

436 A Proof of Theorem 3.3

Theorem 3.3. *Let \mathcal{A} be a proper procedure for testing property \mathcal{P} as defined in Definition 3.2. Suppose the expected number of test samples, s , is bounded from below:*

$$s \geq \Theta \left(\frac{n' \cdot \sqrt{\mathbf{E}_F[\min(\|p^{(F)}\|_2^2, \|q^{(F)}\|_2^2)]}}{\epsilon^2} + \frac{\sqrt{n' \Delta(\bar{Z})}}{\epsilon \sqrt{\xi}} \right).$$

437 Then Algorithm 1 is ξ -differentially private $(\epsilon, 3/4)$ -tester for testing property \mathcal{P} .

438 **Proof:** First, we prove that the algorithm is ξ -differentially private. Note that the statistic we use is
 439 \hat{Z} which is equal to \bar{Z} , as defined in Equation 2, plus a Laplace random variable with mean $\Delta(\bar{Z})/\xi$.
 440 According to the Laplace mechanism \hat{Z} is a ξ -differentially private quantity, so the output of the
 441 algorithm is private.

442 Now, we prove that the algorithm is an $(\epsilon, 3/4)$ tester as well. At a high level, the expected value
 443 of \bar{Z} is zero when $p = q$; whereas it is larger than $\Theta(\tau)$ when p and q are ϵ -far from each other.
 444 By analyzing the variance of \bar{Z} , and using Chebyshev's inequality, we show that \bar{Z} is close to its
 445 expectation. More specifically, we prove the following claims:

- 446 • Completeness case: If p is equal to q , then \bar{Z} at most $\tau/2$ with high probability.
- 447 • Soundness case: If p is ϵ -far from q , then \bar{Z} is at least $3\tau/2$ with high probability.

In addition, the magnitude of the of the Laplace noise we add to \bar{Z} is small with high probability. Set τ , the threshold we used in the algorithm, to be $2c_0 s^2 \epsilon^2 / n'$. Using the CDF of Laplace distribution, we have

$$\Pr \left[|\hat{Z} - \bar{Z}| \geq \frac{\tau}{2} \right] \leq \exp \left(\frac{c_0 s^2 \epsilon^2 \xi}{n' \Delta(\bar{Z})} \right) \leq 0.01$$

where the last inequality is true if s is bounded from below as follows for a sufficiently large constant c_2 :

$$s \geq c_2 \cdot \frac{\sqrt{n' \Delta(\bar{Z})}}{\epsilon \sqrt{\xi}}.$$

448 If $|\hat{Z} - \bar{Z}|$ is less than $\tau/2$, then our claim above is sufficient to prove that the algorithm is $(\epsilon, 3/4)$ -
 449 tester. In the completeness case, with high probability, we will have $\hat{Z} < \bar{Z} + \tau/2 \leq \tau$, and in the
 450 soundness case, with high probability, we will have $\hat{Z} > \bar{Z} - \tau/2 \geq \tau$. Thus, in both case, the \hat{Z}
 451 will be on the correct side of the threshold. Hence, for the rest of the proof, we focus on the proof
 452 of the two claims we state.

453 To show the bounds for \bar{Z} , we introduce an auxiliary random variable W . We analyze the expected
 454 value, and the variance of W , and show that with high probability W is around its expectation by
 455 Chebyshev's inequality. Then, we use this fact about W to prove that \bar{Z} must be around its expected
 456 value as well, and achieve the desired bound for \bar{Z} with high probability.

More specifically, for a given X and a given π , we define W as:

$$W = Z - s^2 \cdot \|p^{(F)} - q^{(F)}\|_2^2,$$

457 where F is the set of flattening samples, $x_{\pi(\hat{s}+1)}, x_{\pi(\hat{s}+2)} \dots, x_{\pi(\hat{s}+\hat{k})}$. We can similarly define \bar{W}
 458 as well:

$$\bar{W}(X) := \mathbf{E}_{\pi, r}[W|X].$$

459 We analyze the expected value and the variance of \bar{W} . First, we define the following notations, $d_{\max}^{(F)}$
 460 and $d_{\min}^{(F)}$ to indicate the maximum and the minimum of the two quantities $\|p^{(F)}\|_2^2$ and $\|q^{(F)}\|_2^2$
 461 respectively. The expected value and the variance of Z , as defined in Equation 1, is given in the
 462 proof of Proposition 3.1 [19], if we fix the set of flattening samples F :

$$\mathbf{E}_{T, r}[Z|F] = s^2 \|p^{(F)} - q^{(F)}\|_2^2, \text{ and } \mathbf{Var}_{T, r}[Z|F] \leq 8s^3 \sqrt{d_{\max}^{(F)}} \|p^{(F)} - q^{(F)}\|_4^2 + 8s^2 d. \quad (6)$$

Using the above equation, we compute the expected value and the variance of \overline{W} . Note that since samples in X are i.i.d, the order of the samples can change neither the expected value nor the variance. Thus, by symmetrization, we can fix an order of the samples, namely π_0 . As mentioned before, the first \hat{s} samples in X are the test samples, and we denote them by T , and the next \hat{k} samples for the flattening samples and are denoted by F . Since T and F are completely separated and independent, by Equation 6, we have:

$$\begin{aligned}\mathbf{E}_X[\overline{W}] &= \mathbf{E}_X[\mathbf{E}_{\pi,r}[W|X]] = \mathbf{E}_X[\mathbf{E}_r[W|X, \pi_0]] \\ &= \mathbf{E}_F\left[\mathbf{E}_{T,r}\left[Z - s^2 \cdot \|p^{(F)} - q^{(F)}\|_2^2 \middle| F\right]\right] = 0\end{aligned}\quad (7)$$

Moreover, given the variance bound in the Equation 6, we obtain the following bound for the variance of \overline{W} , we have:

$$\begin{aligned}\mathbf{Var}_X[\overline{W}] &= \mathbf{Var}_X[\mathbf{E}_{\pi,r}[W|X]] = \mathbf{Var}_X\left[\sum_{\pi} \mathbf{E}_r[W|X, \pi] \cdot \Pr[\pi]\right] \\ &= \sum_{\pi_1} \sum_{\pi_2} \Pr[\pi_1] \cdot \Pr[\pi_2] \cdot \mathbf{Cov}_X(\mathbf{E}_r[W|X, \pi_1], \mathbf{E}_r[W|X, \pi_2]) \\ &\leq \frac{1}{2} \sum_{\pi_1} \sum_{\pi_2} \Pr[\pi_1] \cdot \Pr[\pi_2] \cdot (\mathbf{Var}_X[\mathbf{E}_r[W|X, \pi_1]] + \mathbf{Var}_X[\mathbf{E}_r[W|X, \pi_2]]) \\ &= \mathbf{Var}_X[\mathbf{E}_r[W|X, \pi_0]] = \mathbf{Var}_{F,T}[\mathbf{E}_r[W|F, T]] \\ &= \mathbf{E}_F[\mathbf{Var}_T[\mathbf{E}_r[W|F]]] + \mathbf{Var}_F[\mathbf{E}_{T,r}[W|F]] \\ &\leq \mathbf{E}_F\left[8s^3 \cdot \sqrt{d_{\max}^{(F)}} \cdot \|p^{(F)} - q^{(F)}\|_4^2 + 8s^2 \cdot d_{\max}^{(F)}\right].\end{aligned}$$

Using Chebyshev's inequality, \overline{W} cannot be far from its expectation which is zero. More precisely, given a constant c_0 , we prove that there exists a sufficiently large constant c_2 such that $\Pr_X\left[|\overline{W}| \geq c_0 \cdot \frac{s^2 \epsilon^4}{n}\right]$ is at most 0.01 assuming we have:

$$s \geq c_2 \cdot \frac{n' \cdot \sqrt{\mathbf{E}_F[d_{\min}^{(F)}]}}{\epsilon^2}. \quad (8)$$

We consider the soundness case and the completeness case below separately.

Completeness Case: If p is equal to q , no matter what F and the bucketing are, $\|p^{(F)} - q^{(F)}\|_2^2$ is zero. Therefore, \overline{W} is always equal to \overline{Z} just by definition of W . In fact, we have $\overline{Z} = \overline{W} = \overline{W} - \mathbf{E}_X[\overline{W}]$. Also, the ℓ_2 -norms of $p^{(F)}$ and $q^{(F)}$ are the same. Thus, the minimum and the maximum of $\|p^{(F)}\|_2^2$ and $\|q^{(F)}\|_2^2$ are equal. Thus, the probability of \overline{Z} be above $\tau/2$ is bounded by 0.01 using the Chebyshev inequality:

$$\begin{aligned}\Pr_X\left[\overline{Z} \geq \frac{\tau}{2}\right] &\leq \Pr_X\left[|\overline{W} - \mathbf{E}_X[\overline{W}]| \geq \frac{c_0 s^2 \epsilon^2}{n'}\right] \leq \frac{n'^2 \cdot \mathbf{E}_F[8s^2 d_{\max}^{(F)}]}{c_0^2 s^4 \epsilon^4} \\ &= \Theta\left(\frac{n'^2 \cdot \mathbf{E}_F[d_{\min}^{(F)}]}{s^2 \epsilon^4}\right) \leq 0.01,\end{aligned}$$

where the last inequality is true assuming Equation 8 for a sufficiently large constant c_2 .

Soundness Case: In this case p is ϵ -far from q . Before showing that \overline{W} is close to zero with high probability, we establish two inequalities as below. First, observe that flattening does not change the ℓ_1 -distance between two distributions due to the following:

$$\|p - q\|_1 = \sum_{i=i}^n |p(i) - q(i)| = \sum_{i=1}^n \sum_{j=1}^{b_i} \frac{|p(i) - q(i)|}{b_i} = \|p^{(F)} - q^{(F)}\|.$$

Thus, we have the following lower bound for the ℓ_2 -distance between $p^{(F)}$ and $q^{(F)}$ for any F :

$$\frac{\epsilon^2}{n'} \leq \frac{\|p - q\|_1^2}{n'} \leq \frac{\|p^{(F)} - q^{(F)}\|_1^2}{n'} \leq \|p^{(F)} - q^{(F)}\|_2^2.$$

Therefore, the above inequality is true in expectation as well:

$$\mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right] \geq \frac{\epsilon^2}{n'}. \quad (9)$$

Second, we provide the following lemma to show a bound for $\mathbf{E}_F \left[d_{\max}^{(F)} \right]$.

Lemma A.1. *Assume F is a random set of samples to be used for flattening. Then, we have:*

$$\mathbf{E}_F \left[d_{\max}^{(F)} \right] \leq \Theta \left(\mathbf{E}_F \left[d_{\min}^{(F)} \right] + \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right] \right)$$

For the proof of the lemma, see Section D.

By the Chebyshev inequality and the above lemma, we show that \bar{W} is close to its expectation, i.e., zero, with high probability. Using Jensen's inequality, Equation 5, Equation 9, and Lemma A.1, we have:

$$\begin{aligned} \Pr_X \left[|\bar{W} - \mathbf{E}_X[\bar{W}]| \geq c_0 \cdot \mathbf{E}_F \left[s^2 \|p^{(F)} - q^{(F)}\|_2^2 \right] \right] &\leq \frac{\mathbf{E}_F \left[8s^3 \sqrt{d_{\max}^{(F)}} \|p^{(F)} - q^{(F)}\|_4^2 + 8s^2 d_{\max}^{(F)} \right]}{c_0^2 \cdot \mathbf{E}_F \left[s^2 \|p^{(F)} - q^{(F)}\|_2^2 \right]^2} \\ &\leq \Theta \left(\frac{\sqrt{\mathbf{E}_F \left[d_{\max}^{(F)} \right]} \cdot \sqrt{\mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_4^4 \right]}}{s \cdot \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]^2} + \frac{\mathbf{E}_F \left[d_{\max}^{(F)} \right]}{s^2 \cdot \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]^2} \right) \\ &\leq \Theta \left(\frac{\sqrt{\mathbf{E}_F \left[d_{\max}^{(F)} \right]}}{s \cdot \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]} + \frac{\mathbf{E}_F \left[d_{\max}^{(F)} \right]}{s^2 \cdot \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]^2} \right) \\ &\leq \Theta \left(\frac{\sqrt{\mathbf{E}_F \left[d_{\max}^{(F)} \right]}}{s \cdot \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]} \right) \\ &\leq \Theta \left(\frac{\sqrt{\mathbf{E}_F \left[d_{\min}^{(F)} \right]} + \sqrt{\mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]}}{s \cdot \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]} \right) \\ &\leq 0.01 \end{aligned}$$

where the last inequality is true when s is larger than the bound given in Equation 8 below for a sufficiently large constant, c_2 .

$$\begin{aligned} s &\geq c_2 \cdot \left(\frac{n' \cdot \sqrt{\mathbf{E}_F \left[d_{\min}^{(F)} \right]}}{\epsilon^2} \right) \geq \Theta \left(\frac{n' \cdot \sqrt{\mathbf{E}_F \left[d_{\min}^{(F)} \right]}}{\epsilon^2} + \frac{\sqrt{n'}}{\epsilon} \right) \geq \\ &\geq \Theta \left(\frac{\sqrt{\mathbf{E}_F \left[d_{\min}^{(F)} \right]}}{\mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]} + \frac{1}{\sqrt{\mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]}} \right) \end{aligned}$$

Now, by definition of \bar{Z} and Equation 4, we show Z has to be at least $3\tau/2$ with high probability.

Therefore, with high probability have:

$$\begin{aligned}\bar{Z} &= \bar{W} + \mathbf{E}_\pi \left[s^2 \|p^{(F)} - q^{(F)}\|_2^2 \right] \geq -c_0 \cdot \mathbf{E}_F \left[s^2 \|p^{(F)} - q^{(F)}\|_2^2 \right] + 4c_0 \cdot \mathbf{E}_F \left[s^2 \|p^{(F)} - q^{(F)}\|_2^2 \right] \\ &\geq 3c_0 \cdot \mathbf{E}_F \left[s^2 \|p^{(F)} - q^{(F)}\|_2^2 \right] \geq \frac{3c_0 s^2 \epsilon^2}{n'}.\end{aligned}$$

By taking the union bound, the probability of having too large Laplace noise or a too large $|\bar{W}|$ is at most 0.02. Moreover, Equation 5 and Equation 4 do not hold with probability at most 0.1. Thus, with probability at least 3/4, the algorithm output the correct answer. \square

B Testing Closeness of Distributions with Unequal Sample Sizes

In this section, we prove that the non-private algorithm for testing closeness using unequal sample sizes, provided in [25] with a small modification, is a proper procedure. Therefore, we can turn it into a private tester using our approach provided in Section 3. We also analyze the sample complexity of the tester.

Assume \mathcal{A} is the non-private procedure for testing closeness of p and q using unequal sample sizes. First, we explain how \mathcal{A} works. To generate a sample from p (or q) the algorithm simply draw an i.i.d. sample from p (or q). Assume k_1 , k_2 , and s are three parameters that we determine later. \hat{k}_1 , \hat{k}_2 , and \hat{s} indicate three Poisson random variables with mean k_1 , k_2 and s respectively. \mathcal{A} draws $\hat{s} + \hat{k}_1$ from p and $\hat{s} + \hat{k}_2$ samples from q . For the number of buckets, \mathcal{A} uses the following process. Let F be the number of a set of \hat{k}_1 samples from p and \hat{k}_2 samples from q . The number of buckets for element i is determined by the number of instances of i in F plus one.

Theorem B.1. *There exists a ξ -differentially private algorithm that uses $k_1 = \Omega(\max(n^{2/3}/\epsilon^{4/3}, \sqrt{n}/\epsilon\sqrt{\xi}))$ samples from p , $\Theta(\max(n/(\epsilon^2 \sqrt{\min(n, k_1)}), \sqrt{n}/\epsilon^2, \sqrt{n}/\epsilon\sqrt{\xi}, 1/\epsilon^2 \xi))$ from both p and q and distinguishes the following cases with probability at least 0.8:*

- *Completeness case:* $p = q$
- *Soundness case:* $\|p - q\|_1 > \epsilon$.

