

“Balls into Bins” — A Simple and Tight Analysis

MARTIN RAAB and ANGELIKA STEGER

Institut für Informatik
Technische Universität München
D-80290 München
{raab|steger}@informatik.tu-muenchen.de

Abstract. Suppose we sequentially throw m balls into n bins. It is a natural question to ask for the maximum number of balls in any bin. In this paper we shall derive sharp upper and lower bounds which are reached with high probability. We prove bounds for all values of $m(n) \geq n/\text{polylog}(n)$ by using the simple and well-known method of the first and second moment.

1 Introduction

Suppose that we sequentially throw m balls into n bins by placing each ball into a bin chosen independently and uniformly at random. It is a natural question to ask for the maximum number of balls in any bin. This very simple model has many applications in computer science. Here we name just two of them:

Hashing: The balls-into-bins model may be used to analyze the efficiency of hashing-algorithms. In the case of the so called *separate chaining*, all keys that hash to the same location in the table are stored in a linked list. It is clear that the lengths of these lists are a measure for the complexity. For a well chosen hash-function (*i.e.* a hash-function which assigns the keys to all locations in the table with the same probability), the lengths of the lists have exactly the same distribution as the number of balls in a bin.

Online Load Balancing: With the growing importance of parallel and distributed computing the load balancing problem has gained considerable attention during the last decade. A typical application for online load balancing is the following scenario: consider n database-servers and m requests which arise independently at different clients and which may be handled by any server. The problem is to assign the requests to the servers in such a way that all servers handle (about) the same number of requests.

Of course, by introducing a central dispatcher one can easily achieve uniform load on the servers. However, within a distributed setting the use of such a central dispatcher is highly undesired. Instead randomized strategies have been applied very successfully for the development of good and efficient load balancing

algorithms. In their simplest version each request is assigned to a server chosen independently and uniformly at random. If all requests are of the same size the maximum load of a server then corresponds exactly to the maximum number of balls in a bin in the balls-into-bins model introduced above.

1.1 Previous results

Balls and bins games have been intensively studied in the literature, cf. e.g. [JK77]. The estimation of the maximum number of balls in any bin was originally mainly studied within the context of hashing functions. In particular, GONNET [Gon81] determined for the case $m = n$ that the expected number of balls in the bin containing the maximum number of balls is $\Gamma^{-1}(n)(1 + O(\frac{1}{\log \Gamma^{-1}(n)}))$. One can check that Gonnet's result implies that the maximum load of any bin is with high probability $\frac{\log n}{\log \log n}(1 + o(1))$. In his dissertation MITZENMACHER [Mit96] also included a simpler proof of the fact that the maximum load is $\Theta(\frac{\log n}{\log \log n})$. He also obtains some results for the case $m < n/\log n$. For the case $m \geq n \log n$ it was well known that the maximum load of any bin is $\Theta(\frac{m}{n})$, *i.e.* of the order of the mean. However, the precise deviation from the mean seems not to have been studied before.

We note that for the online load balancing also different models of balls into bin games have been studied. We note in particular the approach of AZAR et al. [ABKU92]. They study the following model: each ball picks d bins uniformly at random and places itself in those bin containing fewest balls. For the case $m = n$ [ABKU92] showed that in this model the maximum load of any bin drops exponentially from $\frac{\log n}{\log \log n}(1 + o(1))$ to $\frac{\log \log n}{\log d}(1 + o(1))$. Compare also CZUMAJ and STEMANN [CS97] for more results in this direction.

1.2 Our results

In this paper we apply the first and second moment method, a well-known tool within the theory of random graphs, cf. *e.g.* [Bol85], to obtain a straightforward proof of the fact that the maximum number of balls in a bin is $\frac{\log n}{\log \log n}(1 + o(1))$ for $m = n$ with probability $1 - o(1)$.

Besides being a lot more elementary than GONNET's proof method the big advantage of our method is that it also easily generalizes to the case where $m \neq n$ balls are placed into n bins. In particular, this allows to also analyze the case $m \gg n$, which can neither be handled by GONNET's approach nor by MITZENMACHER's. (Both are based on approximating the Binomial distribution $B(m, \frac{1}{n})$ by a Poisson distribution, which only gives tight bounds if $m \cdot \frac{1}{n}$ is a constant.) The case $m \gg n$ is particularly important for the load-balancing scenario mentioned above. Here it *e.g.* measures how the unsymmetry between different servers grows over time when more and more requests arrive.

