

Evita Nesteridi - Coupon collecting and strong stationary times

Lesson 1 -

n coupons

Every day a collector buys a coupon.

All coupons appear with some probability. ~~Let~~ Let T be the first time we have collected them all. Want to know how big T must be.

- Lemma:
- a) $P(T > t) \leq n(1 - \frac{1}{n})^t$
 - b) $P(T > n \log n + cn) \leq e^{-c}$ (doesn't take longer than $n \log n$ and a little bit)
 - c) $P(T < n \log n - cn) \geq e^{-c}$ [we are concentrated along $n \log n$]

Proof: a) $P(\text{coupon } i \text{ was purchased at time } t) = \frac{1}{n}$ All the purchases are independent
 $\Rightarrow P(\text{not getting it at time } t) = 1 - \frac{1}{n}$

$$P(T > t) = P\left(\bigcup_{i=1}^n \left\{ \begin{array}{l} \text{coupon } i \text{ has} \\ \text{not been collected} \\ \text{at time } t \end{array} \right\}\right) \leq \sum P(\text{coupon } i \text{ not collected by time } t) = \sum_{i=1}^n \left(1 - \frac{1}{n}\right)^t = n\left(1 - \frac{1}{n}\right)^t$$

b) $\left(1 - \frac{1}{n}\right)^n \leq \frac{1}{e}$

$$P(T > n \log n + cn) \leq n\left(1 - \frac{1}{n}\right)^{n[\log n + c]} \leq n\left(\frac{1}{e}\right)^{\log n + c} \leq e^{-c}$$

c) let $T_i =$ ~~first time~~ time that i -th time we got a new coupon

$$T = T_n = \sum_{i=1}^n (T_i - T_{i-1}) \text{ as } T_0 = 0$$

Markov's inequality: X random variable $P(X > aE(X)) \leq \frac{1}{a}$
← previous trick - s to be chosen later.

$$P(T < t) = P(-sT > -st) = P(e^{-sT} > e^{-st}) = P\left(\frac{e^{-sT}}{E(e^{-sT})} > \frac{e^{-st}}{E(e^{-sT})}\right) \leq \frac{e^{-st}}{E(e^{-sT})} E(e^{-sT})$$

$$= e^{-st} E\left(e^{-s \sum_{i=1}^n (T_i - T_{i-1})}\right) = e^{-st} E\left(e^{-sX_1} e^{-sX_2} \dots e^{-sX_n}\right) = e^{-st} \prod_{i=1}^n E(e^{-sX_i})$$

$X_i = T_i - T_{i-1}$ (how long it took for the i -th coupon to appear after the $(i-1)$ st)

Geometric distribution \rightarrow Geometric $\left(\frac{n-i+1}{n}\right)$

$$= e^{-st} \prod_{i=1}^n \frac{\frac{i-1}{n}}{e^{-s} - 1 + \frac{i-1}{n}}$$

← check in exercise session

If $t = n \log n - cn$, $s = \frac{1}{n}$, $e^s = e^{\frac{1}{n}} = 1 + \frac{1}{n} + \dots \geq 1 + \frac{1}{n}$

$$\Rightarrow P(T < t) \leq e^{\log n - c} \prod_{i=2}^{n+1} \frac{(i-1)/n}{i/n} = e^{\log n - c} \prod_{i=2}^{n+1} \frac{i-1}{i} = \frac{n}{n+1} e^{-c} < e^{-c} \quad \square$$

Will be using this result about coupon collecting over and over again.

This type of argument can be found on Feller - Introduction to probability (also stack exchange is a good resource)

Card shuffling: Random-to-top: deck of n cards. Pick a card uniformly at random and move it to the top. How many times do we need to do this ~~before~~ before the cards have been well shuffled.

People don't shuffle like this but computers often perform tasks ~~like this~~ in this way.

1	2	3	
2	3	2	
3	⋮	4	↪ (123)
⋮	⋮	⋮	
n	n	n	

↪ (12) ∈ S_n |S_n| = n!

