

Chapter 39

Double sampling

By Sarel Har-Peled, April 26, 2022^①

“What does not work when you apply force, would work when you apply even more force.”

, Anonymous

39.1. Double sampling

Double sampling is the idea that two random independent samples should look similar, and should *not* be completely different in the way they intersect a certain set. We use the following sampling model, which makes the computations somewhat easier.

Definition 39.1.1. Let $S = \{f_1, \dots, f_n\}$ be a set of objects, where the i th object has weight $\omega_i > 0$, for all i . Let $W = \sum_i \omega_i$. For a target size ρ , a **ρ -sample** is a random sample $R \subseteq S$, where object f_i is picked independently with probability $\rho\omega_i/W$. To simplify the discussion, we assume that $\rho\omega_i/W < 1$. Handling the more general case is easy if somewhat tedious.

Lemma 39.1.2. Let R_1 and R_2 be two ρ -samples, and consider the merged sample $R = R_1 \cup R_2$. Let $T \subseteq S$ be a set of m elements. Then, we have that

$$\mathbb{P}[T \subseteq R_1 \mid T \subseteq R] \geq \frac{1}{2^m} \quad \text{and} \quad \mathbb{P}[T \subseteq R_1 \text{ and } T \cap R_2 = \emptyset \mid T \subseteq R] \leq \frac{1}{2^m}.$$

Proof: Consider an object $f \in T$, and observe that $\mathbb{P}[f \in R_1 \text{ or } f \in R_2 \mid f \in R] = 1$. As such, by symmetry

$$\mathbb{P}[f \in R_1 \mid f \in R] = \mathbb{P}[f \in R_2 \mid f \in R] \geq 1/2,$$

Now, let $T = \{f_1, \dots, f_m\}$. Since R_1 and R_2 are independent, and each element is being picked independently, we have that

$$\begin{aligned} \mathbb{P}[T \subseteq R_1 \mid T \subseteq R] &= \mathbb{P}[f_1, \dots, f_m \in R_1 \mid f_1, \dots, f_m \in R] = \prod_{i=1}^m \mathbb{P}[f_i \in R_1 \mid f_1, \dots, f_m \in R] \\ &= \prod_{i=1}^m \mathbb{P}[f_i \in R_1 \mid f_i \in R] \geq \frac{1}{2^m}. \end{aligned}$$

For the second claim, observe that again, by symmetry, we have that

$$\mathbb{P}[f \in R_1 \text{ and } f \notin R_2 \mid f \in R] = \mathbb{P}[f \notin R_1 \text{ and } f \in R_2 \mid f \in R] \leq 1/2,$$

as the two events are disjoint. Now, the claim follows by arguing as above. ■

^①This work is licensed under the Creative Commons Attribution-Noncommercial 3.0 License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc/3.0/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.

39.1.1. Disagreement between samples on a specific set

We provide three proofs of the following lemma – the constants are somewhat different for each version.

Lemma 39.1.3. *Let R_1 and R_2 be two ρ -samples from a ground set S , and consider a fixed set $T \subseteq S$. We have that*

$$\mathbb{P}\left[\left| |R_1 \cap T| - |R_2 \cap T| \right| > \varepsilon \rho \right] \leq 3 \exp(-\varepsilon^2 \rho / 2).$$

Proof: (Simplest proof.) By Chernoff's inequality, for $\delta \in (0, 1)$, we have

$$\mathbb{P}\left[\left| |R_1| - \rho \right| \geq (\varepsilon/2)\rho \right] \leq 2 \exp(-(\varepsilon/2)^2 \rho / 4) = 2 \exp(-\varepsilon^2 \rho / 16).$$