Proof: The goal is to transform the problem to the generate tester we provided in Section 3. First, in Lemma B.2 we show that the non-private algorithm in [25] is a “proper procedure”. Using Theorem 3.3, the existence of the tester with the sample complexity s for the test part is immediate where s is at least the bound below

$$s \geq \Theta \left(\frac{n' \cdot \sqrt{\mathbf{E}_F [\min(\|p^{(F)}\|_2^2, \|q^{(F)}\|_2^2)]}}{\epsilon^2} + \frac{\sqrt{n' \Delta(\bar{Z})}}{\epsilon \sqrt{\xi}} \right). \quad (10)$$

We first show the relationship between s above and the rest of the parameters we have. Then we set the parameters k_1 , k_2 , and s and analyze the sample complexity. Without loss of generality assume $k_1 \geq k_2$. Note that after flattening the size of the domain increases to $n' = \Theta(n + k_1 + k_2)$ with high probability. Then, in Lemma B.5, we show that the proposed statistic, \bar{Z} , has a bounded sensitivity:

$$\Delta(\bar{Z}) \leq \Theta \left(\frac{k_1}{k_1 + s} \cdot \left(\frac{s + k_2}{k_2} \right)^2 \right).$$

In addition, it is shown in [25] that the probability of the expected ℓ_2 -norm of p after flattening is at most $1/k_1$. Moreover, adding more flattening samples does not increase this quantity. Hence, we have:

$$\mathbf{E}_F \left[\min(\|p^{(F)}\|_2^2, \|q^{(F)}\|_2^2) \right] \leq \frac{1}{k_1}$$

We consider the following cases:

- **case 1:** $\epsilon = \Omega(n^{-1/4})$ and $\epsilon^2 \xi = \Omega(n^{-1})$: In this case, we have the following properties:

$$\Theta\left(\frac{\sqrt{n}}{\epsilon^2}\right) \leq \Theta\left(\frac{n^{2/3}}{\epsilon^{4/3}}\right) \leq \Theta(n), \quad \text{and} \quad \Theta\left(\frac{1}{\epsilon^2 \xi}\right) \leq \Theta\left(\frac{\sqrt{n}}{\epsilon \sqrt{\xi}}\right) \leq \Theta(n).$$

Let k_1 be a number in the range below:

$$\Theta\left(\max\left(\frac{n^{2/3}}{\epsilon^{4/3}} + \frac{\sqrt{n}}{\epsilon \sqrt{\xi}}\right)\right) \leq k_1 \leq \Theta(n).$$

Hence, n' is $\Theta(n)$. Then, we set s and k_2 as follows:

$$k_2 := s := \Theta\left(\max\left(\frac{n}{\epsilon^2 \sqrt{k_1}}, \frac{\sqrt{n}}{\epsilon \sqrt{\xi}}\right)\right)$$

519 $\Delta(\bar{Z})$ is $O(1)$ in this case. Therefore, s is $\Omega(n/\epsilon^2 \sqrt{k_1} + \sqrt{n \Delta(\bar{Z})}/(\epsilon \sqrt{\xi}))$ and the condi-
520 tion in Equation 10 holds.

- **case 2:** $\epsilon = o(n^{-1/4})$: In this case, \sqrt{n}/ϵ^2 is $\Omega(n)$. Thus, we cannot avoid sample complexity of $\Omega(n)$. We set k_1 and k_2 to be equal to n , and we set s to be the following:

$$s := \Theta\left(\max\left(\frac{\sqrt{n}}{\epsilon^2}, \frac{1}{\epsilon^2 \xi}\right)\right).$$

521 Clearly n' is still $\Theta(n)$, and s is $\Omega(n/\sqrt{k_1 \epsilon^2})$. In this case $\Delta(\bar{Z})$ is $\Theta(s/n)$. Hence, in
522 order to have s at least $\Omega(\sqrt{n' \Delta(\bar{Z})}/\epsilon \sqrt{\xi})$, it suffices to have $s = \Omega(1/\epsilon \sqrt{\xi})$.

523

□

524 B.1 Non-Private Closeness Tester Is a Proper Procedure

525 **Lemma B.2.** *Procedure A explained above is a proper procedure according to Definition 3.2.*

526 **Proof:** First, we show the number of samples we generate is not too far from their expectation.
527 Hence X can be a set with a bounded number of samples. In the following lemma we show if the
528 means are larger than a fixed constant, then with probability 0.01 we can assume the number of
529 samples from each of distributions is at most three times larger than their means.

530 **Lemma B.3.** *Assume random variable x is drawn from $\text{Poi}(\lambda)$. If λ is at least $1.5 \cdot \ln(1/c)$, then
531 the probability of x being larger than 3λ is at most $1 - c$.*

532 Now, we only need to show that inequalities in Equation 4 and Equation 5 are correct. Before
533 proving the equations, we provide an insightful information about the distribution over the b_i 's. It is
534 clear that for a fixed i , $b_i - 1$ is an independent Poisson random variable with mean $k_1 p(i) + k_2 q(i)$.
535 More precisely, we can think of $b_i - 1$ as the sum of two random variables $b_{i,1} \sim \text{Poi}(k_1 p(i))$ and
536 $b_{i,2} \sim \text{Poi}(k_2 q(i))$ plus one. However, assume a random set of samples, X , is given to us with $t_{i,1}$
537 instances of i from p , and $t_{i,2}$ instances of i from q . Then, considering a random permutation of
538 samples, then $b_{i,j}$ is a binomial random variable from $\text{Bin}(t_{i,j}, k_j/(k_j + s))$ for $j = 1, 2$.

539 Now, we focus on proving Equation 4. Fix a set of sample X and a domain element i . Using Jensen's
540 inequality, we have

$$\mathbf{E}_\pi \left[\frac{1}{b_i(X, \pi)} \right] = \mathbf{E}_{b_{i,1}, b_{i,2}} \left[\frac{1}{b_i} \right] \geq \frac{1}{\mathbf{E}_{b_{i,1}, b_{i,2}}[b_i]} = \frac{1}{t_{i,1} k_1/(k_1 + s) + t_{i,2} k_2/(k_2 + s) + 1}.$$

541 Note that by Markov's inequality the probability of any of $t_{i,1}$ or $t_{i,2}$ being 50 times³ larger than
542 their expectations is at most 0.04. Therefore, with probability 0.96 assume they are at most 50 times
543 their expectation. Since $t_{i,1}$ and $t_{i,2}$ are two Poisson random variables with means $p(i)(k_1 + s)$ and
544 $q(i)(k_2 + s)$, we can bound the above quantity as follows:

$$\mathbf{E}_\pi \left[\frac{1}{b_i(\pi)} \right] \geq \frac{1}{50 p(i) k_1 + 50 q(i) k_2 + 1} = \frac{1}{50 \lambda + 1}$$

³Needless to say, we are not optimizing constants here.

where we use λ to denote $p(i)k_1 + q(i)k_2$. On the other hand, the expected value of $1/b_i$ when X has not been observed is the following:

$$\mathbf{E}_F \left[\frac{1}{b_i} \right] = \mathbf{E}_{x \sim \text{Poi}(\lambda)} \left[\frac{1}{x+1} \right] = \frac{1 - e^{-\lambda}}{\lambda} \leq \frac{65}{50\lambda + 1}.$$

Putting all of the above facts together, we conclude:

$$\Pr_{t_{i,1}, t_{i,2}} \left[\mathbf{E}_\pi \left[\frac{(p(i) - q(i))^2}{b_i(\pi)} \right] \geq \frac{1}{65} \cdot \mathbf{E}_F \left[\frac{(p(i) - q(i))^2}{b_i} \right] \right] \geq 0.96. \quad (11)$$

We define a random variable x_i over the randomness of $t_{i,1}$ and $t_{i,2}$ to be the following:

$$x_i := \mathbf{E}_\pi \left[\frac{(p(i) - q(i))^2}{b_i(\pi)} \right]$$

Note that by the *Poissonization method*, all the number of instances of a particular element are independent from the rest. Hence, x_i 's are independent given the independence of $t_{i,j}$'s. In addition, we prove the following lemma, to bound the sum of x_i 's from below:

Lemma B.4. Assume we have n independent random variables x_1, x_2, \dots, x_n in the range $[0, +\infty)$. Suppose each x_i is at least A_i with probability $p \geq 0.95$ where A_i is a fixed number. Then, with probability at least 0.9, $\sum_{i=1}^n x_i$ is at least $0.1 \sum_{i=1}^n A_i$.

For the proof of the lemma, see Section D.

Using Equation 11, and Lemma B.4, with probability 0.9 we have:

$$\begin{aligned} \mathbf{E}_\pi \left[\|p^{(F)} - q^{(F)}\|_2^2 \right] &= \mathbf{E}_\pi \left[\frac{(p(i) - q(i))^2}{b_i(\pi)} \right] \geq \frac{1}{650} \cdot \sum_{i=1}^n \mathbf{E}_F \left[\frac{(p(i) - q(i))^2}{b_i} \right] \\ &= 4c_0 \cdot \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]^2. \end{aligned}$$

where $c_0 = 1/26000$. Hence, the proof of Equation 4 is complete.

Now, we focus on proving Equation 5. To prove the inequality, it suffices to show that $\mathbf{E}_F[1/b_i^3]$ is $O(\mathbf{E}_F[1/b_i]^2)$. Note one can think of b_i to be equal to $x+1$ where x is a Poisson random variable with mean $\lambda' = p(i)k_1 + q(i)k_2$. Thus, we have:

$$\begin{aligned} \mathbf{E}_F \left[\frac{1}{b_i^3} \right] &= \mathbf{E}_{x \sim \text{Poi}(\lambda')} \left[\frac{1}{(x+1)^3} \right] \leq \mathbf{E} \left[\frac{6}{(x+1)(x+2)(x+3)} \right] \leq 6 \cdot \sum_{x=0}^{\infty} \frac{e^{-\lambda'} \lambda'^x}{(x+3)!} \\ &= \frac{6}{\lambda'^3} \sum_{y=3}^{\infty} \frac{e^{-\lambda'} \lambda'^y}{y!} = \frac{6(1 - e^{-\lambda'} - e^{-\lambda'} \lambda' - e^{-\lambda'} \lambda'^2/2)}{\lambda'^3} \leq 6 \cdot \left(\frac{1 - e^{-\lambda'}}{\lambda'} \right)^2. \end{aligned} \quad (12)$$

On the other hand, we can compute the expected value of $1/b_i$ as follows:

$$\left(\mathbf{E}_F \left[\frac{1}{b_i} \right] \right)^2 = \left(\sum_{x=0}^{\infty} \frac{e^{-\lambda'} \lambda'^x}{(x+1)!} \right)^2 = \left(\frac{1}{\lambda'} \sum_{y=1}^{\infty} \frac{e^{-\lambda'} \lambda'^y}{y!} \right)^2 = \left(\frac{1 - e^{-\lambda'}}{\lambda'} \right)^2.$$

Putting these two equations together, one can conclude the Equation 5:

$$\begin{aligned} \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_4^4 \right] &= \sum_{i=1}^n \mathbf{E}_F \left[\frac{(p(i) - q(i))^4}{b_i^4} \right] \leq 6 \cdot \sum_{i=1}^n \left(\mathbf{E}_F \left[\frac{(p(i) - q(i))^2}{b_i} \right] \right)^2 \\ &\leq 6 \cdot \left(\sum_{i=1}^n \mathbf{E}_F \left[\frac{(p(i) - q(i))^2}{b_i} \right] \right)^2 = 6 \cdot \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]^2. \end{aligned}$$

Therefore, the statement of the lemma is concluded. \square

561 B.2 Bounding the Sensitivity

562 In this section, we provide an upper bound for the amount that the statistic changes if we change
 563 one sample in the input. In other words, we find an upper bound for the sensitivity of \bar{Z} as define in
 564 Equation 2.

565 We start off with defining the notation we use in this section. Let $X = (X_1, X_2)$ be a set of samples
 566 consist of m_1 samples from p and m_2 samples from q . As we had before, $\hat{k}_j = \text{Poi}(k_j)$ and
 567 $\hat{s} = \text{Poi}(s)$ to denote the number of samples for test and the flattening respectively. Note that now
 568 that $m_1 = \hat{k}_j + \hat{s}$ is observed, then \hat{k}_j is a binomial random variable, $\text{Bin}(m, k_j / (k_j + s))$. Let F_1
 569 and F_2 denote the set of sample from distribution p and q we use for the flattening. Given X , F_1
 570 and F_2 are determined by the \hat{k}_j 's, and the order of the samples which is determined by π . Although
 571 the F_j 's are two sets of ordered samples, the order of the element in them does not matter. We may
 572 consider them equivalent to set of *indices* I_j , such that the r -th sample is in F_j if and only if r is
 573 in I_j . In this sense, we can define the probability of an index set I_j to be the probability of taking
 574 a flattening set F_j , such that F_j is equivalent to I_j . The randomness in this probability is taken
 575 over the choice of the permutation π and \hat{k}_j . In addition, for each element i , we use the following
 576 notation:

577 • $t_{i,j}$ is the number of instances of element i in the set X_j .

578 • $k_{i,j}$ is the number of instances of element i in the set F_j .

579 • $s_{i,j}$ is the number of instances of element i in the set $X_j \setminus F_j$.

580 Based on these notations, the statistic \bar{Z} is the following, we use the equivalent intermediate statistic
 581 as indicated in Lemma ??, and denote it by $z_i(X, I_1, I_2)$

$$\begin{aligned} \bar{Z} &= \mathbf{E}_{\pi, r}[Z] = \sum_{I_1, I_2} \mathbf{Pr}[I_1] \cdot \mathbf{Pr}[I_2] \cdot \sum_{i=1}^n z_i(X, I_1, I_2) \\ &= \sum_{I_1, I_2} \mathbf{Pr}[I_1] \cdot \mathbf{Pr}[I_2] \cdot \sum_{i=1}^n \frac{(s_{i,1} - s_{i,2})^2 - s_{i,1} - s_{i,2}}{k_{i,1} + k_{i,2} + 1} \end{aligned}$$

582 Assume $X' = (X'_1, X'_2)$ is also a set of samples such that it differs in exactly one sample compared
 583 to X . We similarly define all the notation for X' as well by adding the prime notation to each letter.
 584 Without loss of generality, we assume that the r -th sample in X_1 , namely x_r , is different from the
 585 r -th sample in X'_1 , namely x'_r , and all other samples are the same. Now, we are ready to bound the
 586 sensitivity of \bar{Z} in the following lemma:

Lemma B.5. *For two neighboring sample sets, X_1, X_2 , and X'_1, X'_2 , and for fixed k_1 and k_2 , we have:*

$$|\bar{Z} - \bar{Z}'| \leq \Theta \left(\frac{k_1}{k_1 + s} \cdot \left(\frac{s + k_2}{k_2} \right)^2 + \frac{k_2}{k_2 + s} \cdot \left(\frac{s + k_1}{k_1} \right)^2 \right)$$

587 **Proof:** By the definition of \bar{Z} , by the triangle inequality and Bayes' law, we can find an upper bound
 588 for the difference of the statistics as follows:

$$\begin{aligned}
 |\bar{Z} - \bar{Z}'| &= \left| \sum_{I_1, I_2} \Pr[I_1] \cdot \Pr[I_2] \cdot \sum_{i=1}^n z_i(X, I_1, I_2) - z'_i(X', I_1, I_2) \right| \\
 &\leq \sum_{i=1}^n \sum_{I_2} \Pr[I_2] \cdot \sum_{I_1} \Pr[I_1] \cdot |z_i(X, I_1, I_2) - z'_i(X', I_1, I_2)| \\
 &\leq \sum_{i=1}^n \sum_{I_2} \Pr[I_2] \cdot \sum_{I_1} \Pr[I_1 | r \in I_1] \cdot \Pr[r \in I_1] \cdot |z_i(X, I_1, I_2) - z'_i(X', I_1, I_2)| \\
 &\quad + \sum_{i=1}^n \sum_{I_2} \Pr[I_2] \cdot \sum_{I_1} \Pr[I_1 | r \notin I_1] \cdot \Pr[r \notin I_1] \cdot |z_i(X, I_1, I_2) - z'_i(X', I_1, I_2)|.
 \end{aligned} \tag{13}$$

589 It is clear that if i is a domain element which none of its instances in the sample set is changed, then
 590 the number of occurrences of i remains unchanged in the same subsets of the X_j 's and the X'_j 's.
 591 Thus, $z_i - z'_i$ is equal to zero for $i \notin \{x_r, x'_r\}$. For now, we assume i is equal to x_r . One can
 592 replicate the same bound when $i = x'_r$.