Our results are summarized in the following theorem:

Theorem 1. *Let M be the random variable that counts the maximum number of balls in any bin, if we throw m balls independently and uniformly at random into n bins. Then $\Pr[M > k_\alpha] = o(1)$ if $\alpha > 1$ and $\Pr[M > k_\alpha] = 1 - o(1)$ if $0 < \alpha < 1$, where*

$$k_\alpha = \begin{cases} \frac{\log n}{\log \frac{n \log n}{m}} \left(1 + \alpha \frac{\log^{(2)} \frac{n \log n}{m}}{\log \frac{n \log n}{m}} \right), & \text{if } \frac{n}{\text{polylog}(n)} \leq m \ll n \log n, \\ (d_c - 1 + \alpha) \log n, & \text{if } m = c \cdot n \log n \text{ for some constant } c, \\ \frac{m}{n} + \alpha \sqrt{2 \frac{m}{n} \log n}, & \text{if } n \log n \ll m \leq n \cdot \text{polylog}(n), \\ \frac{m}{n} + \sqrt{\frac{2m \log n}{n} \left(1 - \frac{1}{\alpha} \frac{\log^{(2)} n}{2 \log n} \right)}, & \text{if } m \gg n \cdot (\log n)^3. \end{cases}$$

Here d_c denotes a suitable constant depending only on c , cf. the proof of Lemma 3.

The paper is organized as follows: in § 2 we give a brief overview of the first and second moment method, in § 3 we show how to apply this method within the balls-into-bins scenario and obtain in § 4 the $\frac{\log n}{\log \log n}(1 + o(1))$ bound for $m = n$. In § 5 we then present some more general tail bounds for Binomial random variables and combine them with the first and second moment method to obtain a proof of Theorem 1.

1.3 Notations

Throughout this paper m denotes the number of balls and n the number of bins. The probability that a ball is thrown into a fixed bin is given by $p := 1/n$. We define q by $q := 1 - p$. We shall denote the iterated log by $\log^{(\cdot)}$, i.e. $\log^{(1)} x = \log x$ and $\log^{(k+1)} x = \log(\log^{(k)} x)$ for all $k \geq 1$. In this paper logarithms are to the base e .

Asymptotic notations ($O(\cdot)$, $o(\cdot)$ and $\omega(\cdot)$) are always with respect to n ; $f \ll g$ means $f = o(g)$ and $f \gg g$ means $f = \omega(g)$. We use the term $\text{polylog}(x)$ to denote the class of functions $\bigcup_{k \geq 1} O((\log x)^k)$. We say that an event \mathcal{E} occurs with high probability if $\Pr[\mathcal{E}] = 1 - o(1)$.

2 The first and second moment method

Let X be a non-negative random variable. Then MARKOV's inequality implies that $\Pr[X \geq 1] \leq \mathbb{E}[X]$. Hence, we have

$$\mathbb{E}[X] = o(1) \quad \implies \quad \Pr[X = 0] = 1 - o(1). \quad (1)$$

Furthermore, CHEBYSHEV's inequality implies that

$$\Pr[X = 0] \leq \Pr[|X - \mathbb{E}[X]| \geq \mathbb{E}[X]] \leq \frac{\text{Var}[X]}{(\mathbb{E}[X])^2} = \frac{\mathbb{E}[X^2]}{(\mathbb{E}[X])^2} - 1.$$

Hence, in order to show that $\Pr[X = 0] = o(1)$ we just have to verify that

$$\mathbb{E}[X^2] = (1 + o(1))(\mathbb{E}[X])^2. \quad (2)$$

While it is often quite tedious to verify (2), it is relatively easy if we can write X as the sum of (not necessarily independent) 0-1 variables X_1, \dots, X_n that satisfy

$$\mathbb{E}[X_j] = \mathbb{E}[X_1] \quad \text{and} \quad \mathbb{E}[X_i X_j] \leq (1 + o(1))(\mathbb{E}[X_1])^2 \quad \forall 1 \leq i < j \leq n. \quad (3)$$

Then

$$\mathbb{E}[X^2] = \mathbb{E}\left[\left(\sum_{i=1}^n X_i\right)^2\right] = \mathbb{E}\left[\sum_i \underbrace{X_i^2}_{=X_i} + \sum_{i \neq j} X_i X_j\right] \leq \mathbb{E}[X] + (1 + o(1))(\mathbb{E}[X])^2$$

and we can combine (1) and (2) to obtain

$$\Pr[X = 0] = \begin{cases} 1 - o(1), & \text{if } \mathbb{E}[X] = o(1), \\ o(1), & \text{if } \mathbb{E}[X] \rightarrow \infty. \end{cases} \quad (4)$$

This is the form which we will use for the analysis of the balls and bins scenario.

