

## Lecture 12. Stationary distribution. Bins and Balls.

### Bucket sort

Martin Klazar

December 22, 2020

### Stationary distribution

In this lecture every Markov chain  $M = (X_0, X_1, \dots)$  is finite and has the set of states  $S := [n] = \{1, 2, \dots, n\}$ ,  $n \in \mathbb{N}$ . The evolution of  $M$  is determined by (i) the initial *distribution*

$$\bar{p} = (p_1, p_2, \dots, p_n) \in \mathbb{R}_{\geq 0}^n, \quad p_i = \Pr(X_0 = i),$$

satisfying  $\sum_i p_i = 1$ , which is the starting distribution of probabilities on the states of  $M$ , and (ii) the stochastic *transition matrix*

$$P = (p_{i,j})_{i,j=1}^n \in \mathbb{R}_{\geq 0}^{n \times n}$$

satisfying  $\sum_j p_{i,j} = 1$  for every  $i \in [n]$ , which records transition probabilities between the states:  $p_{i,j}$  is the conditional probability that  $M$  evolves in one step from the state  $i$  to the state  $j$ , so  $p_{i,j} = \Pr(X_{t+1} = j \mid X_t = i)$  for every  $t \in \mathbb{N}_0$  if the right side is defined.

**Exercise (on matrices).** If in matrices  $A \in \mathbb{R}^{k \times l}$  and  $B \in \mathbb{R}^{l \times m}$  every row sum equals 1, then so does in their product  $AB$ .  $\square$

This part of the lecture will actually be more a lecture in linear algebra or linear programming than a lecture in the probabilistic method. Our goal is to show that for any  $P$  satisfying certain obvious necessary conditions there exists a unique *stationary distribution*  $\bar{p}$  that is moreover *attractive*. A distribution  $\bar{p}$  is stationary if it is a fixed point in the evolution of  $M$ ,

$$\bar{p}P = \bar{p}$$

for  $\bar{p}$  written as a  $1 \times n$  row vector. Attractivity of  $\bar{p}$  means that  $M$  evolves towards it from any start, for every distribution  $\bar{q}$  we have

that, coordinatewise,

$$\lim_{n \rightarrow \infty} \bar{q}P^n = \bar{p} .$$

As before we associate with  $P$  the weighted *transition digraph*

$$D_P = ([n], E, h)$$

where  $e = (i, j) \in E$  iff  $p_{i,j} > 0$ , and the edge (arrow)  $e$  has weight  $h(e) = p_{i,j}$ . In any digraph  $D = (V, E)$  (so  $E \subset V \times V$ ), a *walk*  $w$  is any sequence of vertices

$$w = (v_0, v_1, \dots, v_n) \subset V, \quad n \in \mathbb{N}_0 ,$$

such that  $(v_{i-1}, v_i) \in E$  for every  $i = 1, 2, \dots, n$ . If  $v_0 = v_n$ , we call  $w$  a *cycle*. We write  $|w| := n$  for the *length* of a walk or a cycle. If  $v_0 = u$  and  $v_n = v$ , we call  $w$  a *u-v walk*. It follows that in any *u-v walk*  $w$  with minimum length no vertex is repeated and therefore  $|w| \leq |V| - 1$ . We call such  $w$  a *shortest path from u to v*.

We may think of  $M$  in terms of the following game in a sandbox. The given distribution  $\bar{p}$  distributes one kilogram of sand into heaps  $p_i$  at the vertices  $i \in [n]$  of  $D_P$ , and the matrix  $P$  is a redistribution rule. In one step we simultaneously move, according to the arrows

$$\bullet \xrightarrow{p_{i,j}} \bullet$$

of  $D_P$ , for every vertex  $i \in [n]$  and any  $j \in [n]$  the fraction  $p_{i,j} = h(i, j)$  of the heap  $p_i$  (that is, the amount  $p_i p_{i,j}$ ) to the vertex  $j$ . By this the old heap  $p_i$  disappears (as  $\sum_j p_{i,j} = 1$ ) but a new heap  $p'_j$  is formed at each vertex  $j \in [n]$ . In this way we redistribute the sand to the vertices and get its new distribution

$$(p'_1, p'_2, \dots, p'_n) = \left( \sum_i p_i p_{i,1}, \sum_i p_i p_{i,2}, \dots, \sum_i p_i p_{i,n} \right) .$$

