

draw $y \sim q(x, \cdot)$

set $x_t = \begin{cases} y & \text{w.p. } \frac{\pi(y) q(x, y)}{\pi(x) q(y, x)} \\ x_{t-1} & \text{w.p. } 1 - * \end{cases}$

Check: reversible w.r.t. π .
application: physics, statistics / ML, theoretical CS.

Example: Shuffling

$$X = S_n$$

$$P(\sigma, \tilde{\sigma}) = \begin{cases} \frac{1}{\binom{n}{2}} & \text{if } \tilde{\sigma} = \tau \circ \sigma \text{ for some transposition} \\ 0 & \text{otherwise} \end{cases}$$

Check: reversible w.r.t. $\text{Unif}(S_n)$.

Example: Kar's walk. on sphere ($X = S^{p-1}$)

$$x_0 = (x_{0,1}, \dots, x_{0,p}) \in \mathbb{R}^p \text{ with } \sum_{i=1}^p x_{0,i}^2 = 1$$

Each time pick two entries (i, j)

generate $\theta \sim \text{Unif}[0, 2\pi)$

$$\text{change } (x_i, x_j) \text{ to } \begin{pmatrix} x_i \\ x_j \end{pmatrix} \xrightarrow{\begin{pmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{pmatrix}} \begin{pmatrix} x_i \\ x_j \end{pmatrix}$$

▷ Spectral gap (AOIP, 2001)

Spectral gap (Acta, 2003)

Mixing time (AAP, 2017).

▷ Application: physics, quantum

1.4 Stationary distribution and convergence.

Thm: Let P be the transition kernel of a M.C.

Suppose it is irreducible and aperiodic, then there exists a unique stationary distribution π .

Pf: Existence: Skipped, see Levin & Peres 1.5.9

or Perron-Frobenius thm.

Uniqueness: we show the following lemma

lem: If $Ph = h$ then h must be const.

Pf: If not, assume h maximized at x_0 and

there is y_0 s.t. $h(y_0) < h(x_0)$ and r.s.t. $P^r(x,y)$

Then

$$(P^r h)(x_0) = \sum_y h(y) P^r(x_0, y) < h(x_0) \text{ contradiction}$$

convince yourself.

This implies

$$\ker(P - I) = \{\lambda \vec{1} \mid \lambda \in \mathbb{R}\}$$

If π_1, π_2 stationary distribution, then

$$\pi_1(P - I) = \pi_2(P - I) = 0$$

$$(P^T - I) \pi_1^T = (P^T - I) \pi_2^T = 0$$

$$\Rightarrow \dim \ker(P^T - I) \geq \dim \{\text{span } \{\pi_1^T, \pi_2^T\}\} \geq 2$$

$$\text{But } \dim(\ker(P^T - I)) = \dim(\ker(P - I)). \quad \square$$