3 Setup for the analysis

Let $Y_i = Y_i(m, n)$ be the random variable which counts the number of balls in the i th bin if we throw m balls independently and uniformly at random into n bins. Clearly, Y_i is a binomially distributed random variable: we express this fact by writing $Y_i \sim B(m, 1/n)$ respectively $\Pr[Y_i = k] = b(k; m, 1/n) := \binom{m}{k} \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{m-k}$. Let $X_i = X_i(m, n, \alpha)$ be the random variable, which indicates if Y_i is at least $k_\alpha = k_\alpha(m, n)$ (the function from Theorem 1) and let $X = X(m, n, \alpha)$ be the sum over all X_i 's, *i.e.*:

$$X := \sum_{i=1}^n X_i \quad \text{and} \quad X_i := \begin{cases} 1, & \text{if } Y_i \geq k_\alpha, \\ 0, & \text{otherwise.} \end{cases}$$

Clearly,

$$\mathbb{E}[X_i] = \Pr[B(m, 1/n) \geq k_\alpha] \quad \text{for all } i = 1, \dots, n,$$

and

$$\mathbb{E}[X] = n \cdot \Pr[B(m, 1/n) \geq k_\alpha]. \quad (5)$$

In order to apply (4) we therefore need good bounds for the tail of the binomial distribution. Before we obtain these for the general case of all $m = m(n)$ we consider the special case $m = n$.

4 The case $m = n$

The aim of this section is to present a self contained proof of the fact that if $m = n$ the maximum number of balls in a bin is $\frac{\log n}{\log^{(2)} n}(1 + o(1))$ with high probability. We will do this by showing that

$$\Pr \left[\exists \text{ at least one bin with } \geq \alpha \frac{\log n}{\log^{(2)} n} \text{ balls} \right] = \begin{cases} 1 - o(1), & \text{if } 0 < \alpha < 1, \\ o(1), & \text{if } \alpha > 1. \end{cases} \quad (6)$$

Note that the claim of equation (6) is slightly weaker than the corresponding one from Theorem 1. We consider this case first, as here the calculations stay slightly simpler. So, in the rest of this section we let $k_\alpha := \alpha \frac{\log n}{\log^{(2)} n}$.

Recall from § 2 that in order to do this we only have to show that condition (3) is satisfied for the random variables X_i introduced in the previous section and that

$$\mathbb{E}[X] = n \cdot \Pr \left[B(n, \frac{1}{n}) \geq \alpha \frac{\log n}{\log^{(2)} n} \right] \rightarrow \begin{cases} \infty, & \text{if } 0 < \alpha < 1, \\ 0, & \text{if } \alpha > 1. \end{cases} \quad (7)$$

The fact that $\mathbb{E}[X_i] = \mathbb{E}[X_1]$ for all $1 \leq i \leq n$ follows immediately from the definition of the X_i 's. The proof of the second part of (3) is deferred to the end of this section. Instead we start with the verification of (7). For that we prove a small lemma on the binomial distribution. We state it in a slightly more general form then necessary, as this version will be helpful later-on.

Lemma 1. *Let $p = p(m)$ depend on m . Then for all $h \geq 1$*

$$\Pr [B(m, p) \geq mp + h] = \left(1 + O\left(\frac{mp}{h}\right)\right) \cdot b(mp + h; m, p).$$

Proof. Observe that for all $k \geq mp + h$:

$$\frac{b(k + 1; m, p)}{b(k; m, p)} = \frac{(m - k)p}{(k + 1)(1 - p)} \leq \frac{((1 - p)m - h)p}{(mp + h + 1)(1 - p)} =: \lambda.$$

One easily checks that $\lambda < 1$ for $h \geq 1$. Thus

$$\sum_{k \geq mp + h} b(k; m, p) \leq b(mp + h; m, p) \cdot \sum_{i \geq 0} \lambda^i = \frac{1}{1 - \lambda} \cdot b(mp + h; m, p).$$