The total amount of sand is indeed preserved because

$$\sum_{j=1}^n \sum_{i=1}^n p_i p_{i,j} = \sum_{i=1}^n p_i \sum_{j=1}^n p_{i,j} = \sum_{i=1}^n p_i = 1 .$$

If this transformation has not changed the distribution of sand,  $\bar{p}' = \bar{p}$ , the distribution  $\bar{p}$  is stationary.

The necessary conditions on  $P$  for existence of a unique stationary distribution that is attractive, mentioned above, are (i) *irreducibility of  $D_P$*  and (ii) *aperiodicity of  $D_P$* . The former means that for every two vertices  $i, j \in [n]$  there exists in  $D_P$  an  $i$ - $j$  walk, and the latter requires that the lengths of cycles in  $D_P$  be altogether coprime: the only  $d \in \mathbb{N}$  dividing all of them is  $d = 1$ . One also says that the matrix  $P$ , or the Markov chain  $M$ , is irreducible, respectively aperiodic.

It is easy to see that (i) and (ii) are in general necessary. Consider the non-irreducible digraph  $D_P = ([3], E, h)$  with arrows  $1 \rightarrow 2$  and  $1 \rightarrow 3$ , loops  $2 \rightarrow 2$  and  $3 \rightarrow 3$ , and weight  $\frac{1}{2}$  on the arrows and 1 on the loops (the corresponding matrix  $P$  is stochastic). A stationary distribution exists but is not unique, both distributions  $(0, 1, 0)$  and  $(0, 0, 1)$  are stationary. Consider the non-aperiodic digraph  $D_P = ([2], E, h)$  with only two arrows  $1 \rightarrow 2$  and  $2 \rightarrow 1$ , weighted by 1 (the corresponding matrix  $P$  is stochastic). There is the unique stationary distribution  $(\frac{1}{2}, \frac{1}{2})$  but it is not attractive, any distribution  $(p, 1 - p)$ ,  $p \in [0, 1]$ , only evolves to  $(1 - p, p)$  and back.

For the proof of the theorem on stationary distribution we need four lemmas. In the proof I follow the lecture notes of J. Sgall.

- *Jiří Sgall* is, by the Czech mutation of Wikipedia, Czech computer scientist and mathematician working in the areas of approximation algorithms, online algorithms and theory of scheduling. Currently he is the vice-dean of the School of Computer Science at the MFF UK (and hence lecturer's superior).

**Lemma 1.** *If the numbers  $a_1, \dots, a_k \in \mathbb{N}$ ,  $k \in \mathbb{N}$ , are altogether coprime then there is an  $n_0 \in \mathbb{N}_0$  such that for every integer  $n \geq n_0$  there exist numbers  $m_1, \dots, m_k \in \mathbb{N}_0$  with*

$$n = \sum_{i=1}^k m_i a_i .$$

**Proof.** It is a well known result in elementary number theory that then  $1 = \sum_{i=1}^k r_i a_i$  for some numbers  $r_i \in \mathbb{Z}$ . We set  $r := \max_i |r_i| \in \mathbb{N}$  and  $n_0 := \sum_{i=1}^k (r a_i) a_i$ . Adding to each of the numbers  $n_0, n_0 + a_1, n_0 + 2a_1, \dots$  repeatedly the linear combination expressing 1, we see that

every natural number from  $n_0$  on is an integral linear combination of the  $a_i$ s with nonnegative (here even positive) coefficients.  $\square$

**Exercise (on compactness).** Show that every set  $A \subset \mathbb{N}$  of altogether coprime numbers (the only  $d \in \mathbb{N}$  dividing all of the numbers is 1) has a finite subset of altogether coprime numbers.  $\square$