593 It is clear that $t_{i,1} = t'_{i,1} + 1$. Fix a subset of indices, $I_2 \subseteq [m_2]$. Since X_2 and X'_2 are the same,
 594 then $s_{i,2} = s'_{i,2}$ and $k_{i,2} = k'_{i,2}$. Let ℓ be an index in $\{1, 2\}$ such that $t_{i,\ell}$ denotes the maximum
 595 of $t_{i,1}$ and $t_{i,2}$. Observe that we always have $s_{i,j} = t_{i,j} - k_{i,j}$ by definition. Now, we rewrite the
 596 difference of the z_i and z'_i as below using the triangle inequality:

$$\begin{aligned}
 |z_i(X, I_1, I_2) - z'_i(X', I_1, I_2)| &= \left| \frac{(s_{i,1} - s_{i,2})^2 - s_{i,1} - s_{i,2}}{k_{i,1} + k_{i,2} + 1} - \frac{(s'_{i,1} - s'_{i,2})^2 - s'_{i,1} - s'_{i,2}}{k'_{i,1} + k'_{i,2} + 1} \right| \\
 &= \left| \frac{(s_{i,1} - s_{i,2})^2 - s_{i,1} - s_{i,2}}{k_{i,1} + k_{i,2} + 1} - \frac{(s'_{i,1} - s_{i,2})^2 - s'_{i,1} - s_{i,2}}{k'_{i,1} + k_{i,2} + 1} \right| \\
 &\leq \left| \frac{(s_{i,1} - s_{i,2})^2 - s_{i,1} - s_{i,2}}{k_{i,1} + k_{i,2} + 1} - \frac{(s'_{i,1} - s_{i,2})^2 - s'_{i,1} - s_{i,2}}{k_{i,1} + k_{i,2} + 1} \right| \\
 &\quad + \left| \frac{(s'_{i,1} - s_{i,2})^2 - s'_{i,1} - s_{i,2}}{k_{i,1} + k_{i,2} + 1} - \frac{(s'_{i,1} - s_{i,2})^2 - s'_{i,1} - s_{i,2}}{k'_{i,1} + k_{i,2} + 1} \right|
 \end{aligned}$$

597 Note that if r is in I_1 , the number of instance of x_r in the flattening set changes. More precisely,
 598 $k_{i,1}$ is $k'_{i,1} + 1$. However, $s_{i,1}$ remains equal to $s'_{i,1}$ and the changed sample does not affect it, so the
 599 second to the last line above is zero. Similarly, if r is not in I_1 , then $k_{i,1}$ remains the same as $k'_{i,1}$,
 600 and $s_{i,1} = s'_{i,1} + 1$, so the last line above will be zero. $s_{i,j}$ and $k_{i,j}$ are at most $t_{i,j}$ by definition
 601 Therefore, if we use the fact that $s_{i,j}$ and $k_{i,j}$ are at most $t_{i,j}$ by definition, then we have:

$$\begin{aligned}
 &\sum_{I_1} \Pr[I_1] \cdot |z_i(X, I_1, I_2) - z'_i(X', I_1, I_2)| \\
 &\leq \sum_{I_1} \Pr[I_1 | r \in I_1] \cdot \Pr[r \in I_1] \cdot 2 t_{i,\ell}^2 \cdot \left| \frac{1}{k_{i,1} + k_{i,2} + 1} - \frac{1}{k'_{i,1} + k_{i,2} + 1} \right| \\
 &\quad + \sum_{I_1} \Pr[I_1 | r \notin I_1] \cdot \Pr[r \notin I_1] \cdot \frac{2 t_{i,\ell}}{k_{i,1} + k_{i,2} + 1}.
 \end{aligned}$$

602 Using the properties of the Poissonization method, given that we observed $t_{i,j} = k_{i,j} + s_{i,j}$, then
 603 $k_{i,j}$ is Binomial random variable: $\text{Bin}(t_{i,j}, k_j/(k_j + s))$. Given this fact, the probability of $r \in I$

604 is $k_1/(k_1 + s)$. Therefore, we have:

$$\begin{aligned}
& \sum_{I_2} \Pr[I_2] \cdot \sum_{I_1} \Pr[I_1] \cdot |z_i(X, I_1, I_2) - z'_i(X', I_1, I_2)| \\
&= \sum_{I_2} \Pr[I_2] \cdot \sum_{I_1} \Pr[I_1 | r \in I_1] \cdot \frac{k_1}{k_1 + s} \cdot \frac{2t_{i,\ell}^2}{(k'_{i,1} + k_{i,2} + 2) \cdot (k'_{i,1} + k_{i,2} + 1)} \\
&+ \sum_{I_2} \Pr[I_2] \cdot \sum_{I_1} \Pr[I_1 | r \notin I_1] \cdot \frac{s}{k_1 + s} \cdot \frac{2t_{i,\ell}}{k_{i,1} + k_{i,2} + 1} \\
&\leq \frac{2k_1 t_{i,\ell}^2}{k_1 + s} \cdot \sum_{I_2} \Pr[I_2] \cdot \sum_{I_1} \Pr[I_1 | r \in I_1] \cdot \frac{1}{(k'_{i,\ell} + 2) \cdot (k'_{i,\ell} + 1)} \\
&+ \frac{2s t_{i,\ell}}{k_1 + s} \cdot \sum_{I_2} \Pr[I_2] \cdot \sum_{I_1} \Pr[I_1 | r \notin I_1] \cdot \frac{1}{k_{i,\ell} + 1} \\
&\leq \frac{2k_1 t_{i,\ell}^2}{k_1 + s} \cdot \mathbf{E}_{k'_{i,\ell} \sim \text{Bin}(t_{i,j} - 1, k_1/(k_1 + s))} \left[\frac{1}{(k'_{i,\ell} + 2)(k'_{i,\ell} + 1)} \right] \\
&+ \frac{2s t_{i,\ell}}{k_1 + s} \cdot \mathbf{E}_{k_{i,\ell} \sim \text{Bin}(t_{i,j} - 1, k_1/(k_1 + s))} \left[\frac{1}{k_{i,\ell} + 1} \right].
\end{aligned}$$

605 Using Lemma D.2 and Lemma D.3, we have:

$$\begin{aligned}
& \sum_{I_2} \Pr[I_2] \cdot \sum_{I_1} \Pr[I_1] \cdot |z_i(X, I_1, I_2) - z'_i(X', I_1, I_2)| \\
&\leq \frac{2k_1}{k_1 + s} \cdot \left(\frac{s + k_\ell}{k_\ell} \right)^2 + \frac{2s}{k_1 + s} \cdot \frac{s + k_\ell}{k_\ell}.
\end{aligned}$$

606 Note that ℓ can be one or two. Also, the upper bound for $\bar{Z} - \bar{Z}'$ is twice as above, since we have
607 to consider the case when $i = x'_r$. Moreover, the bound we get is based on this assumption that the
608 changed sample is in X_1 , so to find the upper bound we need to consider both cases. Hence, we
609 have the following bound:

$$|\bar{Z} - \bar{Z}'| \leq \Theta \left(\frac{k_1}{k_1 + s} \cdot \left(\frac{s + k_2}{k_2} \right)^2 + \frac{k_2}{k_2 + s} \cdot \left(\frac{s + k_1}{k_1} \right)^2 \right)$$

610 and the proof is complete. \square

611 C Testing independence of two random variables

612 In this section, we provide a ξ -differentially private tester for testing independence of two random
613 variables. The idea is to reduce the optimal non-private tester, delivered in [25], to a private one
614 using the technique we explained in Section 3.

615 We start off with defining the problem and the non-private procedure \mathcal{A} that reduces testing inde-
616 pendence to the testing closeness of two distributions. Assume p is a distribution over $[n] \times [m]$.
617 Without loss of generality, we assume $m \leq n$. Suppose we receive samples (x, y) from p . We say
618 distribution p is an independent distribution, if the x 's and the y 's are independent from each other.
619 The goal is to distinguish whether p is an independent distribution or is it ϵ -far from any independent
620 distribution over $[n] \times [m]$. It is known that if p is an independent distribution, then p is equal to
621 $p_1 \times p_2$, and if p is ϵ -far from being independent, then p is $\epsilon/3$ -far from $p_1 \times p_2$ where p_1 and
622 p_2 are the marginal distributions of p . Using this fact, the non-private tester reduces the problem to
623 testing the closeness of p and $q := p_1 \times p_2$ [6].

624 Here, we describe a proper procedure, say \mathcal{A} , for reducing testing independence of p to the testing
625 closeness of p and q , so it can be turned into a private algorithm using the method explained in
626 Section 3. First, we describe the sampling scheme of the procedure: For every sample that the
627 algorithm needs, it draws two samples, and puts them in a *block*. We denote a block of two samples

628 (x_1, y_1) and (x_2, y_2) by $\langle (x_1, y_1), (x_2, y_2) \rangle$. We can use the samples in block to obtain a sample
 629 from p, p_1, p_2 , and q as follows. To get a sample from p , we always take the first samples, (x_1, y_1) ,
 630 in the block. We take x_1, y_2 as two the samples from p_1 and p_2 . In addition, since x_1 and y_2
 631 are two independent random variables, (x_1, y_2) is a sample from q . Also, we use “dot notation” to
 632 indicate an arbitrary element in the domain. For example, for a given x , (x, \cdot) is a sample that its first
 633 coordinate is x and the second coordinate can be any y in $[m]$. Similarly, we use the same notation
 634 to refer to a block, for example for a given x and y , $\langle (x, \cdot), (\cdot, y) \rangle$ indicates a block that the first
 635 coordinate of the first sample is x and the second coordinate of the second sample is y , and the two
 636 other coordinates can be arbitrary elements in $[m]$ and $[n]$. Let X denotes the set of all blocks that
 637 are available to the procedure. We use f to denote a frequency of the blocks with a certain format
 638 in X , e.g., $f_{\langle (x, \cdot), (\cdot, y) \rangle}$ is the number of blocks in X that the first coordinate of the first sample is x .
 639 For the rest of this section, we focus on blocks and use them accordingly to extract a sample.

640 Suppose we have sample access to p , and we can draw blocks of samples from it. Procedure \mathcal{A}
 641 uses the blocks for the following purposes: the *flattening samples* are used to determine the number
 642 of buckets for each domain element. They are designed to make sure that the ℓ_2 -norm of q after
 643 flattening is low. Also, the *test samples* are used to generate samples from two distributions p and q ,
 644 which we test their closeness. Below is how the algorithm will determine these samples. For now,
 645 assume $k^{(p_1)}, k^{(p_2)}, k^{(p)}, k^{(q)}$, and s are parameters that we determine later.

646 **Flattening samples and the number of buckets:** We flatten distribution p using samples from
 647 the marginal distributions p_1 and p_2 , and also samples from p itself. More specifically, we draw four
 648 sets of blocks from p , namely $F^{(p_1)}, F^{(p_2)}, F^{(p)}, F^{(q)}$, which contain $\text{Poi}(k^{(p_1)}), \text{Poi}(k^{(p_2)}),$
 649 $\text{Poi}(k^{(p)}), \text{Poi}(k^{(q)})$ blocks respectively. We refer to the samples in these sets as flattening sam-
 650 ples, and denote the collection of them by F . As we discuss earlier, we extract samples from the
 651 blocks in these sets to obtain samples from p_1, p_2, p , and q . More specifically, we use the following
 652 notation for the number of occurrences of each sample obtained from each set:

- 653 • $k_x^{(p_1)}$ denotes the number of occurrences of the blocks of the form $\langle (x, \cdot), (\cdot, \cdot) \rangle$ in the
 654 flattening set $F^{(p_1)}$.
- 655 • $k_y^{(p_2)}$ denotes the number of occurrences of the blocks of the form $\langle (\cdot, \cdot), (\cdot, y) \rangle$ in the
 656 flattening set $F^{(p_2)}$.
- 657 • $k_{(x,y)}^{(p)}$ denotes the number of occurrences of the blocks of the form $\langle (x, y), (\cdot, \cdot) \rangle$ in the
 658 flattening set $F^{(p)}$.
- 659 • $k_{(x,y)}^{(q)}$ denotes the number of occurrences of the blocks of the form $\langle (x, \cdot), (\cdot, y) \rangle$ in the
 660 flattening set $F^{(q)}$.

Our procedure uses $b_{(x,y)}$ many buckets for a domain element (x, y) , where $b_{(x,y)}$ is defined as
 follows:

$$b_{(x,y)} = (k_x^{(p_1)} + 1)(k_y^{(p_2)} + 1) + k_{(x,y)}^{(p)} + k_{(x,y)}^{(q)}.$$

661 It is worth to note that $k_x^{(p_1)}$ is always determined by the first samples in the blocks, whereas $k_y^{(p_2)}$ is
 662 determined by the second samples in the blocks. Therefore, for all x 's and y 's, these quantities are
 663 independent of each other.

664 **Test samples:** To determine the test samples, we draw two sets of blocks $T^{(p)}$ and $T^{(q)}$. Each set
 665 contains $\text{Poi}(s)$ many blocks. The samples in $T^{(p)}$ and $T^{(q)}$ are our test samples, and we denote
 666 the collection of these two sets by T . The blocks in $T^{(p)}$ are used to obtain samples from p , and the
 667 blocks in $T^{(q)}$ are used to collect samples from q . In particular, we use the following notation for
 668 the number of occurrences of each domain element (x, y) .

- 669 • $s_{(x,y)}^{(p)}$ denotes the number of occurrences of the blocks of the form $\langle (x, y), (\cdot, \cdot) \rangle$ in the test
 670 set $T^{(p)}$.
- 671 • $s_{(x,y)}^{(q)}$ denotes the number of occurrences of the blocks of the form $\langle (x, \cdot), (\cdot, y) \rangle$ in the test
 672 set $T^{(q)}$.

Now, that we showed how procedure \mathcal{A} determines the number of samples, we prove it yields to a ξ -differentially private tester as well. At a high level, we first show that \mathcal{A} is a proper procedure for testing independence (Section C.1), then use our general closeness tester to achieve a ξ -differentially private tester. Since the sample complexity of the private tester depends on the sensitivity of the statistic we are using, we analyze the sensitivity of the statistic (Section C.2). In particular, we show if the number of occurrences of certain blocks in the sample set is "as expected," then the sensitivity is low, which results in a nearly optimal sample complexity for the private tester. Next, we develop a framework to extend the input domain of the private algorithm to any sample set (Section C.3). More formally, we have the following theorem:

Theorem C.1. *Let p be a distribution over $[n] \times [m]$. There exists a ξ -differentially private $(\epsilon, 2/3)$ tester for the testing independence of p that uses $\Theta(s)$ samples where s is:*

$$s = \Theta \left(\frac{n^{2/3} m^{1/3}}{\epsilon^{4/3}} + \frac{(m n)^{1/2}}{\epsilon^2} + \frac{(m n \log n)^{1/2}}{\epsilon \sqrt{\xi}} + \frac{\log n}{\epsilon^2 \xi} \right).$$

Proof: We first set up the parameters we use:

$$k^{(p_2)} = m, \quad k^{(p_1)} = \min(n, n^{2/3} m^{1/3} / \epsilon^{4/3}), \quad k^{(p)} = k^{(q)} = \min(m \cdot n, s),$$

$$\text{and} \quad s = c \cdot \left(\frac{n^{2/3} m^{1/3}}{\epsilon^{4/3}} + \frac{(m n)^{1/2}}{\epsilon^2} + \frac{(m n \log n)^{1/2}}{\epsilon \sqrt{\xi}} + \frac{\log n}{\epsilon^2 \xi} \right).$$

where c is a large constant. For sufficiently large m and n , with probability 0.99 the number of blocks in each set in $\mathcal{S} = \{F^{(p_1)}, F^{(p_2)}, F^{(p)}, F^{(q)}, T^{(p)}, T^{(q)}\}$ is within a constant factor of its expectation via Lemma B.3. Now, we show that \mathcal{A} that we describe above is a proper procedure:

Lemma C.2. *Procedure \mathcal{A} explained above is a proper procedure according to Definition 3.2 for testing independence of two random variable.*

The proof of the above Lemma is in Section C.1. Now, using Theorem 3.3, there exists a ξ -differentially private tester for the independence property which uses the following number of test samples:

$$s' := \Theta \left(\frac{n' \cdot \sqrt{\mathbf{E}_F[\min(\|p^{(F)}\|_2^2, \|q^{(F)}\|_2^2)]}}{\epsilon^2} + \frac{\sqrt{n' \Delta(\bar{Z})}}{\epsilon \sqrt{\xi}} \right).$$

Here, we show that $s \geq s'$, thus the number of samples that the procedure provides is enough by bounding n' , $\mathbf{E}_F[\min(\|p^{(F)}\|_2^2, \|q^{(F)}\|_2^2)]$, and $\Delta(\bar{Z})$. Note that n' is the new domain size which is equal to $\sum_{(x,y)} b_{(x,y)} = \Theta(m n)$. The expected of minimum of the ℓ_2 -norm of p and q is bounded using the result in Lemma 2.6 in [25].

$$\begin{aligned} \mathbf{E}_{F^{(p_1)}, F^{(p_2)}, F^{(p)}} \left[\min \left(\|p^{(F)}\|_2^2, \|q^{(F)}\|_2^2 \right) \right] &\leq \mathbf{E}_F \left[\|q^{(F)}\|_2^2 \right] \leq \sum_{x=1}^n \sum_{y=1}^m \frac{q(x, y)^2}{b_{(x, y)}} \\ &\leq \sum_{x=1}^n \sum_{y=1}^m \frac{p_1(x)^2 p_2(y)^2}{(k_x^{(p_1)} + 1) \cdot (k_y^{(p_2)} + 1)} \leq \|p_1^{(F^{(p_1)})}\| \cdot \|p_2^{(F^{(p_2)})}\| \leq \frac{1}{k^{(p_1)} k^{(p_2)}}. \end{aligned}$$

Moreover, in Section C.2, we provide the following bound for the sensitivity of the statistic:

Lemma C.3. *Given that the size of all flattening and test samples are within the constant factor of their expectations, the sensitivity of the statistic Z is bounded as follows:*

$$\Theta \left(\frac{s}{k^{(q)}} + \frac{s}{k^{(p)}} + \frac{s}{k^{(p)}} \cdot \frac{f_{\langle (.,b), (.,.) \rangle}}{f_{\langle (.,.), (.,b) \rangle} + 1} \right)$$