The same holds for R_2 , and as such we have

$$\begin{aligned} \mathbb{P}\left[\left| |R_1| - |R_2| \right| \geq \varepsilon \rho \right] &\leq \mathbb{P}\left[\left| |R_1| - \rho \right| + \left| \rho - |R_2| \right| \geq \varepsilon \rho \right] \\ &\leq \mathbb{P}\left[\left| |R_1| - \rho \right| \geq (\varepsilon/2)\rho \right] + \mathbb{P}\left[\left| \rho - |R_2| \right| \geq (\varepsilon/2)\rho \right] \leq 4 \exp(-\varepsilon^2 \rho / 16) \quad \blacksquare \end{aligned}$$

Proof: For an object $f_i \in S$, let X_i be a random variable, where

$$X_i = \begin{cases} 1 & f_i \in R_1 \text{ and } f_i \notin R_2 \\ -1 & f_i \notin R_1 \text{ and } f_i \in R_2 \\ 0 & \text{otherwise.} \end{cases}$$

We have that $p_i = \mathbb{P}[X_i = 1] = \mathbb{P}[X_i = -1] = (\rho \omega_i / W)(1 - \rho \omega_i / W)$ and $\mathbb{E}[X_i] = 0$. Applying the regular concentration inequalities in this case is not immediate, since there are many X_i s that are zero. To overcome this, let T be a random variable that is the number of variables in X_1, \dots, X_n that are non-zero. We have that T is a sum of n independent 0/1 random variables, where $\mathbb{E}[T] = \sum_i 2p_i = 2\rho$. In particular, by [Chernoff's inequality](#), we have that

$$q_1 = \mathbb{P}[T > (1 + \varepsilon)2\rho] \leq \exp(-2\rho\varepsilon^2/4) = \exp(-\rho\varepsilon^2/2).$$

and assume this happens. In particular, let Z_1, \dots, Z_T be the non-zero variables in X_1, \dots, X_n , and observe that $\mathbb{P}[Z_i = 1] = \mathbb{P}[Z_i = -1] = 1/2$. Let $Y = \sum_i X_i = \sum_i Z_i$. Observe that $\mathbb{E}[Y] = 0$, and by [Chernoff inequality](#), we have that

$$\begin{aligned} q_2 &= \mathbb{P}\left[\left| |R_1 \cap S| - |R_2 \cap S| \right| > \varepsilon \rho \right] = \mathbb{P}\left[|Y - \mathbb{E}[Y]| \geq \varepsilon \rho \right] \leq \mathbb{P}\left[\left| \sum_i Z_i - 0 \right| \geq \varepsilon \rho \right] \\ &\leq 2 \exp\left(-2 \frac{(\varepsilon \rho)^2}{2T}\right) \leq 2 \exp\left(-2 \frac{(\varepsilon \rho)^2}{2(1 + \varepsilon)\rho}\right) + q_1 = 2 \exp\left(-\frac{\varepsilon^2 \rho}{1 + \varepsilon}\right) + q_1 \leq 3 \exp(-\varepsilon^2 \rho / 2), \end{aligned}$$

using $T \leq (1 + \varepsilon)2\rho$. ■

39.1.2. Exponential decay for a single set

Lemma 39.1.4. *Consider a set S of m objects, where every object $f_i \in S$ has weight $\omega_i > 0$, and $W = \sum_{i=1}^m \omega_i$. Next, consider a set $\mathbf{r} \subseteq S$ such that $\omega(\mathbf{r}) \geq tW/\rho$ (such a set is [t-heavy](#)). Let R be a ρ -sample from S . Then, the probability that R misses \mathbf{r} is at most e^{-t} . Formally, we have $\mathbb{P}[\mathbf{r} \cap R = \emptyset] \leq \exp(-t)$.*

Proof: Let $\mathbf{r} = \{f_1, \dots, f_k\}$. Clearly, the probability that R fails to pick one of these conflicting objects, is bounded by $\mathbb{P}[\mathbf{r} \cap R = \emptyset] = \mathbb{P}[\forall i \in \{1, \dots, k\} \quad f_i \notin R] = \prod_{i=1}^k (1 - \rho \frac{\omega_i}{W}) \leq \prod_{i=1}^k \exp(-\rho \frac{\omega_i}{W}) = \exp(-\frac{\rho}{W} \sum_i \omega_i) \leq \exp(-\frac{\rho}{W} \cdot t \frac{W}{\rho}) = \exp(-t)$. ■