**Lemma 2.** *For every irreducible and aperiodic stochastic matrix  $P = (p_{i,j})$  in  $\mathbb{R}_{\geq 0}^{n \times n}$ ,  $n \in \mathbb{N}$ , there exists a number  $k \in \mathbb{N}$  such that every entry in the matrix  $P^k$  is positive, and hence every entry in  $P^k$  is at least some  $\delta > 0$ .*

**Proof.** Recall that by a proposition in the last lecture, for every  $i, j \in [n]$  the  $i, j$ -entry in  $P^k$  is the sum of weights of all  $i$ - $j$  walks in  $D_P$  with length  $k$ , where the weight of a walk in  $D_P$  is the product of weights of its edges. Thus it suffices to show that there is a  $k \in \mathbb{N}$  such that for every  $i, j \in [n]$  there is an  $i$ - $j$  walk in  $D_P$  with length  $k$ . By aperiodicity of  $P$  and the above exercise on compactness there exist cycles  $C_1, \dots, C_r$ ,  $r \in \mathbb{N}$ , in  $D_P$  with altogether coprime lengths. For each cycle  $C_l$  we fix a vertex  $v_l \in C_l$ . It is clear by irreducibility of  $D_P$  that for any given  $i, j \in [n]$  there is an  $i$ - $j$  walk  $w$  in  $D_P$  such that  $|w| \leq (r+1)(n-1)$  and  $w$  goes through all vertices  $v_1, \dots, v_r$ : we take the shortest path from  $i$  to  $v_1$ , then from  $v_1$  to  $v_2, \dots$ , and finally from  $v_r$  to  $j$ . Let  $n_0 \in \mathbb{N}$  be the number guaranteed by Lemma 1 for  $a_l := |C_l|$ ,  $l = 1, 2, \dots, r$ . We take the walk  $w$  and extend it by an appropriate detour over the cycles  $C_l$  to an  $i$ - $j$  walk  $w'$  in  $D_P$  with length

$$|w'| = k := (r+1)(n-1) + n_0.$$

$\square$

For  $\bar{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$  and  $x \in \mathbb{R}$  we introduce the notation  $|\bar{x}|_1 := \sum_i |x_i|$  ( $L_1$  norm),  $x^+ := \max(0, x)$  and  $x^- := \max(0, -x)$ . Then  $x = x^+ - x^-$ . Further,

$$\bar{x}^+ := (x_1^+, x_2^+, \dots, x_n^+)$$

and similarly for  $\bar{x}^-$ . The vectors  $\bar{x}^+$  and  $\bar{x}^-$  have disjoint supports (sets of indices of the nonzero entries),  $\bar{x} = \bar{x}^+ - \bar{x}^-$  and  $|\bar{x}|_1 = |\bar{x}^+|_1 + |\bar{x}^-|_1$ .

**Lemma 3.** *If  $\bar{v}, \bar{w} \in \mathbb{R}_{\geq 0}^n$ ,  $n \in \mathbb{N}$ , are vectors with each entry at least  $\alpha \geq 0$ , then*

$$|\bar{v} - \bar{w}|_1 \leq |\bar{v}|_1 + |\bar{w}|_1 - 2n\alpha .$$

**Proof.** It suffices to prove this inequality only for  $n = 1$ : if  $0 \leq \alpha \leq a \leq b$  then  $b - a \leq a + b - 2\alpha$ , which clearly holds as it is equivalent with  $2\alpha \leq 2a$ .  $\square$

**Lemma 4.** *Let  $\delta > 0$  be a real number,  $M$  be a matrix in  $[\delta, +\infty)^{n \times n}$ ,  $n \in \mathbb{N}$ , and with every row sum equal to 1, and let  $\bar{x} \in \mathbb{R}^n$ . Then the following hold.*