As $\frac{1}{1 - \lambda} \leq 1 + \frac{mp}{h}$ the claim of the lemma follows. \square

We apply Lemma 1 for “ m ” = n , “ p ” = $\frac{1}{n}$ and “ $mp + h$ ” = k_α . Subsequently, we use STIRLING’s formula $x! = (1 + o(1))\sqrt{2\pi x}e^{-x}x^x$ to estimate the binomial coefficient. Together we obtain:

$$\begin{aligned} \mathbb{E}[X] &= n \cdot \Pr \left[B(n, \frac{1}{n}) \geq k_\alpha \right] = n \cdot (1 + o(1)) \cdot b(k_\alpha; n, \frac{1}{n}) \\ &= n \cdot (1 + o(1)) \binom{n}{k_\alpha} \left(\frac{1}{n}\right)^{k_\alpha} \left(1 - \frac{1}{n}\right)^{n - k_\alpha} \end{aligned}$$

$$\begin{aligned}
&= n \cdot (1 + o(1)) \frac{1}{e\sqrt{2\pi k_\alpha}} \left(\frac{e}{k_\alpha}\right)^{k_\alpha} \\
&= n \cdot e^{\alpha \frac{\log n}{\log^{(2)} n} \cdot (1 - \log \alpha - \log^{(2)} n + \log^{(3)} n + o(1))} \\
&= n^{1 - \alpha + o(1)},
\end{aligned}$$

which implies the statement of equation (7).

To complete the proof for the case $m = n$ we still have to verify that $\mathbb{E}[X_i X_j] \leq (1 + o(1))(\mathbb{E}[X_1])^2$ for all $i \neq j$. In order to keep the proof elementary we shall proceed similarly as in the proof of Lemma 1. A more elegant version can be found in § 5.

$$\begin{aligned}
\mathbb{E}[X_i X_j] &= \Pr[Y_i \geq k_\alpha \wedge Y_j \geq k_\alpha] \\
&= \sum_{k_1=k_\alpha}^{n-k_\alpha} \sum_{k_2=k_\alpha}^{n-k_1} \binom{n}{k_1} \underbrace{\binom{n-k_1}{k_2}}_{\leq \binom{n}{k_2}} \left(\frac{1}{n}\right)^{k_1+k_2} \underbrace{\left(1 - \frac{2}{n}\right)^{n-(k_1+k_2)}}_{\leq \left(1 - \frac{1}{n}\right)^2} \\
&\leq \sum_{k_1=k_\alpha}^n \sum_{k_2=k_\alpha}^n \binom{n}{k_1} \binom{n}{k_2} \left(\frac{1}{n}\right)^{k_1+k_2} \left(1 - \frac{1}{n}\right)^{2n-2(k_1+k_2)} \\
&\leq \left[b(k_\alpha; n, \frac{1}{n}) \left(1 - \frac{1}{n}\right)^{-k_\alpha} \sum_{i=0}^{\infty} \lambda^i \right]^2
\end{aligned}$$

where λ is defined as $\lambda := \frac{b(k_\alpha+1; n, \frac{1}{n})(1 - \frac{1}{n})^{k_\alpha}}{b(k_\alpha; n, \frac{1}{n})(1 - \frac{1}{n})^{k_\alpha+1}}$. As $\lambda = o(1)$ and $b(k_\alpha; n, \frac{1}{n}) = (1 + o(1))\mathbb{E}[X_1]$ (cf. Lemma 1) this concludes the proof of (6).

5 The general case

For the proof of Theorem 1 we will follow the same pattern as in the proof of the previous section. The main difference is that in various parts we need better bounds. We start by collecting some bounds on the tails of the binomial distribution.

5.1 Tails of the binomial distribution

The binomial distribution is very well studied. In particular it is well-known that the binomial distribution $B(m, p)$ tends to the normal distribution if $0 < p < 1$ is a fixed constant and m tends to infinity. If on the other hand $p = p(m)$ depends on m in such a way that mp converges to a constant λ for m tending to infinity, then the corresponding binomial distribution $B(m, p)$ tends to the Poisson distribution with parameter λ . For these two extreme cases also very good bounds on the tail of the binomial distributions are known. In the context of our “balls and bins” scenario, however, we are interested in the whole spectrum of values $p = p(m)$.

In this section we collect some bounds on the tails of the binomial distribution which are tailored to the proof of Theorem 1.