To get a bound on sensitivity, for now, suppose all the input block sets X has a desired property that the ratio between $f_{\langle (.,b), (.,.) \rangle}$ and $f_{\langle (.,.), (.,b) \rangle} + 1$ are bounded:

$$X \in \mathcal{X}^* := \left\{ X : \frac{f_{\langle(\cdot, b), (\cdot, \cdot)\rangle}}{f_{\langle(\cdot, \cdot), (\cdot, b)\rangle} + 1} \leq \tau \right\}$$

where $\tau = 1200 \ln n$. Thus, using the fact that X is in \mathcal{X}^* , one can obtain:

$$\Delta(\bar{Z}) \leq \Theta \left(\frac{s \log n}{mn} + \log n \right)$$

Now, we are ready to show that $s' \leq s$ implying that we have enough samples for the ξ -private tester. It is not hard to see that we have the following bounds (up to a constant factors):

$$\begin{aligned} s' &\leq \Theta \left(\frac{n' \cdot \sqrt{\mathbf{E}_F[\min(\|p^{(F)}\|_2^2, \|q^{(F)}\|_2^2)]}}{\epsilon^2} + \frac{\sqrt{n' \Delta(\bar{Z})}}{\epsilon \sqrt{\xi}} \right) \\ &\leq \Theta \left(\frac{mn}{\epsilon^2 \sqrt{k^{(p_1)} k^{(p_2)}}} + \frac{\sqrt{mn \Delta(\bar{Z})}}{\epsilon \sqrt{\xi}} \right) \\ &\leq \Theta \left(\frac{mn}{\epsilon^2 \sqrt{m \min(n, n^{2/3} m^{1/3} / \epsilon^{4/3})}} + \frac{\sqrt{mn}}{\epsilon \sqrt{\xi}} \cdot \sqrt{\frac{s \log n}{mn} + \log n} \right) \\ &\leq \Theta \left(\frac{n^{2/3} m^{1/3}}{\epsilon^{4/3}} + \frac{(mn)^{1/2}}{\epsilon^2} + \frac{\sqrt{mn}}{\epsilon \sqrt{\xi}} \cdot \sqrt{\frac{s \log n}{mn} + \log n} \right) \\ &\leq s \end{aligned}$$

Thus, given that X is in \mathcal{X}^* , there exists a ξ -differentially private tester that outputs the right answer with probability 0.8. This is sufficient to show that there exists an ξ -differentially private algorithm with asymptotically the same number of samples via Lemma C.6. \square

C.1 Non-Private Independence Tester is a proper procedure

Lemma C.2. *Procedure \mathcal{A} explained above is a proper procedure according to Definition 3.2 for testing independence of two random variable.*

Proof: Let X be the set of all blocks we received. Since the number of blocks in each of the flattening set and the test set is Poisson random variable, by Lemma B.3, we can conclude with probability 1-0.01 we draw at most three times more samples than what is expected. Hence X is a set with a bounded number of samples.

We start proving Equation 4 by recalling a fact about the Poissonization method. The number of blocks of a certain form in one of the flattening and test sets is a Binomial random variable with the bias that is proportional to the expected size of each set. For example, if X contains t blocks of the form $\langle(x, \cdot), (\cdot, y)\rangle$, the number of the blocks of the form $\langle(x, \cdot), (\cdot, y)\rangle$ in $T^{(q)}$, namely r , is $\text{Bin}(t, s/(k^{(p_1)} + k^{(p_2)} + k^{(p)} + 2s))$. Moreover, the probability of getting a block of this type is $q(x, y) = p_1(x) \cdot p_2(y)$, so t is a Poisson random variable with mean $q(x, y) \cdot (k^{(p_1)} + k^{(p_2)} + k^{(p)} + 2s)$ over the randomness of X . By Markov's inequality, with probability $1 - 1/c$, we may assume t is at most c times its expectation. As a consequence, $\mathbf{E}[r]$ is at most $c \cdot \mathbf{E}[t] \cdot s/(k^{(p_1)} + k^{(p_2)} + k^{(p)} + 2s) = c q(x, y) s$. Note that we can extend this example further to any type of block and any test or flattening sets.

723 Given X , for a domain element (x, y) , the following holds using Jensen's inequality:

$$\begin{aligned}
\mathbf{E}_\pi \left[\frac{1}{b_{(x,y)}(X, \pi)} \right] &= \mathbf{E}_{k_x^{(p_1)}, k_y^{(p_2)}, k_{(x,y)}^{(p)}} \left[\frac{1}{b_{(x,y)}} \right] \geq \frac{1}{\mathbf{E}_{k_x^{(p_1)}, k_y^{(p_2)}, k_{(x,y)}^{(p)}} [b_{(x,y)}]} \\
&= \frac{1}{(\mathbf{E} [k_x^{(p_1)}] + 1) \cdot (\mathbf{E} [k_y^{(p_2)}] + 1) + \mathbf{E} [k_{(x,y)}^{(p)}]} \\
&\geq \frac{1}{(50 p_1(x) k^{(p_1)} + 1) \cdot (50 p_2(y) k^{(p_2)} + 1) + 100 p(x, y) k^{(p)}} \\
&\geq \frac{1}{2500} \cdot \frac{1}{(p_1(x) k^{(p_1)} + 1) \cdot (p_2(y) k^{(p_2)} + 1) + p(x, y) k^{(p)}}.
\end{aligned}$$

724 where the second to last inequality holds with probability 0.95.

725 On the other hand, we find an upper bound of $\mathbf{E}[1/b_{(x,y)}]$ over the randomness of all variables.
726 Note that now that X has not been observed, so the number of each block type in each set is an
727 Poisson random variable, and it is independent from the rest. Let F denote the set of all flattening
728 blocks. We denote $p_1(x) k^{(p_1)}$, $p_2(y) k^{(p_2)}$, and $p(x, y) k^{(p)}$ by λ_1 , λ_2 , and λ_3 respectively. The
729 expected value of $1/b_{(x,y)}$ can be rewritten as:

$$\begin{aligned}
\mathbf{E}_F \left[\frac{1}{b_{(x,y)}} \right] &= \mathbf{E}_F \left[\frac{1}{(k_x^{(p_1)} + 1)(k_y^{(p_2)} + 1) + k_{(x,y)}^{(p)}} \right] \\
&\leq \mathbf{E}_F \left[\min \left(\frac{1}{(k_x^{(p_1)} + 1)(k_y^{(p_2)} + 1)}, \frac{1}{k_{(x,y)}^{(p)} + 1} \right) \right] \\
&\leq \min \left(\mathbf{E}_F \left[\frac{1}{k_x^{(p_1)} + 1} \right] \cdot \mathbf{E}_F \left[\frac{1}{k_y^{(p_2)} + 1} \right], \mathbf{E}_F \left[\frac{1}{k_{(x,y)}^{(p)} + 1} \right] \right) \\
&\leq \min \left(\frac{1 - e^{-\lambda_1}}{\lambda_1} \cdot \frac{1 - e^{-\lambda_2}}{\lambda_2}, \frac{1 - e^{-\lambda_3}}{\lambda_3} \right) \leq \min \left(\frac{4}{(\lambda_1 + 1)(\lambda_2 + 1)}, \frac{2}{\lambda_3 + 1} \right) \\
&\leq \frac{8}{(\lambda_1 + 1)(\lambda_2 + 1) + \lambda_3} \\
&= \frac{8}{(p_1(x) k^{(p_1)} + 1) \cdot (p_2(y) k^{(p_2)} + 1) + p(x, y) k^{(p)}}
\end{aligned}$$

730 Governed by the previous equations, we obtain:

$$\Pr_X \left[\mathbf{E}_\pi \left[\frac{1}{b_{(x,y)}(X, \pi)} \right] \geq \frac{1}{20000} \cdot \mathbf{E}_F \left[\frac{1}{b_{(x,y)}} \right] \right] \geq 0.96,$$

731 which is equivalent to

$$\Pr_X \left[\mathbf{E}_\pi \left[\frac{(p(x, y) - q(x, y))^2}{b_{(x,y)}(X, \pi)} \right] < \frac{1}{20000} \cdot \mathbf{E}_F \left[\frac{(p(x, y) - q(x, y))^2}{b_{(x,y)}} \right] \right] \leq 0.04.$$

732 To prove Equation 4, we need to show that the above equation holds even for the sum of the quantities
733 over all (x, y) with high probability. We show the claim in the following lemma. The proof is in
734 Section D:

Lemma C.4. *Let x_1, x_2, \dots, x_n be n non-negative random variables. Suppose there exist two constants c and p , both at most one, such that for each random variable x_i , we have:*

$$\Pr[x_i < c \cdot \mathbf{E}[x_i]] \leq p,$$

Then, one can show:

$$\Pr \left[\sum_{i=1}^n x_i < \frac{c \cdot \sum_{i=1}^n \mathbf{E}[x_i]}{10} \right] \leq \frac{10p}{9}.$$

Now, we focus on proving Equation 5. To prove the inequality, it suffices to show that $\mathbf{E}_F \left[1/b_{(x,y)}^3 \right]$ is $O \left(\mathbf{E}_F \left[1/b_{(x,y)}^2 \right]^2 \right)$. Again, note that we can think of $b_{(x,y)}$ to be equal to $(X+1)(Y+1)+Z$ where X , Y and Z are three Poisson random variables with means $\lambda_1 = p_1(x)k^{(p_1)}$, $\lambda_2 = p_2(y)k^{(p_2)}$, and $\lambda_3 =$ respectively. In Equation 12 we show that:

$$\mathbf{E}_{W \sim \text{Poi}(\lambda)} \left[\frac{1}{W^3} \right] \leq 6 \cdot \left(\frac{1 - e^{-\lambda}}{\lambda} \right)^2$$

Thus, we obtain an upper bound for the expected value of $1/b_{(x,y)}^3$ as follows:

$$\begin{aligned} \mathbf{E}_F \left[\frac{1}{b_{(x,y)}^3} \right] &= \mathbf{E}_{X,Y,Z} \left[\frac{1}{((X+1) \cdot (Y+1) + Z)^3} \right] \leq \mathbf{E}_{X,Y,Z} \left[\min \left(\frac{1}{(X+1)^3} \cdot \frac{1}{(Y+1)^3}, \frac{1}{(Z+1)^3} \right) \right] \\ &\leq \min \left(\mathbf{E}_X \left[\frac{1}{(X+1)^3} \right] \cdot \mathbf{E}_Y \left[\frac{1}{(Y+1)^3} \right], \mathbf{E}_Z \left[\frac{1}{(Z+1)^3} \right] \right) \\ &\leq 36 \min \left(\left(\frac{1 - e^{-\lambda_1}}{\lambda_1} \right)^2 \cdot \left(\frac{1 - e^{-\lambda_2}}{\lambda_2} \right)^2, \left(\frac{1 - e^{-\lambda_3}}{\lambda_3} \right)^2 \right). \end{aligned}$$

Note that in the case that one of the λ 's is equal to zero, one can replace $1 - e^{-\lambda}/\lambda$ by one in the rest of the proof. On the other hand, we can find a lower bound for $1/b_{(x,y)}$ by Jensen's inequality:

$$\begin{aligned} \left(\mathbf{E}_F \left[\frac{1}{b_{(x,y)}} \right] \right)^2 &\geq \left(\frac{1}{\mathbf{E}_F[b_i]} \right)^2 = \left(\frac{1}{(\lambda_1 + 1)(\lambda_2 + 1) + \lambda_3} \right)^2 \\ &\geq \left(\frac{1}{2} \min \left(\frac{1}{\lambda_1 + 1} \cdot \frac{1}{\lambda_2 + 1}, \frac{1}{\lambda_3 + 1} \right) \right)^2 \\ &\geq \frac{1}{4} \min \left(\left(\frac{1}{\lambda_1 + 1} \cdot \frac{1}{\lambda_2 + 1} \right)^2, \left(\frac{1}{\lambda_3 + 1} \right)^2 \right) \\ &\geq \frac{1}{64} \min \left(\left(\frac{1 - e^{-\lambda_1}}{\lambda_1} \cdot \frac{1 - e^{-\lambda_2}}{\lambda_2} \right)^2, \left(\frac{1 - e^{-\lambda_3}}{\lambda_3} \right)^2 \right) \end{aligned}$$

where the last inequality is due to the fact that $(1 - e^{-t})/t$ is at most $2/(t+1)$ for a non-negative number t . Putting these two equations together, one can conclude Equation 5:

$$\begin{aligned} \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_4^4 \right] &= \sum_{x=1}^n \sum_{y=1}^m \mathbf{E}_F \left[\frac{(p(x,y) - q(x,y))^4}{b_{(x,y)}^3} \right] \leq 36 \cdot 64 \cdot \sum_{x=1}^n \sum_{y=1}^m \left(\mathbf{E}_F \left[\frac{(p(x,y) - q(x,y))^2}{b_{x,y}} \right] \right)^2 \\ &\leq 2304 \cdot \left(\sum_{x=1}^n \sum_{y=1}^m \mathbf{E}_F \left[\frac{(p(x,y) - q(x,y))^2}{b_{x,y}} \right] \right)^2 = 2304 \cdot \mathbf{E}_F \left[\|p^{(F)} - q^{(F)}\|_2^2 \right]^2. \end{aligned}$$

Therefore, the statement of the lemma is concluded. \square

3

C.2 sensitivity of the statistic for the independence problem

In this section, we give an upper bound for the sensitivity of the independence statistic: the amount that the statistic changes if we change one sample in the input.

Let X denote a set of block that the algorithm received as the input. Assume we permute the blocks in X using a permutation π . Note that if we fix the size of each flattening set and sample set, \hat{s}_1, \hat{s}_2 , etc., one can deterministically find $s_{(x,y)}^{(p)}, s_{(x,y)}^{(q)}$, etc.. Thus, given X , π , and sizes of sets, one can compute the following statistic:

$$Z(X, \pi) := \sum_{x=1}^m \sum_{y=1}^n \frac{(s_{(x,y)}^{(p)} - s_{(x,y)}^{(q)})^2 - s_{(x,y)}^{(p)} - s_{(x,y)}^{(q)}}{(k_x^{(p_1)} + 1)(k_y^{(p_2)} + 1) + k_{(x,y)}^{(p)} + k_{(x,y)}^{(q)}}$$

We denote the average of Z over all π by $\overline{Z}(X)$. Our goal here is to calculate

$$\Delta Z = \max X, X' |\overline{Z}(X) - \overline{Z}(X')|$$

where X and X' are two neighboring data sets that they differ in exactly one element.

Through this section, we use an important property of Poissonization method: Let A and B be two sets with $\hat{n}_1 = \text{Poi}(n_1)$ and $\hat{n}_2 = \text{Poi}(n_2)$ samples. Given that there are k instance of element i in A and B together, the number of instances of element i in A is a Binomial random variable: $\text{Bin}(\hat{n}_1 + \hat{n}_2, n_1/(n_1 + n_2))$.

Lemma C.3. *Given that the size of all flattening and test samples are within the constant factor of their expectations, the sensitivity of the statistic Z is bounded as follows:*

$$\Theta \left(\frac{s}{k^{(q)}} + \frac{s}{k^{(p)}} + \frac{s}{k^{(p)}} \cdot \frac{f_{\langle(.,b),(\cdot,\cdot)\rangle}}{f_{\langle(.,\cdot),(\cdot,b)\rangle} + 1} \right)$$

Proof: In this proof, we assume X and X' are fully given, and X and X' are only different in the r -th block of the samples. Note that when we permute the elements in X and X' , we only permute the blocks and do not change the order of the samples within each block. The expectations in the this proof are taken over the random choice of a permutation π . As we mentioned earlier, we partition the blocks into the following sets, and the number of occurrences of each block types in each set determines the statistic:

$$\mathcal{S} = \left\{ F^{(p_1)}, F^{(p_2)}, F^{(p)}, F^{(q)}, T^{(p)}, T^{(q)} \right\}$$

We can separate our calculation based on where the r -th block is:

$$\begin{aligned} \Delta(\overline{Z}) &= |\overline{Z}(X) - \overline{Z}(X')| \\ &= |\mathbf{E}_\pi[Z(X, \pi) - Z(X', \pi)]| \leq \mathbf{E}_\pi[|Z(X, \pi) - Z(X', \pi)|] \\ &= \sum_{S \in \mathcal{S}} \Pr_\pi[r \in S] \cdot |\mathbf{E}_\pi[|Z(X, \pi) - Z(X', \pi)| \mid r \in S]| \end{aligned}$$

Now, we consider each term separate.