1.  $|\bar{x}M|_1 \leq |\bar{x}|_1 - 2\delta n \min(|\bar{x}^+|_1, |\bar{x}^-|_1)$ .
2. *If  $\bar{x}$  has both positive and negative entries then  $\bar{x}M \neq \bar{x}$ .*
3. *If  $\sum_i x_i = 0$  then  $|\bar{x}M|_1 \leq (1 - \delta n)|\bar{x}|_1$ .*

**Proof.** 1. We define  $\bar{v} := \bar{x}^+M$  and  $\bar{w} := \bar{x}^-M$ . Since the vectors  $\bar{x}^+$ ,  $\bar{x}^-$ ,  $\bar{v}$  and  $\bar{w}$  have nonnegative entries and the row sums of  $M$  are 1, we have that  $|\bar{v}|_1 = |\bar{x}^+|_1$  and  $|\bar{w}|_1 = |\bar{x}^-|_1$ . Hence  $|\bar{v}|_1 + |\bar{w}|_1 = |\bar{x}|_1$ . The entries in  $M$  are at least  $\delta$  and therefore we have for every  $i \in [n]$  that  $v_i \geq \delta \sum_i x_i^+ = \delta |\bar{x}^+|_1$  and, similarly,  $w_i \geq \delta |\bar{x}^-|_1$ . Applying Lemma 3 with  $\alpha := \delta \min(|\bar{x}^+|_1, |\bar{x}^-|_1)$ , we get that

$$\begin{aligned} |\bar{x}M|_1 &= |\bar{v} - \bar{w}|_1 \\ &\leq |\bar{v}|_1 + |\bar{w}|_1 - 2\delta n \min(|\bar{x}^+|_1, |\bar{x}^-|_1) \\ &= |\bar{x}|_1 - 2\delta n \min(|\bar{x}^+|_1, |\bar{x}^-|_1) . \end{aligned}$$

2. By the assumption,  $\min(|\bar{x}^+|_1, |\bar{x}^-|_1) > 0$ . Thus by part 1,  $|\bar{x}M|_1 < |\bar{x}|_1$  and  $\bar{x}M \neq \bar{x}$ .

3. By the assumption,  $|\bar{x}^+|_1 = |\bar{x}^-|_1 = |\bar{x}|_1/2$ . Thus by part 1,

$$|\bar{x}M|_1 \leq |\bar{x}|_1 - 2\delta n \min(|\bar{x}^+|_1, |\bar{x}^-|_1) = (1 - \delta n)|\bar{x}|_1 .$$

$\square$

Now we can prove the main theorem.

**Theorem (on stationary distribution).** *Any irreducible and aperiodic Markov chain  $M$  with the states  $[n]$  has a unique stationary distribution  $\bar{p}$  which is moreover attractive. If  $P$  is the transition matrix then for every  $i, j \in [n]$  one has that*

$$p_j = \lim_{n \rightarrow \infty} (P^n)_{i,j} .$$

**Proof.** Let  $k \in \mathbb{N}$  and  $\delta > 0$  be as in Lemma 2. We claim that the system

$$\bar{x}P = \bar{x}$$

of  $n$  homogeneous linear equations with  $n$  unknowns  $x_1, \dots, x_n$  has a nontrivial solution, different from  $\bar{0} = (0, \dots, 0)$ . Indeed, if we move the right sides to the left to get the canonical form  $\bar{x}(P - I) = \bar{0}$ , where  $I$  is the identity  $n \times n$  matrix, the obtained matrix  $P - I$  is singular because the sum of its columns is the zero column ( $P$  is stochastic). We denote this nontrivial solution by  $\bar{x}$ . From  $\bar{x}P = \bar{x}$  we get that also  $\bar{x}P^k = \bar{x}$  and see by part 2 of Lemma 4 and the exercise on matrices that  $\bar{x} \geq \bar{0}$  (the vector  $\bar{x}$  has only nonnegative entries) or  $\bar{x} \leq \bar{0}$ . Thus we may define the distribution

$$\bar{p} := (\sum_i x_i)^{-1} \cdot \bar{x} .$$

We show that  $\bar{p}$  is the unique stationary distribution of  $M$ . Since  $\bar{x}P = \bar{x}$ , also  $\bar{p}P = \bar{p}$ . If  $\bar{q} \neq \bar{p}$  were another distribution satisfying  $\bar{q}P = \bar{q}$ , we would have also

$$(\bar{p} - \bar{q})P^k = \bar{p} - \bar{q} .$$

But this contradicts part 2 of Lemma 4 because the vector  $\bar{p} - \bar{q}$  has both positive and negative entries.