For values of $p(m)$ such that mp tends to infinity one can analyze the proof of the theorem of DEMOIVRE–LAPLACE to get asymptotic formulas for $\Pr [B(m, p) \geq mp + h]$ for all values h that are not “too“ large:

Theorem 2 (DeMoivre–Laplace). *Assume $0 < p < 1$ depends on n such that $pqm = p(1 - p)m \rightarrow \infty$ for $m \rightarrow \infty$. If $0 < h = x(pqm)^{1/2} = o((pqm)^{2/3})$ and $x \rightarrow \infty$ then*

$$\Pr [B(m, p) \geq mp + h] = (1 + o(1)) \cdot \frac{1}{x\sqrt{2\pi}} e^{-\frac{x^2}{2}}.$$

For an explicit proof of this version of the DEMOIVRE–LAPLACE Theorem see *e.g.* [Bol85].

The probability that a binomial distributed random variable $B(m, p)$ obtains a value of size at least $mp(1 + \epsilon)$ for some constant $\epsilon > 0$ is usually estimated using the so-called CHERNOFF bounds. Recall, however, that CHERNOFF bounds provide only an upper bound.

With the help of Lemma 1 in the previous section we are now in the position to prove the tail bounds for those special cases of the binomial distribution which we will need further-on.

Lemma 2. *a) If $mp + 1 \leq t \leq (\log m)^\ell$, for some positive constant ℓ , then*

$$\Pr [B(m, p) \geq t] = e^{t(\log mp - \log t + 1) - mp + O(\log^{(2)} m)}.$$

b) If $t = mp + o\left((pqm)^{\frac{2}{3}}\right)$ and $x := \frac{t - mp}{\sqrt{pqm}}$ tends to infinity, then

$$\Pr [B(m, p) \geq t] = e^{-\frac{x^2}{2} - \log x - \frac{1}{2} \log 2\pi + o(1)}.$$

Proof. a) Using STIRLING’s formula $x! = (1 + o(1))\sqrt{2\pi x}e^{-x}x^x$ we obtain:

$$b(t; m, p) = (1 + o(1)) \frac{1}{\sqrt{2\pi t}} \left(\frac{mp}{t}\right)^t \left(1 + \frac{t - mp}{m - t}\right)^{m-t}.$$

Together with Lemma 1 we thus get for $\log \Pr [B(m, p) \geq t]$ the following expression:

$$\log \left(1 + O\left(\frac{mp}{t - mp}\right)\right) + t(\log mp - \log t + 1) - mp - \frac{\log t}{2} + O\left(\frac{(t - mp)^2}{m - t}\right) - O(1)$$

The term $O\left(\frac{(t - mp)^2}{m - t}\right)$ gets arbitrarily small if $mp \leq t = o(\sqrt{m})$ because $\frac{(t - mp)^2}{m - t} \leq \frac{t^2}{m(1 - o(1))} = o(1)$. By assumption $mp + 1 \leq (\log m)^\ell$. That is, $\log \left(1 + O\left(\frac{mp}{t - mp}\right)\right) = O(\log \log m)$.

b) This case is simply a reformulation of the DEMOIVRE–LAPLACE theorem. \square

5.2 Proof of Theorem 1

We follow the setup outlined in § 2. That is, we only have to show that the variables X_i satisfy condition (3) and that the expectation of $X = \sum X_i$ tends either to infinity or to zero depending on the fact whether α is smaller or greater than 1. We start with the later.

Lemma 3. *Let k_α be defined as in Theorem 1. Then*

$$\log \mathbb{E}[X] \rightarrow \begin{cases} \infty, & \text{if } 0 < \alpha < 1, \\ -\infty, & \text{if } \alpha > 1, \end{cases}$$

for all values $m = m(n) \geq n/\text{polylog}(n)$.

Proof. The case $\frac{n}{\text{polylog}(n)} \leq m \ll n \log n$.

We first note that it suffices to consider non-negative α 's. Assume that $m = \frac{n \log n}{g}$ where $g = g(n)$ tends to infinity arbitrarily slowly and $g(n) \leq \text{polylog}(n)$. Then

$$k_\alpha = \frac{\log n}{\log g} \left(1 + \alpha \frac{\log^{(2)} g}{\log g} \right).$$

From equation (5) and Lemma 2 case a) (we leave it to the reader to verify that this case may be applied) it follows that

$$\begin{aligned} \log \mathbb{E}[X] &= \log n + k_\alpha \left(\log^{(2)} n - \log g - \log k_\alpha + 1 \right) - \frac{\log n}{g} + O\left(\log^{(2)} m\right) \\ &= \frac{\log n}{\log g} \left(\log g + \left(1 + \alpha \frac{\log^{(2)} g}{\log g} \right) \left(1 - \log g + \log^{(2)} g \right) + o(1) \right) \\ &= (1 - \alpha + o(1)) \frac{\log n \cdot \log^{(2)} g}{\log g}, \end{aligned}$$

which yields the desired result.