39.1.3. Moments of the sample size

Lemma 39.1.5. *Let R an m -sample. And let $f(t) \leq \alpha t^\beta$, where $\alpha \geq 1$ and $\beta \geq 1$ are constants, such that $m \geq 16\beta$. Then $U(m) = \mathbb{E}[f(|R|)] \leq 2\alpha(2m)^\beta$.*

Proof: The proof follows from Chernoff's inequality and some tedious but straightforward calculations. The reader is as such encouraged to skip reading it.

Let $X = |R|$. This is a sum of 0/1 random variables with expectation m . As such, we have

$$\nu = \mathbb{E}[f(|R|)] \leq \sum_{i=0}^{\infty} \mathbb{P}[X = i] f(i) \leq \alpha \sum_{i=0}^{\infty} \mathbb{P}[X = i] i^\beta.$$

Considering the last sum, we have

$$\sum_{i=0}^{\infty} \mathbb{P}[X = i] i^\beta \leq \sum_{j=0}^{\infty} \mathbb{P}[X \geq jm] ((j+1)m)^\beta \leq (2m)^\beta + m^\beta \sum_{j=2}^{\infty} \mathbb{P}[X \geq jm] (j+1)^\beta.$$

We bound the last summation using Chernoff's inequality (see [Theorem 39.3.2](#)), we have

$$\begin{aligned} \tau &= \sum_{j=2}^5 \mathbb{P}[X \geq jm] (j+1)^\beta + \sum_{j=6}^{\infty} \mathbb{P}[X \geq jm] (j+1)^\beta \\ &\leq \sum_{j=2}^5 \exp\left(-\frac{m(j-1)^2}{4}\right) (j+1)^\beta + \sum_{j=6}^{\infty} 2^{-jm} (j+1)^\beta \\ &\leq \exp\left(-\frac{m}{4}\right) 3^\beta + \exp(-m) 4^\beta + \exp(-2m) 5^\beta + \exp(-4m) 6^\beta + \sum_{j=6}^{\infty} 2^{-jm} (j+1)^\beta < 1, \end{aligned}$$

since $m \geq 16\beta$. We conclude that $\nu \leq \alpha(2m)^\beta + \alpha m^\beta \tau \leq 2\alpha(2m)^\beta$. ■

Remark 39.1.6. The constant 16 in the above lemma is somewhat strange. A better constant can be derived by breaking the range of sizes into smaller intervals and using the right Chernoff inequality. Since this is somewhat tangential to the point of this write-up, we leave it as is (i.e., this constant is not critical to our discussion).

39.1.4. Growth function

The **growth function** $\mathcal{G}_\delta(n)$ is the maximum number of ranges in a range space with VC dimension δ , and with n elements. By Sauer's lemma, it is known that

$$\mathcal{G}_\delta(n) = \sum_{i=0}^{\delta} \binom{n}{i} \leq \sum_{i=0}^{\delta} \frac{n^i}{i!} \leq n^\delta, \quad (39.1)$$

The following is well known (the estimates are somewhat tedious to prove):

Lemma 39.1.7 ([\[Har11\]](#)). *For $n \geq 2\delta$ and $\delta \geq 1$, we have $\left(\frac{n}{\delta}\right)^\delta \leq \mathcal{G}_\delta(n) \leq 2\left(\frac{ne}{\delta}\right)^\delta$, where $\mathcal{G}_\delta(n) = \sum_{i=0}^{\delta} \binom{n}{i}$.*

Lemma 39.1.8. *Let R and R' be two independent m -samples from x . Assume that $m \geq \delta$. Then $\mathbb{E}[\mathcal{G}_\delta(|R| + |R'|)] \leq G_\delta(2m)$, where $G_\delta(2m) = 4(4em/\delta)^\delta$.*