We show that  $\bar{p}$  is attractive: for any distribution  $\bar{q}$  we have coordinatewise that

$$\lim_{n \rightarrow \infty} \bar{q}P^n = \bar{p} .$$

Then the vector  $\bar{q} = (0, \dots, 0, 1, 0, \dots, 0)$  with 1 at the  $i$ -th place yields the formula in the statement of the theorem.

Let  $\bar{a}$  be any distribution. For  $s = 0, 1, 2, \dots$  let  $\bar{v}^{(s)} := \bar{a}P^{sk} - \bar{p}$ . Stationarity of  $\bar{p}$  implies that  $\bar{v}^{(s)} = \bar{a}P^{sk} - \bar{p}P^{sk} = (\bar{a} - \bar{p})P^{sk}$  and

therefore  $\bar{v}^{(s+1)} = \bar{v}^{(s)}P^k$ . By the exercise on matrices the coordinates of each vector  $\bar{v}^{(s)}$  sum to 0. Since each entry in  $P^k$  is at least  $\delta$ , by part 3 of Lemma 4 one has that

$$|\bar{v}^{(s+1)}|_1 = |\bar{v}^{(s)}P^k|_1 \leq (1 - \delta n)|\bar{v}^{(s)}|_1 .$$

So  $|\bar{v}^{(s)}|_1 \leq (1 - \delta n)^s |\bar{v}^{(0)}|_1$  and  $|\bar{v}^{(s)}|_1 \rightarrow 0$  for  $s \rightarrow \infty$ . Here  $1 - \delta n$  is in  $[0, 1)$  because every entry in  $P^k$  is at least  $\delta > 0$  and every row sum of the  $n \times n$  matrix  $P^k$  is 1. Therefore

$$\lim_{s \rightarrow \infty} \bar{a}P^{sk} = \bar{p} .$$

We use this result with  $\bar{a} := \bar{q}P^l$  for any  $l = 0, 1, \dots, k - 1$  and get that

$$\lim_{s \rightarrow \infty} \bar{a}P^{sk} = \lim_{s \rightarrow \infty} \bar{q}P^{sk+l} = \bar{p} .$$

Thus  $\lim_{n \rightarrow \infty} \bar{q}P^n = \bar{p}$ . □

**Example (more on the Ehrenfest model).** I will write here on the Diaconis–Shakshahani result on the mixing time in the Laplace–Bernoulli diffusion. □

### Balls and bins. Bucket sort

An intuitive probabilistic model widely used in computer science is that of  $m$  balls and  $n$  bins,  $m, n \in \mathbb{N}$ , where the balls are placed randomly and mutually independently in the bins. One describes by it various situations in design of randomized algorithms. One can investigate various parameters of the (random) distribution of balls in bins.

For example, if  $m \leq n$ , or better if  $m = o(n)$ , then one can estimate for  $n \rightarrow \infty$  the probability  $B_m(n)$  that no bin contains two or more balls:

$$B_m(n) = \prod_{j=1}^{m-1} \left(1 - \frac{j}{n}\right) \approx \prod_{j=1}^{m-1} e^{-j/n} = e^{-m(m-1)/2n} \approx e^{-m^2/2n} .$$

The first equality here follows by placing the first ball in any of the  $n$  bins, the second ball in any of the remaining  $n - 1$  empty bins, the third ball in any of the remaining  $n - 2$  empty bins, and so on.