The case $m = c \cdot n \log n$.

Let $k_\alpha := (d_c - 1 + \alpha)$. By Lemma 2 we get:

$$\log \mathbb{E}[X] = \log n (1 + (d_c - 1 + \alpha) (\log c - \log (d_c - 1 + \alpha) + 1) - c + o(1)).$$

As a consequence, for $\alpha = 1$ $\log \mathbb{E}[X]$ is exactly then $o(\log n)$ when d_c is a solution of

$$f_c(x) := 1 + x (\log c - \log x + 1) - c = 0.$$

For all $c > 0$ this equation admits exactly two real zeros x_1, x_2 . One of these solutions is smaller than c and is therefore not the one we are looking for. That is, we define d_c as the (unique) solution of $f_c(x) = 0$ that is greater than c . In the neighborhood of the solutions x_1 and x_2 , $f_c(x)$ changes its sign. This means that for $d_c - 1 + \alpha$ for a given $\alpha > 1$, $\log \mathbb{E}[X]$ tends to $-\infty$, whereas for $d_c - 1 + \alpha$ for an $0 < \alpha < 1$, $\log \mathbb{E}[X]$ tends to ∞ .

The case $n \log n \ll m \leq n \cdot \text{polylog}(n)$.

Assume that $m = gn \log n$ where $g = g(n) \leq \text{polylog}(n)$ tends to infinity arbitrarily slowly. Then

$$k_\alpha = g \log n \left(1 + \alpha \sqrt{\frac{2}{g}} \right).$$

From Lemma 2 case a) it follows that

$$\begin{aligned} \log \mathbb{E}[X] &= \log n + k_\alpha \left(\log g + \log^{(2)} n - \log k_\alpha + 1 \right) - g \log n + O\left(\log^{(2)} n\right) \\ &= g \log n \left(\frac{1}{g} + \left(1 + \alpha \sqrt{\frac{2}{g}} \right) \left(1 - \alpha \sqrt{\frac{2}{g}} + \frac{\alpha^2}{g} + o\left(\frac{1}{g}\right) \right) - 1 + o\left(\frac{1}{g}\right) \right) \\ &= \log n \left(1 - \alpha^2 + o(1) \right). \end{aligned}$$

One easily checks, that we didn't hurt the conditions of Lemma 2.

The case $m \gg n(\log n)^3$.

For this case we shall use the theorem of DEMOIVRE-LAPLACE. Recall that in this case

$$k_\alpha = \frac{m}{n} + \sqrt{\frac{2m \log n}{n} \left(1 - \frac{1}{\alpha} \frac{\log^{(2)} n}{2 \log n} \right)}.$$

Using the notations of Lemma 2 case b) we set

$$x := \frac{k_\alpha - mp}{\sqrt{pqm}} = \sqrt{2 \log n \left(1 - \frac{1}{2\alpha} \frac{\log^{(2)} n}{\log n} \right) \left(1 + \frac{1}{n-1} \right)}.$$

Applying DEMOIVRE-LAPLACE we obtain:

$$\begin{aligned} \log \mathbb{E}[X] &= \log n - \frac{x^2}{2} - \log x - \log \sqrt{2\pi} + o(1) \\ &= \log^{(2)} n \cdot \left(\frac{1}{2\alpha} - \frac{1}{2} + o(1) \right). \end{aligned}$$

We still need to check that we didn't violate the conditions of DEMOIVRE-LAPLACE, i.e. that $k_\alpha - \frac{m}{n} = o\left((pqm)^{\frac{2}{3}}\right)$, but this is true if $\frac{m}{n \log n} = \omega(\log^2 n)$.
□

In order to show that the variables X_i satisfy the second part of condition (3) (note that the first part is trivially true) we start with two simple lemmas.