1. **Block r is in $F^{(p_1)}$:** Suppose the types of the r -th block in X and X' are $\langle(a, \cdot), (\cdot, \cdot)\rangle$ and $\langle(a', \cdot), (\cdot, \cdot)\rangle$ respectively. If a and a' are equal, then the statistic will remains unchanged. Otherwise, $k_a^{(p_1)}$ and $k_{a'}^{(p_1)}$ is changed by one. First, we simplify the term $|Z(X, \pi) - Z(X', \pi)|$ for a given π :

$$\begin{aligned} |Z(X, \pi) - Z(X', \pi)| &= \left| \sum_{x=1}^n \sum_{y=1}^m \frac{(s_{(x,y)}^{(p)} - s_{(x,y)}^{(q)})^2 - s_{(x,y)}^{(p)} - s_{(x,y)}^{(q)}}{(k_x^{(p_1)} + 1)(k_y^{(p_2)} + 1) + k_{(x,y)}^{(p)} + k_{(x,y)}^{(q)}} \right. \\ &\quad \left. - \sum_{x=1}^n \sum_{y=1}^m \frac{(s'_{(x,y)}^{(p)} - s'_{(x,y)}^{(q)})^2 - s'_{(x,y)}^{(p)} - s'_{(x,y)}^{(q)}}{(k'_x{}^{(p_1)} + 1)(k'_y{}^{(p_2)} + 1) + k'_{(x,y)}{}^{(p)} + k'_{(x,y)}{}^{(q)}} \right| \\ &\leq \sum_{x \in \{a, a'\}} \left| \sum_{y=1}^m \frac{(s_{(x,y)}^{(p)} - s_{(x,y)}^{(q)})^2 - s_{(x,y)}^{(p)} - s_{(x,y)}^{(q)}}{(k_x^{(p_1)} + 1)(k_y^{(p_2)} + 1) + k_{(x,y)}^{(p)} + k_{(x,y)}^{(q)}} \right. \\ &\quad \left. - \frac{(s'_{(x,y)}^{(p)} - s'_{(x,y)}^{(q)})^2 - s'_{(x,y)}^{(p)} - s'_{(x,y)}^{(q)}}{(k'_x{}^{(p_1)} + 1)(k'_y{}^{(p_2)} + 1) + k'_{(x,y)}{}^{(p)} + k'_{(x,y)}{}^{(q)}} \right| \end{aligned}$$

For the rest of the proof, we focus on the term above when $x = a$. The other term can be upper bounded similarly, and at the end we multiply our final bound by two.

$$\begin{aligned}
& \left| \sum_{y=1}^m \frac{(s_{(a,y)}^{(p)} - s_{(a,y)}^{(q)})^2 - s_{(a,y)}^{(p)} - s_{(a,y)}^{(q)}}{(k_a^{(p_1)} + 1)(k_y^{(p_2)} + 1) + k_{(a,y)}^{(p)} + k_{(x,y)}^{(q)}} - \frac{(s'_{(a,y)}^{(p)} - s'_{(a,y)}^{(q)})^2 - s'_{(a,y)}^{(p)} - s'_{(a,y)}^{(q)}}{(k'_a(p_1) + 1)(k'_y(p_2) + 1) + k'_{(a,y)}^{(p)} + k'_{(x,y)}^{(q)}} \right| \\
& \leq \sum_{y=1}^m \left| \frac{(s_{(a,y)}^{(p)} - s_{(a,y)}^{(q)})^2 - s_{(a,y)}^{(p)} - s_{(a,y)}^{(q)}}{(k'_a(p_1) + 2)(k_y^{(p_2)} + 1) + k_{(a,y)}^{(p)} + k_{(x,y)}^{(q)}} - \frac{(s_{(a,y)}^{(p)} - s_{(a,y)}^{(q)})^2 - s_{(a,y)}^{(p)} - s_{(a,y)}^{(q)}}{(k'_a(p_1) + 1)(k_y^{(p_2)} + 1) + k_{(a,y)}^{(p)} + k_{(x,y)}^{(q)}} \right| \\
& \leq \sum_{y=1}^m \frac{(k_y^{(p_2)} + 1) \cdot |(s_{(a,y)}^{(p)} - s_{(a,y)}^{(q)})^2 - s_{(a,y)}^{(p)} - s_{(a,y)}^{(q)}|}{\left((k'_a(p_1) + 2)(k_y^{(p_2)} + 1) + k_{(a,y)}^{(p)} + k_{(x,y)}^{(q)} \right) \cdot \left((k'_a(p_1) + 1)(k_y^{(p_2)} + 1) + k_{(a,y)}^{(p)} + k_{(x,y)}^{(q)} \right)} \\
& \leq \sum_{y=1}^m \frac{(k_y^{(p_2)} + 1) \cdot ((s_{(a,y)}^{(p)})^2 + (s_{(a,y)}^{(q)})^2)}{\left((k'_a(p_1) + 2)(k_y^{(p_2)} + 1) + k_{(a,y)}^{(p)} + k_{(x,y)}^{(q)} \right) \cdot \left((k'_a(p_1) + 1)(k_y^{(p_2)} + 1) + k_{(a,y)}^{(p)} + k_{(x,y)}^{(q)} \right)} \\
& \leq \sum_{y=1}^m \frac{1}{k_a^{(p_1)} + 1} \cdot \frac{(s_{(a,y)}^{(p)})^2}{k_{(a,y)}^{(p)} + 1} + \frac{1}{k_a^{(p_1)} + 1} \cdot \frac{(s_{(a,y)}^{(q)})^2}{k_{(a,y)}^{(q)} + 1}
\end{aligned}$$

For brevity's sake, let v denote the following expectation:

$$v := \mathbf{E}_\pi \left[\sum_{y=1}^m \frac{1}{k_a^{(p_1)} + 1} \cdot \frac{(s_{(a,y)}^{(p)})^2}{k_{(a,y)}^{(p)} + 1} + \frac{1}{k_a^{(p_1)} + 1} \cdot \frac{(s_{(a,y)}^{(q)})^2}{k_{(a,y)}^{(q)} + 1} \middle| r \in F^{(p_1)} \right]$$

Using the tower rule, we achieve:

$$\begin{aligned}
v & \leq \mathbf{E}_{F^{(p_1)}} \left[\frac{1}{(k_a^{(p_1)} + 1)} \cdot \mathbf{E}_\pi \left[\sum_{y=1}^m \frac{(s_{(a,y)}^{(p)})^2}{(k_{(a,y)}^{(p)} + 1)} \middle| r \in F^{(p_1)}, F^{(p_1)} \right] \middle| r \in F^{(p_1)} \right] \\
& + \mathbf{E}_{F^{(p_1)}} \left[\frac{1}{k_a^{(p_1)} + 1} \cdot \mathbf{E}_\pi \left[\sum_{y=1}^m \frac{(s_{(a,y)}^{(q)})^2}{k_{(a,y)}^{(q)} + 1} \middle| r \in F^{(p_1)}, F^{(p_1)} \right] \middle| r \in F^{(p_1)} \right]
\end{aligned}$$

Let $f_{\langle(a,y),(\cdot,\cdot)\rangle}$, $f_{\langle(a,\cdot),(\cdot,y)\rangle}$, and $f_{\langle(a,\cdot),(\cdot,\cdot)\rangle}$ be the numbers of blocks of the forms $\langle(a,y),(\cdot,\cdot)\rangle$, $\langle(a,\cdot),(\cdot,y)\rangle$, and $\langle(a,\cdot),(\cdot,\cdot)\rangle$ in X respectively. Using Lemma D.4, one can bound the terms inside the expectations as below:

$$\begin{aligned}
\mathbf{E} \left[\frac{(s_{(a,y)}^{(p)})^2}{(k_{(a,y)}^{(p)} + 1)} \right] & \leq \min \left(\frac{2(s-1)f_{\langle(a,y),(\cdot,\cdot)\rangle}}{(k^{(p)} + 1)}, 2f_{\langle(a,y),(\cdot,\cdot)\rangle}^2 \right) + f_{\langle(a,y),(\cdot,\cdot)\rangle} \\
& \leq \min \left(\left(\frac{2(s-1)}{(k^{(p)} + 1)} + 1 \right) \cdot f_{\langle(a,y),(\cdot,\cdot)\rangle}, 3f_{\langle(a,y),(\cdot,\cdot)\rangle}^2 \right),
\end{aligned}$$

$$\begin{aligned}
\mathbf{E} \left[\frac{(s_{(a,y)}^{(q)})^2}{(k_{(a,y)}^{(q)} + 1)} \right] & \leq \min \left(\frac{2(s-1)f_{\langle(a,\cdot),(\cdot,y)\rangle}}{(k^{(q)} + 1)}, 2f_{\langle(a,\cdot),(\cdot,y)\rangle}^2 \right) + f_{\langle(a,\cdot),(\cdot,y)\rangle} \\
& \leq \min \left(\left(\frac{2(s-1)}{(k^{(q)} + 1)} + 1 \right) \cdot f_{\langle(a,\cdot),(\cdot,y)\rangle}, 3f_{\langle(a,\cdot),(\cdot,y)\rangle}^2 \right).
\end{aligned}$$

766
767

Observe that $\sum_y f_{\langle(a,y),(\cdot,\cdot)\rangle}$ and $\sum_y f_{\langle(a,\cdot),(\cdot,y)\rangle}$ are equal to $f_{\langle(a,\cdot),(\cdot,\cdot)\rangle}$. Thus, using Lemma D.5, and Lemma D.6, we have:

$$\begin{aligned}
v &\leq \mathbf{E}_{F^{(p_1)}} \left[\frac{1}{(k_a^{(p_1)} + 1)} \middle| r \in F^{(p_1)} \right] \cdot \sum_{y=1}^m \min \left(\left(\frac{2(s-1)}{(k^{(p)} + 1)} + 1 \right) \cdot f_{\langle(a,y),(\cdot,\cdot)\rangle}, 3f_{\langle(a,y),(\cdot,\cdot)\rangle}^2 \right) \\
&\quad + \mathbf{E}_{F^{(p_1)}} \left[\frac{1}{k_a^{(p_1)} + 1} \middle| r \in F^{(p_1)} \right] \cdot \sum_{y=1}^m \min \left(\left(\frac{2(s-1)}{(k^{(q)} + 1)} + 1 \right) \cdot f_{\langle(a,\cdot),(\cdot,y)\rangle}, 3f_{\langle(a,\cdot),(\cdot,y)\rangle}^2 \right) \\
&\leq \min \left(1, \frac{|X|}{f_{\langle(a,\cdot),(\cdot,\cdot)\rangle} k^{(p_1)}} \right) \cdot \left(\frac{2(s-1)}{(k^{(p)} + 1)} + 1 \right) \cdot f_{\langle(a,\cdot),(\cdot,\cdot)\rangle} \\
&\quad + \min \left(1, \frac{|X|}{f_{\langle(a,\cdot),(\cdot,\cdot)\rangle} k^{(p_1)}} \right) \cdot \left(\frac{2(s-1)}{(k^{(q)} + 1)} + 1 \right) \cdot f_{\langle(a,\cdot),(\cdot,\cdot)\rangle} \\
&\leq \Theta \left(\frac{|X|}{k^{(p_1)}} \cdot \left(\frac{s}{k^{(p)}} + \frac{s}{k^{(q)}} \right) \right)
\end{aligned}$$

768

Using the above calculation, it is not hard to see that the following holds

$$\begin{aligned}
&\Pr_{\pi} \left[r \in F^{(p_1)} \right] \cdot \left| \mathbf{E}_{\pi} \left[|Z(X, \pi) - Z(X', \pi)| \middle| r \in F^{(p_1)} \right] \right| \\
&\leq 2 \cdot \frac{k^{(p_1)}}{|X|} \cdot \left(\frac{|X|}{k^{(p_1)}} \cdot \left(\frac{s}{k^{(p)}} + \frac{s}{k^{(q)}} \right) \right) \\
&\leq \Theta \left(\frac{s}{k^{(p)}} + \frac{s}{k^{(q)}} \right).
\end{aligned}$$

769
770
771

2. **Block r is in $F^{(p_2)}$:** Suppose the the r -th block in X and X' are of the forms $\langle(\cdot, \cdot), (\cdot, b)\rangle$ and $\langle(\cdot, \cdot), (\cdot, b')\rangle$ respectively. Using the symmetry of this case and the previous case, we take the same approach.

$$\begin{aligned}
&\mathbf{E}_{\pi} \left[\sum_{x=1}^n \frac{1}{k_b^{(p_2)} + 1} \cdot \frac{(s_{(x,b)}^{(p)})^2}{k_{(x,b)}^{(p)} + 1} + \frac{1}{k_b^{(p_2)} + 1} \cdot \frac{(s_{(x,b)}^{(q)})^2}{k_{(x,b)}^{(q)} + 1} \middle| r \in F^{(p_2)} \right] \\
&\leq \mathbf{E}_{F^{(p_2)}} \left[\frac{1}{k_b^{(p_2)} + 1} \cdot \mathbf{E}_{\pi} \left[\sum_{x=1}^n \frac{(s_{(x,b)}^{(p)})^2}{k_{(x,b)}^{(p)} + 1} \middle| r \in F^{(p_2)}, F^{(p_2)} \right] \middle| r \in F^{(p_2)} \right] \\
&\quad + \mathbf{E}_{F^{(p_2)}} \left[\frac{1}{k_b^{(p_2)} + 1} \cdot \mathbf{E}_{\pi} \left[\sum_{x=1}^n \frac{(s_{(x,b)}^{(q)})^2}{k_{(x,b)}^{(q)} + 1} \middle| r \in F^{(p_2)}, F^{(p_2)} \right] \middle| r \in F^{(p_2)} \right] \\
&\leq \mathbf{E}_{F^{(p_2)}} \left[\frac{1}{k_b^{(p_2)} + 1} \middle| r \in F^{(p_2)} \right] \cdot \sum_{x=1}^n \min \left(\left(\frac{2(s-1)}{(k^{(p)} + 1)} + 1 \right) \cdot f_{\langle(x,b),(\cdot,\cdot)\rangle}, 3f_{\langle(x,b),(\cdot,\cdot)\rangle}^2 \right) \\
&\quad + \mathbf{E}_{F^{(p_2)}} \left[\frac{1}{k_b^{(p_2)} + 1} \middle| r \in F^{(p_2)} \right] \cdot \sum_{x=1}^n \min \left(\left(\frac{2(s-1)}{(k^{(q)} + 1)} + 1 \right) \cdot f_{\langle(x,\cdot),(\cdot,b)\rangle}, 3f_{\langle(x,\cdot),(\cdot,b)\rangle}^2 \right) \\
&\leq \min \left(1, \frac{|X|}{f_{\langle(\cdot,\cdot),(\cdot,b)\rangle} k^{(p_2)}} \right) \cdot \left(\frac{2(s-1)}{(k^{(p)} + 1)} + 1 \right) \cdot f_{\langle(\cdot,b),(\cdot,\cdot)\rangle} \\
&\quad + \min \left(1, \frac{|X|}{f_{\langle(\cdot,\cdot),(\cdot,b)\rangle} k^{(p_2)}} \right) \cdot \left(\frac{2(s-1)}{(k^{(q)} + 1)} + 1 \right) \cdot f_{\langle(\cdot,\cdot),(\cdot,b)\rangle} \\
&\leq \frac{|X|}{k^{(p_2)}} \cdot \left(\frac{2(s-1)}{(k^{(p)} + 1)} + 1 \right) \cdot \frac{f_{\langle(\cdot,b),(\cdot,\cdot)\rangle}}{f_{\langle(\cdot,\cdot),(\cdot,b)\rangle} + 1} + \frac{|X|}{k^{(p_2)}} \cdot \left(\frac{2(s-1)}{(k^{(q)} + 1)} + 1 \right) \\
&\leq \Theta \left(\frac{|X|}{k^{(p_2)}} \cdot \left(\frac{s}{k^{(p)}} \cdot \frac{f_{\langle(\cdot,b),(\cdot,\cdot)\rangle}}{f_{\langle(\cdot,\cdot),(\cdot,b)\rangle} + 1} + \frac{s}{k^{(q)}} \right) \right)
\end{aligned}$$

772

where the last inequality is due to the fact that $\min_{x>0}(\alpha/x, x) \leq \sqrt{\alpha}$.

$$\begin{aligned} & \Pr_{\pi} \left[r \in F^{(p_2)} \right] \cdot \left| \mathbf{E}_{\pi} \left[|Z(X, \pi) - Z(X', \pi)| \mid r \in F^{(p_2)} \right] \right| \\ & \leq 2 \cdot \frac{k^{(p_2)}}{|X|} \cdot \Theta \left(\frac{|X|}{k^{(p_2)}} \cdot \left(\frac{s}{k^{(p)}} \cdot \frac{f_{\langle(.,b),(\cdot,\cdot)\rangle}}{f_{\langle(.,\cdot),(\cdot,b)\rangle} + 1} + \frac{s}{k^{(q)}} \right) \right) \\ & \leq \Theta \left(\frac{s}{k^{(p)}} \cdot \frac{f_{\langle(.,b),(\cdot,\cdot)\rangle}}{f_{\langle(.,\cdot),(\cdot,b)\rangle} + 1} + \frac{s}{k^{(q)}} \right) \end{aligned}$$

773

3. **Block r is in $F^{(p)}$ or $F^{(q)}$:** Here we assume r is in $F^{(p)}$. Very similar calculation, yield the same bound if r is in $F^{(q)}$. Suppose the the r -th block in X and X' are of the forms $\langle(a, b), (\cdot, \cdot)\rangle$ and $\langle(a', b'), (\cdot, \cdot)\rangle$ respectively. Note that in this case, only two terms will be different, so we have:

774

775

776

$$\begin{aligned} |Z(X, \pi) - Z(X', \pi)| \leq & \sum_{\substack{(x,y) \in \\ \{(a,b), (a',b')\}}} \left| \frac{(s_{(x,y)}^{(p)} - s_{(x,y)}^{(q)})^2 - s_{(x,y)}^{(p)} - s_{(x,y)}^{(q)}}{(k_x^{(p_1)} + 1)(k_y^{(p_2)} + 1) + k_{(x,y)}^{(p)} + k_{(x,y)}^{(q)}} \right. \\ & \left. - \frac{(s'_{(x,y)}^{(p)} - s'_{(x,y)}^{(q)})^2 - s'_{(x,y)}^{(p)} - s'_{(x,y)}^{(q)}}{(k'_x{}^{(p_1)} + 1)(k'_y{}^{(p_2)} + 1) + k'_{(x,y)}{}^{(p)} + k'_{(x,y)}{}^{(q)}} \right| \end{aligned}$$

777

Now, we focus on the term where (x, y) is equal to (a, b) . The other term can be bounded similarly.