**Exercise.** Recall that for two functions  $f, g: \mathbb{N} \rightarrow \mathbb{C}$  the notation  $f \sim g$  means that  $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 1$ . Show that for  $m = o(n)$  the last  $\approx$  is indeed  $\sim$ . Show that for  $m = o(n^{2/3})$  the first  $\approx$  is  $\sim$ .  $\square$

We can interpret  $B_m(n)$  as the probability that if there are  $n$  possibilities for birthdays (usually  $n = 365$ ) then in a group of  $m$  people on a party no two of them have birthdays on the same day. (Assuming that the people and their birth dates are uncorrelated; for example, the party is not the first-day-of-the-year-baby convention.) For example, for  $m \sim \sqrt{2(\log 2)n}$  we get (by the exercise) that  $B_m(n) \sim \frac{1}{2}$ .

As a sample result (taken from Mitzenmacher and Upfal) in the balls and bins model we prove a probabilistic bound on the maximum load of a bin.

**Proposition (maximum load).** *Let  $m = n$  and  $X_{\max}$  be the random variable recording the maximum number of balls in a bin (“maximum load”). Then*

$$\Pr(X_{\max} > 3 \frac{\log n}{\log \log n}) \leq \frac{1}{n}, \quad n \geq n_0 .$$

**Proof.** For any  $M \in \mathbb{N}$  we get by the union bound that

$$\Pr(\geq M \text{ balls in bin } 1) \leq \binom{n}{M} \left(\frac{1}{n}\right)^M \leq \frac{1}{M!} \leq \left(\frac{e}{M}\right)^M .$$

Thus, again by the union bound, for  $M \geq \frac{3 \log n}{\log \log n}$  we have that

$$\begin{aligned} \Pr(\geq M \text{ balls in any bin}) &\leq n \left(\frac{e \log \log n}{3 \log n}\right)^{3 \log n / \log \log n} \\ &\leq n \left(\frac{\log_{(2)} n}{\log n}\right)^{3 \log n / \log_{(2)} n} \\ &= \exp(\log n) \exp\left((\log_{(3)} n - \log_{(2)} n) \frac{3 \log n}{\log_{(2)} n}\right) \\ &= \exp\left(-2 \log n + \frac{3 \log n \log_{(3)} n}{\log_{(2)} n}\right) \\ &\leq \frac{1}{n}, \quad n \geq n_0 . \end{aligned}$$

$\square$



**Bucket sort.** We sort  $n = 2^m$  elements which are random integers from  $[0, 2^k)$  where  $k \geq m$ . Using bucket sort we sort them in expected time  $O(n)$ . The algorithm is deterministic but the input is random.

*First stage.* We place the elements in  $n$  buckets  $b_1, \dots, b_n$  where we put in  $b_j$  the elements whose first  $m$  binary digits give number  $j$ . Thus  $b_1 < b_2 < \dots < b_n$  meaning that for  $k < l$  any element in  $b_k$  is smaller than any element in  $b_l$ . We assume that each element can be placed in its bucket in  $O(1)$  time. This stage therefore takes  $O(n)$  time. Note that the random variable  $X$  recording the number of elements in a fixed bucket has the binomial distribution  $B(n, 1/n)$  (see Lecture 9 for its definition).

*Second stage.* The elements in every bucket are sorted by some algorithm, for example bubblesort, in quadratic time. We estimate the expected running time of the second stage. If  $X_j$  is the random variable recording the number of elements in  $b_j$ , they are sorted in time  $\leq cX_j^2$  for some constant  $c > 0$ . Thus the expected running time of the second stage is at most

$$\mathbb{E}\left(\sum_{j=1}^n cX_j^2\right) = c \sum_{j=1}^n \mathbb{E}X_j^2 = cn\mathbb{E}X_1^2 = cn(2 - 1/n) = O(n) .$$

We used linearity of expectation and the formula for the second moment of a binomial distribution, computed in Lecture 9. In total, both stages together, bucket sort runs in time  $O(n)$ .

**More on bucket sort.** Time permitting, I write here what DEK writes on it in *TAOCP, Volume 3 (Sorting and Searching)*.

**Thank you!**

(semifinal version — more information on the Ehrenfest model and bucket sort will be added — of January 13, 2021)