Lemma 4. *Let $p \leq \frac{1}{4}$ and m be such that $p^2 m = o(1)$. Then*

$$\Pr \left[B(m(1-p), \frac{p}{1-p}) \geq t \right] \leq (1 + o(1)) \cdot \Pr [B(m, p) \geq t] \quad \text{for all } 0 \leq t \leq m.$$

Proof. We will show that for all $1 \leq t \leq m$ we have $b(t; m(1-p), \frac{p}{1-p}) \leq (1 + o(1))b(t; m, p)$. Clearly, this then completes the proof of the lemma. So consider an arbitrary, but fixed $1 \leq t \leq m$. For $t > m(1-p)$ we have $b(t; m(1-p), \frac{p}{1-p}) = 0$ so we might as well assume that $t \leq m(1-p)$. Then

$$\begin{aligned}
b(t; m(1-p), \frac{p}{1-p}) &= \binom{m(1-p)}{t} \cdot \left(\frac{p}{1-p}\right)^t \cdot \left(1 - \frac{p}{1-p}\right)^{m(1-p)-t} \\
&= \binom{m}{t} \cdot \prod_{i=0}^{t-1} \underbrace{\frac{m(1-p)-i}{m-i}}_{\leq 1-p} \cdot \left(\frac{p}{1-p}\right)^t \cdot \underbrace{\left(1 - \frac{p}{1-p}\right)^{m(1-p)}}_{\leq (1-p)^m} \cdot \left(1 - \frac{p}{1-p}\right)^{-t} \\
&\leq \binom{m}{t} \cdot p^t \cdot (1-p)^{m-t} \cdot \left(\frac{1-p}{1-\frac{p}{1-p}}\right)^t \\
&= b(t; m, p) \cdot \underbrace{\left(1 + \frac{p^2}{1-2p}\right)^t}_{\leq e^{2p^2m}} \\
&= b(t; m, p) \cdot (1 + o(1)).
\end{aligned}$$

□

Lemma 5. Let $p = o(1)$ and m, t be such that $x := \frac{t-mp}{\sqrt{mp(1-p)}}$ satisfies $x \rightarrow \infty$, $x = o((mp(1-p))^{1/6})$ and $xp = o(1)$. Then

$$\Pr \left[B(m(1-p), \frac{p}{1-p}) \geq t \right] \leq (1 + o(1)) \Pr [B(m, p) \geq t].$$

Proof. Observe that the assumptions of the lemma are such that we may apply case b) of Lemma 2 to compute $\Pr [B(m, p) \geq t]$. Observe furthermore that we may also apply this case of Lemma 2 to bound $\Pr \left[B(m(1-p), \frac{p}{1-p}) \geq t \right]$, as here the corresponding x -value is

$$\bar{x} = \frac{t - m(1-p) \cdot \frac{p}{1-p}}{\sqrt{m(1-p) \cdot \frac{p}{1-p} \cdot \left(1 - \frac{p}{1-p}\right)}} = \frac{t - mp}{\sqrt{mp(1-p)}} \cdot \sqrt{\frac{1-p}{1-\frac{p}{1-p}}} = x \cdot \sqrt{1 + \frac{p^2}{1-2p}}.$$

Together we deduce

$$\begin{aligned}
\Pr \left[B(m(1-p), \frac{p}{1-p}) \geq t \right] &= e^{-\frac{\bar{x}^2}{2} (1 + \frac{p^2}{1-2p}) - \log x - \frac{1}{2} \log(1 + \frac{p^2}{1-2p}) - \frac{1}{2} \log 2\pi + o(1)} \\
&= \Pr [B(m, p) \geq t] \cdot e^{-O(p^2 x^2) - O(p^2) + o(1)} \\
&= \Pr [B(m, p) \geq t] \cdot (1 + o(1)).
\end{aligned}$$

□

Corollary 1. Let $m = m(n)$ and $p = \frac{1}{n}$ be such that $m \geq \log n$, and let k_α denote the value from Theorem 1. Then

$$\Pr \left[B(m(1-p), \frac{p}{1-p}) \geq k_\alpha \right] \leq (1 + o(1)) \cdot \Pr [B(m, p) \geq k_\alpha].$$

Proof. One easily checks that for all $m = o(n^2)$ Lemma 4 applies and that for all $m \gg n(\log n)^3$ Lemma 5 applies. \square