Transition matrix : $x, y \in S_n$ $P(x, y)$ = probability of $x \rightarrow y$ after one step

$$= \begin{cases} \frac{1}{n} & \text{if } y = x(12 \dots i) \text{ for } i=1, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

We are interested in $P_x^t(x, y)$ = probability $x \rightarrow y$ after t steps. Now we can't write down a formula.
↑ t-th power It's useful to compute eigenvectors and eigenvalues

$P_x^t(y)$ is a probability measure on S_n ($\sum_x P_x^t(y) = 1$).

Tomorrow we'll see this converges to the uniform measure as $t \rightarrow \infty$. We'll study the convergence

Used metric is the total variation distance Persi introduced

$$\|P_x^t - U\|_{TV} = \frac{1}{2} \sum_{y \in S_n} |P_x^t(y) - \frac{1}{n!}|$$

uniform measure on S_n

Separation distance : $S_x(t) = \max_{y \in S_n} \left\{ 1 - \frac{P_x^t(y)}{U(y)} \right\}$ exercise: this is positive (but not a measure)

$S(t) = \max_x \{ S_x(t) \}$

Lemma (Aldous - Diaconis) $\|P_x^t - U\|_{TV} \leq S(t)$ for every $x \in S_n$

~~Def~~ The mixing time is defined as

$$t_{\text{mix}}(\epsilon) = \min \{ t > 0 : \|P_x^t - U\|_{TV} \leq \epsilon \}$$

Each time we bring k distinct cards to the top, ~~then~~ ^{orderings} all $k!$ ~~orderings~~ of the first k cards should appear with some probability. How long should we wait until we have all n cards? $n \log n!$ (coupon collecting).

T = first time that all cards have been selected

T_i = first time that the i -th new card appeared

What we have said can be expressed as

$$P_x^t(y | T \leq t) = \frac{1}{n!}$$

Went to look this to $P_x^t(y) = P_x^t(y | T \leq t) P(T \leq t) + P_x^t(y | T > t) P(T > t)$ (Bayes formula)

$$\geq P_x^t(y | T \leq t) P(T \leq t) = \frac{1}{n!} P(T \leq t)$$

Plugging this into the separation distance we have $S_x(t) \leq \max_{y \in S_n} \left\{ 1 - \frac{1/n! P(T \leq t)}{1/n!} \right\} = P(T > t)$

So we have proved

Theorem (Aldous - Diaconis) If $t_{n,c} = n \log n + cn$ then $S(t_{n,c}) \leq e^{-c}$



Suppose we have a process $x_1, x_2, \dots, x_t, \dots$ (sequence of random variables)

$$\text{if } P(X_{t+1} = x \mid X_1 = x_1, \dots, X_t = x_t) = P(X_{t+1} = x \mid X_t = x_t)$$

we say this is a Markov process on Ω (before $\Omega = \mathbb{S}_n$)

$$P = (P(x, y))_{x, y \in \Omega} \quad \text{if } \exists \text{ measure } \pi: \Omega \rightarrow [0, 1] \text{ such that } P_x^t \rightarrow \pi \text{ as } t \rightarrow \infty.$$

$$t_{\text{mix}}(\epsilon) \stackrel{\text{def}}{=} \min \{t \geq 0 : \max_x d_x(t) < \epsilon\}$$

time distance

$T: \Omega \times \dots \times \Omega \times \dots \rightarrow \mathbb{N}$ random variable, is a stopping time for (X_t) if $\{T \leq t\}$ depends only on x_1, x_2, \dots, x_t .

Example: $T =$ first time that card 1 is on top.

Definition: T is a strong stationary time for (X_t) if $P_x(X_t = y \mid T \leq t) = \pi(y)$

Lemma (A-D) $S(t) \leq P(T > t)$.

The proof is exactly the same as in the special case discussed above.

Tomorrow: 1) Eigenvalues and eigenfunctions of transition matrix in mixing
2) Comparison