778

$$\begin{aligned} & \left| \frac{(s_{(a,b)}^{(p)} - s_{(a,b)}^{(q)})^2 - s_{(a,b)}^{(p)} - s_{(a,b)}^{(q)}}{(k_a^{(p_1)} + 1)(k_b^{(p_2)} + 1) + k_{(a,b)}^{(p)} + k_{(a,b)}^{(q)}} - \frac{(s'_{(a,b)}^{(p)} - s'_{(a,b)}^{(q)})^2 - s'_{(a,b)}^{(p)} - s'_{(a,b)}^{(q)}}{(k'_a{}^{(p_1)} + 1)(k'_b{}^{(p_2)} + 1) + k'_{(a,b)}{}^{(p)} + k'_{(a,b)}{}^{(q)}} \right| \\ & = \left| \frac{(s_{(a,b)}^{(p)} - s_{(a,b)}^{(q)})^2 - s_{(a,b)}^{(p)} - s_{(a,b)}^{(q)}}{(k_a^{(p_1)} + 1)(k_b^{(p_2)} + 1) + k_{(a,b)}^{(p)} + k_{(a,b)}^{(q)}} - \frac{(s_{(a,b)}^{(p)} - s_{(a,b)}^{(q)})^2 - s_{(a,b)}^{(p)} - s_{(a,b)}^{(q)}}{(k_a^{(p_1)} + 1)(k_b^{(p_2)} + 1) + k_{(a,b)}^{(p)} + k_{(a,b)}^{(q)} + 1} \right| \\ & \leq \frac{(s_{(a,b)}^{(p)})^2 + (s_{(a,b)}^{(q)})^2}{\left((k_a^{(p_1)} + 1)(k_b^{(p_2)} + 1) + k_{(a,b)}^{(p)} + k_{(a,b)}^{(q)} \right) \cdot \left((k_a^{(p_1)} + 1)(k_b^{(p_2)} + 1) + k_{(a,b)}^{(p)} + k_{(a,b)}^{(q)} + 1 \right)} \\ & \leq \frac{(s_{(a,b)}^{(p)})^2}{k_{(a,b)}^{(p)} (k_{(a,b)}^{(p)} + 1)} + \frac{(s_{(a,b)}^{(q)})^2}{k_a^{(q)} (k_a^{(q)} + 1)} \end{aligned}$$

779

Now, using Lemma D.6, we bound the expected value of the above quantity from above:

$$\begin{aligned} & \mathbf{E}_{\pi} \left[\frac{(s_{(a,b)}^{(p)})^2}{k_{(a,b)}^{(p)} (k_{(a,b)}^{(p)} + 1)} + \frac{(s_{(a,b)}^{(q)})^2}{k_a^{(q)} (k_a^{(q)} + 1)} \mid r \in F^{(p)} \right] \\ & \leq f_{\langle(a,b),(\cdot,\cdot)\rangle}^2 \cdot \mathbf{E}_{\pi} \left[\frac{1}{k_{(a,b)}^{(p)} (k_{(a,b)}^{(p)} + 1)} \mid r \in F^{(p)} \right] \\ & \quad + f_{\langle(a,\cdot),(\cdot,b)\rangle}^2 \cdot \mathbf{E}_{\pi} \left[\frac{1}{k_a^{(q)} (k_a^{(q)} + 1)} \mid r \in F^{(p)} \right] \\ & \leq \frac{|X|(|X| + 1)}{k^{(p)} (k^{(p)} + 1)} + \frac{|X|(|X| + 1)}{k^{(q)} (k^{(q)} + 1)} \end{aligned}$$

780

Using the above equation, and the fact that $|X| = \Theta(s)$, it is not hard to see that

$$\begin{aligned} & \Pr_{\pi} \left[r \in F^{(p)} \right] \cdot \left| \mathbf{E}_{\pi} \left[|Z(X, \pi) - Z(X', \pi)| \mid r \in F^{(p)} \right] \right| \\ & \leq 2 \cdot \frac{k^{(p)}}{|X|} \cdot \Theta \left(\frac{|X|(|X|+1)}{k^{(p)}(k^{(p)}+1)} + \frac{|X|(|X|+1)}{k^{(q)}(k^{(q)}+1)} \right) \\ & \leq \Theta \left(\frac{s}{k^{(p)}} + \frac{s}{k^{(q)}} \right) \end{aligned}$$

781

Note that the factor of two in the above inequality, comes from including the symmetric term for (a', b') .

782

783

4. **Block r is in $T^{(p)}$ or $T^{(q)}$:** Suppose the the r -th block in X and X' are of the forms $\langle (a, b), (\cdot, \cdot) \rangle$ and $\langle (a', b'), (\cdot, \cdot) \rangle$ respectively. Note that in this case, only two terms will be different for the two datasets, so we have:

784

785

$$\begin{aligned} |Z(X, \pi) - Z(X', \pi)| \leq & \sum_{\substack{(x,y) \in \\ \{(a,b), (a',b')\}}} \left| \frac{(s_{(x,y)}^{(p)} - s_{(x,y)}^{(q)})^2 - s_{(x,y)}^{(p)} - s_{(x,y)}^{(q)}}{(k_x^{(p_1)} + 1)(k_y^{(p_2)} + 1) + k_{(x,y)}^{(p)} + k_{(x,y)}^{(q)}} \right. \\ & \left. - \frac{(s'_{(x,y)}^{(p)} - s'_{(x,y)}^{(q)})^2 - s'_{(x,y)}^{(p)} - s'_{(x,y)}^{(q)}}{(k'_x{}^{(p_1)} + 1)(k'_y{}^{(p_2)} + 1) + k'_{(x,y)}{}^{(p)} + k'_{(x,y)}{}^{(q)}} \right| \end{aligned}$$

786

Below, we assume r is in $T^{(p)}$. However, the calculation will be the same if r was in $T^{(q)}$. Now, we focus on the term where (x, y) is equal to (a, b) . The other term can be bounded similarly.

787

788

$$\begin{aligned} & \left| \frac{(s_{(a,b)}^{(p)} - s_{(a,b)}^{(q)})^2 - s_{(a,b)}^{(p)} - s_{(a,b)}^{(q)}}{(k_a^{(p_1)} + 1)(k_b^{(p_2)} + 1) + k_{(a,b)}^{(p)} + k_{(a,b)}^{(q)}} - \frac{(s'_{(a,b)}^{(p)} - s'_{(a,b)}^{(q)})^2 - s'_{(a,b)}^{(p)} - s'_{(a,b)}^{(q)}}{(k'_a{}^{(p_1)} + 1)(k'_b{}^{(p_2)} + 1) + k'_{(a,b)}{}^{(p)} + k'_{(a,b)}{}^{(q)}} \right| \\ & = \left| \frac{(s_{(a,b)}^{(p)} - s_{(a,b)}^{(q)})^2 - s_{(a,b)}^{(p)} - s_{(a,b)}^{(q)}}{(k_a^{(p_1)} + 1)(k_b^{(p_2)} + 1) + k_{(a,b)}^{(p)} + k_{(a,b)}^{(q)}} - \frac{(s_{(a,b)}^{(p)} - 1 - s_{(a,b)}^{(q)})^2 - (s_{(a,b)}^{(p)} - 1) - s_{(a,b)}^{(q)}}{(k_a^{(p_1)} + 1)(k_b^{(p_2)} + 1) + k_{(a,b)}^{(p)} + k_{(a,b)}^{(q)}} \right| \\ & \leq \frac{2s_{(a,b)}^{(p)} + 2s_{(a,b)}^{(q)}}{(k_a^{(p_1)} + 1)(k_b^{(p_2)} + 1) + k_{(a,b)}^{(p)} + k_{(a,b)}^{(q)}} \leq \frac{2s_{(a,b)}^{(p)}}{k_{(a,b)}^{(p)} + 1} + \frac{2s_{(a,b)}^{(q)}}{k_a^{(q)} + 1} \end{aligned}$$

789

Now, using Lemma D.5, we bound the expected value of the above quantity from above:

$$\begin{aligned} & \mathbf{E}_{\pi} \left[\frac{2s_{(a,b)}^{(p)}}{k_{(a,b)}^{(p)} + 1} + \frac{2s_{(a,b)}^{(q)}}{k_a^{(q)} + 1} \mid r \in F^{(p)} \right] \\ & \leq 2 f_{\langle (a,b), (\cdot, \cdot) \rangle} \cdot \mathbf{E}_{\pi} \left[\frac{1}{k_{(a,b)}^{(p)} + 1} \mid r \in F^{(p)} \right] \\ & \quad + 2 f_{\langle (a, \cdot), (\cdot, b) \rangle} \cdot \mathbf{E}_{\pi} \left[\frac{1}{k_a^{(q)} + 1} \mid r \in F^{(p)} \right] \\ & \leq \frac{|X|}{k^{(p)} + 1} + \frac{|X|}{k^{(q)} + 1} \end{aligned}$$

790

Using the above equation, and the fact that $|X| = \Theta(s)$, it is not hard to see that

$$\begin{aligned} & \Pr_{\pi} \left[r \in F^{(p)} \right] \cdot \left| \mathbf{E}_{\pi} \left[|Z(X, \pi) - Z(X', \pi)| \mid r \in F^{(p)} \right] \right| \\ & \leq 2 \cdot \frac{k^{(p)}}{|X|} \cdot \Theta \left(\frac{|X|}{k^{(p)} + 1} + \frac{|X|}{k^{(q)} + 1} \right) \leq \Theta(1) \end{aligned}$$

791

Note that the factor of two in the above inequality, comes from including the symmetric term for (a', b') .

792

793 Putting all the terms computed above together, we obtain:

$$\begin{aligned} |\bar{Z}(X) - \bar{Z}(X')| &= \sum_{S \in \mathcal{S}} \Pr_{\pi}[r \in S] \cdot |\mathbf{E}_{\pi}[|Z(X, \pi) - Z(X', \pi)| \mid r \in S]| \\ &\leq \Theta \left(\frac{s}{k(q)} + \frac{s}{k(p)} + \frac{s}{k(p)} \cdot \frac{f_{\langle(\cdot, b), (\cdot, \cdot)\rangle}}{f_{\langle(\cdot, \cdot), (\cdot, b)\rangle} + 1} \right) \end{aligned}$$

794

□

795 C.3 Stretching the domain of a private algorithm

In this section, we investigate whether we can extend the domain of a differentially private tester under certain conditions. We start off by defining the domains. The input of a differential private tester is a sample set from a universe Ω . Suppose we have a dataset X of $2s$ samples from $[n]$ that are arranged in two rows each of size s (namely top and bottom rows). Let the domain of a differential algorithm, denoted by \mathcal{X} , be the set of all such pairs of rows, namely $[n]^{2s}$. We denote the frequency of an element $i \in [n]$ in the top row by $t_i(X)$, and in the bottom row by $b_i(X)$. A desired property of a dataset is that the ratio of the frequencies in the row is bounded by a fixed parameter $A \geq 2$. More precisely, we define the subset of \mathcal{X} which contains the data sets with this property as:

$$\mathcal{X}^* = \left\{ X \mid \forall i \in [n] : \frac{t_i(X)}{b_i(X) + 1} \leq A \right\}$$

796 Let \mathcal{A} be a tester that receives a set of samples, X , as its input, and outputs $\mathcal{A}(X)$ which is known
797 to be correct with probability at least $1 - \delta$. Suppose \mathcal{A} is ξ -differentially private when X is in \mathcal{X}^* .
798 Our goal here is to design a $\Theta(\xi)$ -differentially private algorithm, namely \mathcal{B} , that takes $X \in \mathcal{X}$ as
799 its input, and outputs $\mathcal{B}(X)$ which is incorrect with probability slightly larger than δ .

800 At a high level, we implement \mathcal{B} by using \mathcal{A} as a blackbox as follows: We first look at the input
801 $X \in \mathcal{X}$. If X is already in \mathcal{X}^* , we output $\mathcal{A}(X)$. Otherwise, if X is in $\mathcal{X} \setminus \mathcal{X}^*$, we pass X through
802 a “filter” and turn it into another dataset Y which is in \mathcal{X}^* . Then, we output $\mathcal{A}(Y)$.

803 To show that \mathcal{B} is the desired algorithm, we have few challenges: (1) we need to show the mapping
804 does not affect the correctness probability by too much. (2) \mathcal{B} is $\Theta(\xi)$ -differentially private although
805 its input may be from $\mathcal{X} \setminus \mathcal{X}^*$. Overcoming the second challenge is closely related to the design
806 of the mapping. If two datasets have Hamming distance one, then we need to make sure they will
807 remain “close”. In the following section, we explain the mapping, and in the next section, we prove
808 that \mathcal{B} is a ξ -differentially private algorithm with large correctness probability.

809 C.4 Mapping datasets in $\mathcal{X} \setminus \mathcal{X}^*$ to datasets in \mathcal{X}^*

810 In this section, we provide a randomized mapping that takes $X \in \mathcal{X}$ as the input, and maps it to
811 randomly selected Y in \mathcal{X}^* with two important properties stated in the following Lemma:

812 **Lemma C.5.** *There exists a randomized mapping that takes $X, X' \in \mathcal{X}$ and maps them to $Y, Y' \in$
813 \mathcal{X}^* respectively with the following property:*

- 814 • *If X is in \mathcal{X}^* , then it will always be mapped to itself. : $Y = X$.*
- 815 • *If the Hamming distance between X and X' is one, then there exists a coupling \mathcal{C} between*
816 *the random outputs of the mapping, Y and Y' , where for any (Y, Y') drawn from \mathcal{C} , the*
817 *Hamming distance between Y and Y' is at most a constant $c = 4$ (independent of A).*

Proof: The main idea is to decrease the ratio $t_i(X)/(b_i(X) + 1)$ by replacing a subset of samples in the bottom row with the copies of i to decrease the ratio without introducing new elements that violate the ratio condition. For a dataset X , we look at each element $i \in [n]$, and see how many copies of i are needed to “fix” the ratio. It is not hard to see that if for each element i , $r_i(X)$ many copies is sufficient where $r_i(X)$ is defined as below:

$$r_i(X) := \max \left(\left\lceil \frac{t_i(X)}{A} \right\rceil - b_i(X) - 1, 0 \right).$$

Let R be a multiset that contains $r_i(X)$ copies of i . We find $|R|$ slots in the bottom row, and replace the samples in those slots with an element in R . If we carefully select the slots and do not replace any copy of i in the bottom row, the new ratio will be: $t_i(X)/(b_i(X) + r_i(X) + 1)$ which is at most A . Now, we focus on finding the slots in the bottom row. We can select a slot containing an instance of i , only if the replacement of i does not increase the ratio of the frequencies above A . For an element i , we may remove at most $s_i(X)$ samples where

$$s_i(X) := \max \left(b_i(X) + 1 - \left\lceil \frac{t_i(X)}{A} \right\rceil, 0 \right).$$

818 For each element i , we mark $s_i(X)$ many slots which contains copy of i in the bottom row as
 819 “available” preferring the slots with the smaller index. Observe that we always have at least $|R|$
 820 many slots since A is at least two:

$$\begin{aligned} |R| &= \sum_{i=1}^n r_i(X) \leq \sum_{i=1}^n \frac{t_i(X)}{A} \leq \frac{s}{A} \leq (A-1) \frac{s}{A} = s - \sum_{i=1}^n \frac{t_i(X)}{A} \leq s - \sum_{i=1}^n \left(\left\lceil \frac{t_i(X)}{A} \right\rceil - 1 \right) \\ &\leq s - \sum_{i=1}^n b_i(X) - s_i(X) = \sum_{i=1}^n s_i(X) \end{aligned}$$

821 We choose the first $|R|$ available slots (i.e. with the smaller indices), and replace the bottom samples
 822 in them by the samples in R randomly. After the replacements, it is clear that we did not remove a
 823 sample where its ratio could go above A , and we fixed all those elements with the ratio above A as
 824 well. Thus, the dataset we get after this process is surely in \mathcal{X}^* . Furthermore, if X is already in \mathcal{X}^* ,
 825 then R is an empty set, and the mapping does not change it, so $Y = X$.