Lemma 6. Let X_1, \dots, X_n be defined as in § 3. Then for all $1 \leq i < j \leq n$

$$\mathbb{E}[X_i X_j] \leq (1 + o(1)) \cdot (\mathbb{E}[X_1])^2.$$

Proof. Using the notation from § 3 we have

$$\begin{aligned} \mathbb{E}[X_i X_j] &= \Pr [Y_i \geq k_\alpha \wedge Y_j \geq k_\alpha] \\ &= \sum_{\substack{k_1+k_2 \leq m \\ k_1, k_2 \geq k_\alpha}} \binom{m}{k_1} \binom{m-k_1}{k_2} p^{k_1+k_2} (1-2p)^{m-k_1-k_2} \\ &= \sum_{k_1=k_\alpha}^{m-k_\alpha} \binom{m}{k_1} p^{k_1} (1-p)^{m-k_1} \cdot \sum_{k_2=k_\alpha}^{m-k_1} \binom{m-k_1}{k_2} \left(\frac{p}{1-p}\right)^{k_2} \left(1 - \frac{p}{1-p}\right)^{m-k_1-k_2} \\ &= \sum_{k_1=k_\alpha}^{m-k_\alpha} \binom{m}{k_1} p^{k_1} (1-p)^{m-k_1} \cdot \Pr \left[B(m-k_1, \frac{p}{1-p}) \geq k_\alpha \right]. \end{aligned}$$

As $k_1 \geq k_\alpha \geq mp$ we observe that

$$\begin{aligned} \Pr \left[B(m-k_1, \frac{p}{1-p}) \geq k_\alpha \right] &\leq \Pr \left[B(m(1-p), \frac{p}{1-p}) \geq k_\alpha \right] \\ &= (1 + o(1)) \cdot \Pr [B(m, p) \geq k_\alpha], \end{aligned}$$

where the last equality follows from Corollary 1. Hence,

$$\begin{aligned} \mathbb{E}[X_i X_j] &= \sum_{k_1=k_\alpha}^{m-k_\alpha} \binom{m}{k_1} p^{k_1} (1-p)^{m-k_1} \cdot (1 + o(1)) \cdot \Pr [B(m, p) \geq k_\alpha] \\ &= (1 + o(1)) \cdot \Pr [B(m, p) \geq k_\alpha] \cdot \underbrace{\sum_{k_1=k_\alpha}^{m-k_\alpha} \binom{m}{k_1} p^{k_1} (1-p)^{m-k_1}}_{\leq \Pr [B(m, p) \geq k_\alpha]} \\ &\leq (1 + o(1)) \cdot (\Pr [B(m, p) \geq k_\alpha])^2. \end{aligned}$$

As $\Pr [B(m, p) \geq k_\alpha] = \mathbb{E}[X_1]$, this completes the proof of the lemma. \square

6 Conclusion

In this paper we derived an asymptotic formula for the maximum number of balls in any bin, if $m = m(n)$ balls are thrown randomly into n bins, for all values of $m \geq n/\text{polylog}(n)$. Our proof is based on the so-called first and second moment.

The result for $m = n$ was well-known before. However, our method gave a much simpler proof compared to those which were previously available in the literature. To the best of our knowledge the result for the case $m \gg n$ is new. In our opinion it is a challenging open problem to study the behavior of the modified balls into bins scenario as introduced in [ABKU92] for the case $m \gg n$ as well. Intensive computational experiments seem to indicate that in this case the difference of the maximum load in any bin from the mean m/n should be independent of m . We intend to settle this problem in a forthcoming paper.

References

- [ABKU92] Y. Azar, A.Z. Broder, A.R. Karlin, and E. Upfal. On-line load balancing (extended abstract). In *33rd Annual Symposium on Foundations of Computer Science*, pages 218–225, Pittsburgh, Pennsylvania, 24–27 October 1992. IEEE.
- [Bol85] B. Bollobás. *Random graphs*. Academic Press, New York-San Francisco-London-San Diego, 1985.
- [CS97] A. Czumaj and V. Stemann. Randomized allocation processes. In *38th Annual Symposium on Foundations of Computer Science*, pages 194–203, 1997.
- [Gon81] G.H. Gonnet. Expected length of the longest probe sequence in hash code searching. *J. ACM*, 28(2):289–304, 1981.
- [JK77] N. Johnson and S. Kotz. *Urn Models and Their Applications*. John Wiley and Sons, 1977.
- [Mit96] M.D. Mitzenmacher. *The Power of Two Choices in Randomized Load Balancing*. PhD thesis, Computer Science Department, University of California at Berkeley, 1996.