826 Now, we focus on the proof of the existence of the coupling. Let S be the indices of the $|R|$ available
 827 slots we select. First note that we consider all the elements in R to be distinct. (even though they
 828 might be different copies of the same sample, we can index $r_i(X)$ copies of i by $1, 2, \dots, r_i(X)$.)
 829 Thus, there are $|R|!$ for assigning the samples in R into the slots in S , and each assigning has
 830 probability $1/|R|!$. Suppose two datasets, X and X' , differ in exactly one sample: X has an extra
 831 copy of i , and X' instead has an extra copy of j . Also, let R' and S' be the equivalents of R and S
 832 respectively for X' . Clearly, we have $|R| = |S|$, and $|R'| = |S'|$. This discrepancy between X and
 833 X' happens in either on the top row or the bottom row. Since the frequency of i and j changes by
 834 at most one, $r_i(X)$, $s_i(X)$, $r_j(X)$, and $s_j(X)$ will change by at most. Without loss of generality, if
 835 we consider all possible cases, it is not hard to see that one of the two following cases happens:

836 **Case 1:** R and R' has the same size, and $|R \cap R'|$ and $|S \cap S'|$ is at least $|R| - 1$. It is not hard to see
 837 that there is a bijection between Y and Y' . Assume there exists a set of replacement that
 838 turns X into Y . We construct the corresponding Y' accordingly. We start off with X' . We
 839 apply the same set of replacements with only two exceptions: Suppose we want to replace
 840 the sample in the slot ℓ with k according to the original set of replacement, then we see if
 841 k is not in R' , we carry on the replacement with $k' = R' \setminus k$. Also, if the ℓ is not in S' ,
 842 we will choose slot $\ell' = S' \setminus S$, pick the slot ℓ' for the replacement. After performing all
 843 the replacement we get Y' which has Hamming distance at most four to Y . It is not hard
 844 to see that we can map Y' to Y similarly, so there exists a matching between the Y 's, and
 845 the Y' 's. We define the coupling \mathcal{C} to be a probability distribution over $\mathcal{X}^* \times \mathcal{X}^*$, where
 846 the probability of (Y, Y') according to the above definition is $1/|R|!$, and it is zero for the
 847 rest of the pairs.

848 **Case 2:** R and S have one extra member: $R' = R \cup \{k\}$, and $S' = S \cup \{\ell\}$. Assume there exists a
 849 set of replacements that turns X into Y . We construct $|R| + 1$ sets of replacements that turn
 850 X' into $Y'_1, Y'_2, \dots, Y'_{|R|+1}$. We start off with X' . We choose one of the replacement in the
 851 set which replaces the sample in slot ℓ' by k' . Then, we perform all the replacements on X'
 852 except the one that is left out. Now, we do the following: We replace the sample in slot ℓ by
 853 k' and the sample in slot ℓ' by k . Clearly, we found an assignment between R' and S' , so
 854 we construct $Y'_1, \dots, Y'_{|R|}$. We also perform all the replacement in the set, and in addition
 855 to that, we replace the sample in slot ℓ by k to obtain $Y'_{|R|+1}$. It is not hard to see that given
 856 Y' , we can construct Y as well, so there is a matching between Y and the Y'_t 's. Also, Y
 857 and the Y' 's have a Hamming distance of at most three. Now, we define the coupling \mathcal{C} .

Algorithm 2 A private procedure for extending the domain

```

1: procedure PRIVATE TESTER( $X, A$ )
2:    $R, S \leftarrow \emptyset$ .
3:   for  $i = 1, 2, \dots, n$  do
4:     if  $r_i(X) \geq 1$  then
5:        $R \leftarrow R \cup \{r_i(X) \text{ copies of } i\}$ 
6:     if  $s_i(X) \geq 1$  then
7:        $S_i \leftarrow$  Set of the smallest  $s_i(X)$  indices of the entries in the bottom row of  $X$  which
       contains  $i$ .
8:        $S \leftarrow S \cup S_i$ .
9:    $S \leftarrow |R|$  smallest element in  $S$ .
10:  for each  $k \in R$  do
11:     $\ell \leftarrow$  a random element in  $S$ .
12:     $S \leftarrow S \setminus \{\ell\}$ .
13:     $X\text{-bottom}(\ell) \leftarrow k$ .
14:  Output  $\mathcal{A}(X)$ .
```

858 We set the probability of the pairs (Y, Y'_t) to be $1/(|R| + 1)!$ for $t = 1, \dots, |R| + 1$. It is
859 clear that each Y appears with probability $(|R| + 1)/(|R| + 1)! = 1/|R|!$. Thus, the desired
860 coupling exists.

861 Note that in both case, there exists a coupling \mathcal{C} such that each pair drawn from \mathcal{C} have a Hamming
862 distance of at most four. Hence the proof is complete. \square

863 C.5 Proving privacy guarantee after extending the domain

864 As we describe \mathcal{B} at a high level before, now we formally described it in Algorithm 2. Below, we
865 formally show that the algorithm is differentially private as well.

866 **Lemma C.6.** Assume \mathcal{A} is a $\xi/4$ -differentially private algorithm over \mathcal{X}^* with parameter $A \geq$
867 $12 \ln n / \delta'$ that output the correct answer with probability at least $1 - \delta$. Algorithm 2 is a ξ -
868 differentially private algorithm over \mathcal{X} . which outputs the correct answer with probability at least
869 $1 - \delta - \delta'$.

870 **Proof:** First, we claim that the algorithm changes X with probability at most δ' . Assume s is a
871 Poisson random variable with parameter λ , and let X be the set of $2s$ samples from a distribution
872 p . Using Poissonization method, we can think of $t_i(X)$ and $b_i(X)$ as two Poisson random variables
873 with mean $\lambda_i := p(i) \cdot \lambda$. Now, we bound the probability that $t_i(X)/(b_i(X) + 1)$ become larger than
874 one. If λ_i is zero, then $t_i(X)$ and $b_i(X)$ must be zero, so the ratio is below A . Let $B = b_i(X)/\lambda_i$.
875 We consider the following cases for λ_i :

876 Case 1: $\lambda_i \leq A/2$. By the concentration of a Poisson random variables, we have the following:

$$\begin{aligned} \Pr \left[\frac{t_i(X)}{b_i(X) + 1} \geq A \right] &\leq \Pr[t_i(X) - \lambda_i \geq A - \lambda_i] \leq \exp \left(-\frac{(A - \lambda_i)^2}{2A} \right) \\ &\leq \exp \left(-\frac{A}{8} \right) \end{aligned}$$

877 Case 2: $\lambda_i > A/2$. Clearly, we have:

$$\Pr \left[\frac{t_i(X)}{b_i(X) + 1} \geq A \right] \leq \Pr[t_i(X) \geq A \cdot b_i(X) + A] = \Pr[t_i(X) \geq A \cdot B \cdot \lambda_i + A]$$

878 Now, if $A \cdot B \geq 2$, we obtain:

$$\begin{aligned} \Pr \left[\frac{t_i(X)}{b_i(X) + 1} \geq A \right] &\leq \Pr[t_i(X) - \lambda_i \geq \lambda_i + A] \leq \exp \left(-\frac{(A + \lambda_i)^2}{2(A + 2\lambda_i)} \right) \\ &\leq \exp \left(-\frac{\lambda_i^2}{6\lambda_i} \right) \leq \exp \left(-\frac{A}{12} \right) \end{aligned}$$

879

If $A \cdot B = A b_i(X)/\lambda_i < 2$, it means that $b_i(X)$ is at most $2\lambda_i/A$. Thus, we have:

$$\begin{aligned} \Pr\left[\frac{t_i(X)}{b_i(X) + 1} \geq A\right] &\leq \Pr\left[b_i(X) \leq \frac{2\lambda_i}{A}\right] = \Pr\left[\lambda_i - b_i(X) \geq \frac{(A-2) \cdot \lambda_i}{A}\right] \\ &\leq \exp\left(-\frac{(A-2)^2 \lambda_i^2}{2A^2 \left(\frac{2A-2}{A} \cdot \lambda_i\right)}\right) = \exp\left(-\frac{(A-2)^2}{2A(2A-2)} \cdot \lambda_i\right) \\ &\leq \exp\left(-\frac{\lambda_i}{6}\right) \leq \exp\left(-\frac{A}{12}\right) \end{aligned}$$

880

where the second to last inequality is true when $A \geq 10$.

881 In all of the cases above, The probability that the ratio associated with element i goes above A is at
 882 most $\exp(-A/12) \leq \delta'/n$. By union bound, the probability of having at least one i with the ratio
 883 above A is at most δ' . Observe that if all the ratios are below A , all the $r_i(X)$'s will be zero. Thus,
 884 the algorithm does not change X with probability $1 - \delta'$. Also, if \mathcal{A} outputs the correct answer with
 885 probability at least $1 - \delta$, then \mathcal{B} outputs the correct answer with probability at least $1 - \delta - \delta'$.

886 Now, we show that \mathcal{B} is private. In Lemma C.5, we show our mapping has the following property:
 887 Let X and X' in \mathcal{X} be two datasets with Hamming distance at most one. Let Y and Y' be the
 888 randomized datasets that X and X' are mapped to. There exists a coupling \mathcal{C} between Y and Y'
 889 where the Hamming distance between any (Y, Y') with non-zero probability in \mathcal{C} is at most four.
 890 The existence of the coupling and the fact that \mathcal{A} is an $\xi/4$ private algorithm help us to prove the
 891 privacy guarantee for \mathcal{B} . Let O be an arbitrary output for \mathcal{B} . In the context of our paper O can be
 892 accept or reject. Below, we show the probability of outputting O on two neighboring dataset X and
 893 X' with Hamming distance one, is the same up to a e^ξ factor.

$$\begin{aligned} \Pr[\mathcal{B}(X) = O] &= \sum_Y \Pr[\mathcal{A}(Y) = O] \cdot \Pr[X \text{ is mapped to } Y] \\ &= \sum_{Y, Y'} \Pr[\mathcal{A}(Y) = O] \cdot \mathcal{C}(Y, Y') \\ &\leq \sum_{Y, Y'} e^{(\xi/c) \cdot |Y - Y'|} \Pr[\mathcal{A}(Y') = O] \cdot \mathcal{C}(Y, Y') \\ &\leq \sum_{Y, Y'} e^\xi \Pr[\mathcal{A}(Y') = O] \cdot \mathcal{C}(Y, Y') \\ &\leq \sum_{Y'} e^\xi \Pr[\mathcal{A}(Y') = O] \cdot \Pr[X' \text{ is mapped to } Y'] \\ &= e^\xi \Pr[\mathcal{B}(X') = O] \end{aligned}$$

894 Therefore, \mathcal{B} is ξ -private on \mathcal{X} . □

895 D Proof of the Lemmas

896 D.1 Proof of Lemma 3.1

897 **Lemma 3.1.** Suppose r , b_i , $s_{i,1}$, $s_{i,2}$, $v_{i,j,1}$, and $v_{i,j,2}$ are quantities defined above. Then, we have:

$$\mathbf{E}_r \left[\sum_{j=1}^{b_i} (v_{i,j,1} - v_{i,j,2})^2 - v_{i,j,1} - v_{i,j,2} \middle| b_i, s_{i,1}, s_{i,2} \right] = \frac{(s_{i,1} - s_{i,2})^2 - s_{i,1} - s_{i,2}}{b_i}.$$

898 **Proof:** Observe that given b_i , $s_{i,1}$, and $s_{i,2}$, the number of instances of element i in each bucket,
 899 $v_{i,j,1}$, is random variables drawn from a binomial distribution $\text{Bin}(s_{i,1}, 1/b_i)$. Similarly, $v_{i,j,2}$ is

900 drawn from $\text{Bin}(s_{i,2}, 1/b_i)$. Thus, we have:

$$\begin{aligned}\mathbf{E}[v_{i,j,1}] &= \frac{s_{i,1}}{b_i}, \quad \mathbf{E}[v_{i,j,1}^2] = \text{Var}[v_{i,j,1}] + \mathbf{E}[v_{i,j,1}]^2 = s_{i,1} \cdot \left(1 - \frac{1}{b_i}\right) \cdot \frac{1}{b_i} + \frac{s_{i,1}^2}{b_i^2} = \frac{s_{i,1}}{b_i} + \frac{s_{i,1}^2 - s_{i,1}}{b_i^2}, \\ \mathbf{E}[v_{i,j,2}] &= \frac{s_{i,2}}{b_i}, \quad \mathbf{E}[v_{i,j,2}^2] = \text{Var}[v_{i,j,2}] + \mathbf{E}[v_{i,j,2}]^2 = s_{i,2} \cdot \left(1 - \frac{1}{b_i}\right) \cdot \frac{1}{b_i} + \frac{s_{i,2}^2}{b_i^2} = \frac{s_{i,2}}{b_i} + \frac{s_{i,2}^2 - s_{i,2}}{b_i^2}.\end{aligned}$$

901 Since $v_{i,j,1}$ is independent from $v_{i,j,2}$, then we have:

$$\begin{aligned}\mathbf{E} \left[\sum_{j=1}^{b_i} (v_{i,j,1} - v_{i,j,2})^2 - v_{i,j,1} - v_{i,j,2} \middle| b_i, s_{i,1}, s_{i,2} \right] \\ = \sum_{j=1}^{b_i} \mathbf{E}[(v_{i,j,1} - v_{i,j,2})^2 - v_{i,j,1} - v_{i,j,2} | b_i, s_{i,1}, s_{i,2}] \\ = b_i \cdot (\mathbf{E}[u_{i,1}^2 + v_{i,1}^2 - 2 \cdot v_{i,j,1} \cdot v_{i,j,2} - v_{i,j,1} - v_{i,j,2} | b_i, s_{i,1}, s_{i,2}]) \\ = b_i \cdot \left(\frac{s_{i,1}^2 - s_{i,1}}{b_i^2} + \frac{s_{i,2}^2 - s_{i,2}}{b_i^2} - 2 \cdot \frac{s_{i,1}}{b_i} \cdot \frac{s_{i,2}}{b_i} \right) \\ = \frac{(s_{i,1} - s_{i,2})^2 - s_{i,1} - s_{i,2}}{b_i}.\end{aligned}$$

902 which completes the proof. \square

903 **Lemma A.1.** Assume F is a random set of samples to be used for flattening. Then, we have:

$$\mathbf{E}_F[d_{\max}^{(F)}] \leq \Theta \left(\mathbf{E}_F[d_{\min}^{(F)}] + \mathbf{E}_F[\|p^{(F)} - q^{(F)}\|_2^2] \right)$$

904 **Proof:** Given a random set F , we the ℓ_2 -norm of p and q are two random variables: $\|p^{(F)}\|$ and
 905 $\|p^{(F)}\|$. Recall that $d_{\max}^{(F)}$ and $d_{\min}^{(F)}$ are the minimum and the maximum of $\|p^{(F)}\|$ and $\|p^{(F)}\|$
 906 respectively. Consider an event, namely E , over the randomness of F that indicates $d_{\max}^{(F)}$ is at most
 907 $3 \cdot d_{\min}^{(F)}$. Also, let \bar{E} indicate the complimentary event, when $d_{\max}^{(F)}$ is greater than $3 \cdot d_{\min}^{(F)}$. Using
 908 Observation D.1, in this latter case, there exists a constant c such that $d_{\max}^{(F)}$ is at most $c \cdot \|p^{(F)} - q^{(F)}\|_2$.
 909 $\|q^{(F)}\|_2^2$.

910 Hence, we have:

$$\begin{aligned}\mathbf{E}_F[d_{\max}^{(F)}] &= \mathbf{E}_F[d_{\max}^{(F)} | E] \cdot \Pr_F[E] + \mathbf{E}_F[d_{\max}^{(F)} | \bar{E}] \cdot \Pr_F[\bar{E}] \\ &\leq 3 \cdot \mathbf{E}_F[d_{\min}^{(F)} | E] \cdot \Pr_F[E] + c \cdot \mathbf{E}_F[\|p^{(F)} - q^{(F)}\|_2^2 | \bar{E}] \cdot \Pr_F[\bar{E}] \\ &\leq 3 \cdot \mathbf{E}_F[d_{\min}^{(F)} | E] \cdot \Pr_F[E] + 3 \cdot \mathbf{E}_F[d_{\min}^{(F)} | \bar{E}] \cdot \Pr_F[\bar{E}] \\ &\quad + c \cdot \mathbf{E}_F[\|p^{(F)} - q^{(F)}\|_2^2 | E] \cdot \Pr_F[E] + c \cdot \mathbf{E}_F[\|p^{(F)} - q^{(F)}\|_2^2 | \bar{E}] \cdot \Pr_F[\bar{E}] \\ &= \Theta \left(\mathbf{E}_F[d_{\min}^{(F)}] + \mathbf{E}_F[\|p^{(F)} - q^{(F)}\|_2^2] \right)\end{aligned}$$

911 Therefore, the proof is complete. \square

912 **Observation D.1.** If $\|p\|_2^2 \geq C \cdot \|q\|_2^2$ for a constant $C > 1$, then $\|p - q\|_2^2 = \Theta(\|p\|_2^2)$.

913 **Proof:** By the Cauchy-Schwarz inequality, we have:

$$\begin{aligned}
& \left(\sum_i (p_i - q_i)^2 \right) \cdot \left(\sum_i (p_i + q_i)^2 \right) \geq \left(\sum_i (p_i - q_i)(p_i + q_i) \right)^2 \Rightarrow \\
& \|p - q\|_2^2 \cdot \left(2 + \frac{2}{C} \right) \cdot \|p\|_2^2 \geq \|p - q\|_2^2 \cdot \left(\sum_i 2(p_i^2 + q_i^2) \right) \geq \|p - q\|_2^2 \cdot \left(\sum_i (p_i + q_i)^2 \right) \geq \left(\sum_i p_i^2 - q_i^2 \right)^2 \Rightarrow \\
& \|p - q\|_2^2 \cdot \left(2 + \frac{2}{C} \right) \cdot \|p\|_2^2 \geq \left(\sum_i p_i^2 - q_i^2 \right)^2 \geq \left(1 - \frac{1}{C} \right)^2 \cdot \|p\|_2^4 \Rightarrow \\
& \|p - q\|_2^2 \geq \left(\frac{(1 - 1/C)^2}{2 + 2/C} \right) \cdot \|p\|_2^2 = \Omega(\|p\|_2^2).
\end{aligned}$$

914 On the other hand, we have:

$$\|p - q\|_2^2 = \sum_i (p_i - q_i)^2 \leq \sum_i p_i^2 + q_i^2 \leq \left(1 + \frac{1}{C} \right) \cdot \|p\|_2^2 = O(\|p\|_2^2).$$

915

□

916 **Lemma B.3.** Assume random variable x is drawn from $\text{Poi}(\lambda)$. If λ is at least $1.5 \cdot \ln(1/c)$, then
917 the probability of x being larger than 3λ is at most $1 - c$.

918 **Proof:** We use the tail bound for the Poisson distribution we have:

$$\Pr_x[x \geq \lambda + 2\lambda] \leq \exp\left(-\frac{(2\lambda)^2}{2 \cdot (2 + 1) \cdot \lambda}\right) \leq \exp(-2\lambda/3) \leq c.$$

919 Thus, the proof is complete.

□

920 **Lemma B.4.** Assume we have n independent random variables x_1, x_2, \dots, x_n in the range $[0, +\infty)$.
921 Suppose each x_i is at least A_i with probability $p \geq 0.95$ where A_i is a fixed number. Then, with
922 probability at least 0.9 , $\sum_{i=1}^n x_i$ is at least $0.1 \sum_{i=1}^n A_i$.

Proof: We define another set of random variables, y_i 's, as follows:

$$y_i = \begin{cases} A_i & \text{with probability } p \\ 0 & \text{with probability } 1-p \end{cases}$$

Clearly, the expected value of $\sum_{i=1}^n y_i$ is $p \cdot \sum_{i=1}^n A_i$. Note that we can see y_i as A_i multiplied by a Bernoulli random variable with bias p . Thus, the variance of $\sum_{i=1}^n y_i$ is:

$$\text{Var}\left[\sum_{i=1}^n y_i\right] = \sum_{i=1}^n \text{Var}[y_i] = \sum_{i=1}^n A_i^2 \text{Var}[\text{Ber}(p)] = p(1-p) \cdot \sum_{i=1}^n A_i^2.$$

923 Now, by the Chebyshev inequality, we can bound the probability of being far from their expectation:

$$\begin{aligned}
\Pr\left[\sum_{i=1}^n y_i \leq 0.1 \cdot \mathbf{E}\left[\sum_{i=1}^n y_i\right]\right] & \leq \Pr\left[\left|\sum_{i=1}^n y_i - \mathbf{E}\left[\sum_{i=1}^n y_i\right]\right| \geq 0.9 \cdot \mathbf{E}\left[\sum_{i=1}^n y_i\right]\right] \\
& \leq \frac{\text{Var}[\sum_{i=1}^n y_i]}{0.9^2 \cdot \mathbf{E}[\sum_{i=1}^n y_i]^2} \leq \frac{p(1-p) \sum_{i=1}^n A_i^2}{0.9^2 p^2 \cdot (\sum_{i=1}^n A_i)^2} \leq \frac{p(1-p)}{0.9^2 p^2} \leq 0.1.
\end{aligned}$$

Observe that y_i 's are defined such that the probability of $\sum_{i=1}^n x_i > a$ for any number a is at least the probability of $\sum_{i=1}^n y_i > a$. Thus, we have:

$$\Pr\left[\sum_{i=1}^n x_i \geq 0.1 \sum_{i=1}^n A_i\right] \geq \Pr\left[\sum_{i=1}^n y_i \geq 0.1 \sum_{i=1}^n A_i\right] \geq 0.9.$$

924 Hence, the proof is complete.

□

Lemma C.4. Let x_1, x_2, \dots, x_n be n non-negative random variables. Suppose there exist two constants c and p , both at most one, such that for each random variable x_i , we have:

$$\Pr[x_i < c \cdot \mathbf{E}[x_i]] \leq p,$$

Then, one can show:

$$\Pr\left[\sum_{i=1}^n x_i < \frac{c \cdot \sum_{i=1}^n \mathbf{E}[x_i]}{10}\right] \leq \frac{10p}{9}.$$

Proof: At a high level, we expect each random variable x_i to “contributes” to the sum of x_i ’s by $\mathbf{E}[x_i]$. If a random variable x_i is at least $c \mathbf{E}[x_i]$, it is contributing “enough” to the sum. While otherwise, the sum “misses” a contribution of amount $\mathbf{E}[x_i]$. The main idea is to show that total amount that the sum misses is not too large.

More Formally, for each i , we define an auxiliary random variables y_i as below. Roughly speaking y_i indicates how much the sum is missing due to a low x_i :

$$y_i = \begin{cases} \mathbf{E}[x_i] & \text{if } x_i < c \cdot \mathbf{E}[x_i] \\ 0 & \text{otherwise} \end{cases}$$

First, we claim that the sum of y_i ’s is not too large since we have:

$$\mathbf{E}\left[\sum_{i=1}^n y_i\right] = \sum_{i=1}^n \mathbf{E}[x_i] \cdot \Pr[x_i < c \cdot \mathbf{E}[x_i]] \leq p \cdot \sum_{i=1}^n \mathbf{E}[x_i].$$

Using Markov’s inequality, the sum of y_i ’s cannot be larger than $0.9 \cdot \sum_{i=1}^n \mathbf{E}[x_i]$ with probability more than $10p/9$. Hence, with probability $1 - 10p/9$, we may assume $\sum_{i=1}^n y_i$ is at most $0.9 \cdot \sum_{i=1}^n \mathbf{E}[x_i]$.

Now, we show that the sum of x_i ’s cannot be too small when the sum of y_i ’s is less than $0.9 \cdot \sum_{i=1}^n \mathbf{E}[x_i]$. To see this, let I be the set of indices i for which $x_i \geq c \cdot \mathbf{E}[x_i]$. Then, one can obtain:

$$\begin{aligned} \sum_{i=1}^n x_i &\geq \sum_{i \in I} x_i \geq c \cdot \sum_{i \in I} \mathbf{E}[x_i] = c \cdot \left(\sum_{i=1}^n \mathbf{E}[x_i] - \sum_{i \notin I} \mathbf{E}[x_i] \right) \\ &= c \cdot \left(\sum_{i=1}^n \mathbf{E}[x_i] - \sum_{i \notin I} y_i \right) = c \cdot \left(\sum_{i=1}^n \mathbf{E}[x_i] - \sum_{i=1}^n y_i \right) \\ &\geq \frac{c \cdot \sum_{i=1}^n \mathbf{E}[x_i]}{10}, \end{aligned}$$

which concludes the lemma. \square

Lemma D.2. Assume x is binomial random variable with n trials and bias p . Then, the following is true.

$$\mathbf{E}_x \left[\frac{1}{x+1} \right] \leq \min \left(\frac{1}{p \cdot (n+1)}, 1 \right)$$

Proof:

$$\begin{aligned} \mathbf{E}_x \left[\frac{1}{x+1} \right] &= \frac{1}{p \cdot (n+1)} \sum_{x=0}^n \frac{n+1}{x+1} \binom{n}{x} p^{x+1} (1-p)^{n-x} \\ &= \frac{1}{p \cdot (n+1)} \sum_{y=1}^n \binom{n+1}{y} p^y (1-p)^{(n+1)-y} = \frac{1 - (1-p)^{n+1}}{p \cdot (n+1)} \\ &\leq \min \left(\frac{1}{p \cdot (n+1)}, 1 \right) \end{aligned}$$

\square

Lemma D.3. Assume x is binomial random variable with n trials and bias p . Then, the following is true.

$$\mathbf{E}_x \left[\frac{1}{(x+2)(x+1)} \right] \leq \min \left(\frac{1}{p^2 \cdot (n+1)(n+2)}, 1 \right)$$

939 **Proof:**

$$\begin{aligned} \mathbf{E}_x \left[\frac{1}{(x+2)(x+1)} \right] &= \frac{1}{p^2 \cdot (n+1)(n+2)} \sum_{x=0}^n \frac{(n+2)(n+1)}{(x+2)(x+1)} \binom{n}{x} p^{x+2} (1-p)^{n-x} \\ &= \frac{1}{p^2 \cdot (n+1)(n+2)} \sum_{y=2}^n \binom{n+2}{y} p^y (1-p)^{(n+1)-y} \\ &= \frac{1 - (1-p)^{n+2} - (n+2)p(1-p)^{n+1}}{p^2 \cdot (n+1)(n+2)} \\ &\leq \min \left(\frac{1}{p^2 \cdot (n+1)(n+2)}, 1 \right) \end{aligned}$$

940

□

Lemma D.4. Suppose we have a bin with m balls where exactly t of them are red. We draw balls from the bin without replacement. Let X be the number of red balls in the first s trials and let Y be the number of red balls in the next k trials. Then, we have:

$$\mathbf{E} \left[\frac{X^2}{Y+1} \right] \leq \min \left(\frac{2(s-1)t}{(k+1)}, 2t^2 \right) + t.$$

941 **Proof:** We write the expectation explicitly:

$$\begin{aligned} \mathbf{E} \left[\frac{X^2}{Y+1} \right] &\leq \sum_a \sum_b \frac{a^2}{b+1} \cdot \Pr[X = a] \cdot \Pr[Y = b] = \sum_a \sum_b \frac{a^2}{b+1} \cdot \frac{\binom{s}{a} \binom{k}{b} \binom{m-s-k}{t-a-b}}{\binom{m}{t}} \\ &= \sum_{a \geq 2} \sum_b \frac{2a(a-1)}{b+1} \frac{\binom{s}{a} \binom{k}{b} \binom{m-s-k}{t-a-b}}{\binom{m}{t}} + s \cdot \sum_b \frac{1}{b+1} \frac{\binom{k}{b} \binom{m-s-k}{t-1-b}}{\binom{m}{t}} \\ &= \sum_{a \geq 2} \sum_b \frac{2a(a-1)}{b+1} \frac{\binom{s}{a} \binom{k}{b} \binom{m-s-k}{t-a-b}}{\binom{m}{t}} + s \cdot \sum_b \frac{1}{b+1} \frac{\binom{k}{b} \binom{m-s-k}{t-1-b}}{\binom{m}{t}} \\ &= \frac{2s(s-1)t}{(k+1)m} \sum_{a \geq 2} \sum_b \frac{\binom{s-2}{a-2} \binom{k+1}{b+1} \binom{m-s-k}{t-a-b}}{\binom{m-1}{t-1}} + \frac{s}{k+1} \cdot \sum_b \frac{\binom{k+1}{b+1} \binom{m-s-k}{t-1-b}}{\binom{m}{t}} \end{aligned}$$

We define the two sums in the last line as A and B :

$$A := \sum_{a \geq 2} \sum_b \frac{\binom{s-2}{a-2} \binom{k+1}{b+1} \binom{m-s-k}{t-a-b}}{\binom{m-1}{t-1}}, \quad B := \sum_b \frac{\binom{k+1}{b+1} \binom{m-s-k}{t-1-b}}{\binom{m}{t}}.$$

942 We claim A and B are two probabilities of the following randomized processes, so we can bound
943 them. Suppose we have an urn with $m-1$ balls, $t-1$ of them are red. A is the probability that
944 we get at least one red ball if we draw $k+1$ balls from the bin without replacement. Let Z be the
945 number of red balls we draw after $k+1$ draws. Using Markov's inequality, we get:

$$A = \Pr[Z \geq 1] \leq \min(1, \mathbf{E}[Z]) \leq \min \left(1, \frac{(t-1) \cdot (k+1)}{(m-1)} \right)$$

Furthermore, we can define B as the following probability: Assume we have an urn of m balls including t red balls. If we draw $(s-1) + (k+1)$ balls from the urn without replacement. B is the probability that non of the $s-1$ draws are red, and there is at least one red draw in the next $k+1$ draws. This is clearly smaller than the probability of seeing at least one red ball in the $k+1$ draws. Thus, similar to the above, we have:

$$B \leq \min \left(1, \frac{t \cdot (k+1)}{m} \right).$$

Now, putting all these together, and using the fact that $s \leq m$, we obtain:

$$\mathbf{E} \left[\frac{X^2}{Y+1} \right] \leq \min \left(\frac{2(s-1)t}{(k+1)}, 2t^2 \right) + t.$$

946

□

Lemma D.5. Assume X is a random variable drawn from $\mathbf{HG}(m, t, k)$, then

$$\mathbf{E}_X \left[\frac{1}{(X+1)} \right] \leq \min \left(1, \frac{(m+1)}{(t+1)(k+1)} \right).$$

947 **Proof:** Clearly, the expectation cannot be larger than one since $X \geq 0$. For the other term, by the
948 definition, we can achieve the following bound:

$$\begin{aligned} \mathbf{E}_x \left[\frac{1}{x+1} \right] &= \sum_{x=\max(0, k-(m-t))}^{\min(t, k)} \mathbf{HG}(x; m, t, k) \cdot \frac{1}{x+1} = \sum_x \frac{\binom{t}{x} \binom{m-t}{k-x}}{\binom{m}{k}} \cdot \frac{1}{x+1} \\ &= \sum_x \frac{m+1}{(t+1)(k+1)} \frac{\frac{t+1}{x+1} \binom{t}{x} \binom{m-t}{k-x}}{\frac{m+1}{k+1} \binom{m}{k}} = \sum_x \frac{m+1}{(t+1)(k+1)} \cdot \frac{\binom{t+1}{x+1} \binom{(m+1)-(t+1)}{(k+1)-(x+1)}}{\binom{m+1}{k+1}} \\ &\leq \frac{m+1}{(t+1)(k+1)} \cdot \sum_x \mathbf{HG}(x+1; m+1, t+1, k+1) \leq \frac{m+1}{(t+1)(k+1)} \end{aligned}$$

949 where the last line is true because the sum of the probabilities according to a distribution is at most
950 one. □

Lemma D.6. Assume X is a random variable drawn from $\mathbf{HG}(m, t, k)$, then

$$\mathbf{E}_X \left[\frac{1}{(X+2)(X+1)} \right] \leq \min \left(1, \frac{(m+2)(m+1)}{(t+2)(t+1)(k+2)(k+1)} \right).$$

951 **Proof:** Clearly, the expectation cannot be larger than one since $X \geq 0$. For the other term, by the
952 definition, we can achieve the following bound:

$$\begin{aligned} \mathbf{E}_X \left[\frac{1}{(X+2)(X+1)} \right] &= \sum_X \mathbf{HG}(X; m, t, k) \cdot \frac{1}{(X+2)(X+1)} \\ &= \sum_X \frac{\binom{t}{X} \binom{m-t}{k-X}}{\binom{m}{k}} \cdot \frac{1}{(X+2)(X+1)} \\ &= \frac{(m+2)(m+1)}{(t+2)(t+1)(k+2)(k+1)} \cdot \sum_X \frac{\binom{t+2}{x+2} \binom{(m+2)-(t+2)}{(k+2)-(x+2)}}{\binom{m+2}{k+2}} \\ &= \frac{(m+2)(m+1)}{(t+2)(t+1)(k+2)(k+1)} \sum_x \mathbf{HG}(x; m+2, t+2, k+2) \\ &\leq \frac{(m+2)(m+1)}{(t+2)(t+1)(k+2)(k+1)} \end{aligned}$$

953 where the last line is true, because the sum of the probabilities in a distribution is at most one. □