Proof: We set $\alpha = 2(\frac{\varepsilon}{\delta})^\delta$, $\beta = \delta$, and $f(n) = \alpha n^\beta$. Duplicate every element in x , and let x' be the resulting set. Clearly, the size of a $2m$ -sample R from x' is the same as $|R| + |T|$. By Lemma 39.1.7, we have $\mathbb{E}[\mathcal{G}_\delta(|R|)] \leq \mathbb{E}[f(|R|)] \leq 2\alpha(4m)^\beta \leq 4(\frac{4em}{\delta})^\delta$. The last inequality follows from Lemma 39.1.5. ■

39.2. Proof of the ε -net theorem

Here we are working in the unweighted settings (i.e., the weight of a single element is one).

Theorem 39.2.1 (ε -net theorem, [HW87]). *Let (X, \mathcal{R}) be a range space of VC dimension δ , let x be a finite subset of X , and suppose that $0 < \varepsilon \leq 1$ and $\varphi < 1$. Let N be an m -sample from x (see Definition 39.1.1), where*

$$m \geq \max\left(\frac{8}{\varepsilon} \lg \frac{4}{\varphi}, \frac{16\delta}{\varepsilon} \lg \frac{16}{\varepsilon}\right). \quad (39.2)$$

Then N is an ε -net for x with probability at least $1 - \varphi$.

39.2.1. The proof

39.2.1.1. Reduction to double sampling

Let $n = |x|$. Let N be the m -sample from x . Let \mathcal{E}_1 be the probability that N fails to be an ε -net. Namely,

$$\mathcal{E}_1 = \{\exists \mathbf{r} \in \mathcal{R} \mid |\mathbf{r} \cap x| \geq \varepsilon n \text{ and } \mathbf{r} \cap N = \emptyset\}.$$

(Namely, there exists a “heavy” range \mathbf{r} that does not contain any point of N .) To complete the proof, we must show that $\mathbb{P}[\mathcal{E}_1] \leq \varphi$. Let T be another m -sample generated in a similar fashion to N . Let \mathcal{E}_2 be the event that N fails but T “works”. Formally

$$\mathcal{E}_2 = \left\{ \exists \mathbf{r} \in \mathcal{R} \mid |\mathbf{r} \cap x| \geq \varepsilon n, \mathbf{r} \cap N = \emptyset, \text{ and } |\mathbf{r} \cap T| \geq \frac{\varepsilon m}{2} \right\}.$$

Intuitively, since $\mathbb{E}[|\mathbf{r} \cap x|] \geq \varepsilon m$, we have that for the range \mathbf{r} that N fails for, it follows with “good” probability that $|\mathbf{r} \cap T| \geq \varepsilon m/2$. Namely, \mathcal{E}_1 and \mathcal{E}_2 have more or less the same probability.

Claim 39.2.2. $\mathbb{P}[\mathcal{E}_2] \leq \mathbb{P}[\mathcal{E}_1] \leq 2\mathbb{P}[\mathcal{E}_2]$.

Proof: Clearly, $\mathcal{E}_2 \subseteq \mathcal{E}_1$, and thus $\mathbb{P}[\mathcal{E}_2] \leq \mathbb{P}[\mathcal{E}_1]$. As for the other part, note that by the definition of conditional probability, we have

$$\mathbb{P}[\mathcal{E}_2 \mid \mathcal{E}_1] = \mathbb{P}[\mathcal{E}_2 \cap \mathcal{E}_1] / \mathbb{P}[\mathcal{E}_1] = \mathbb{P}[\mathcal{E}_2] / \mathbb{P}[\mathcal{E}_1].$$

It is thus enough to show that $\mathbb{P}[\mathcal{E}_2 \mid \mathcal{E}_1] \geq 1/2$.

Assume that \mathcal{E}_1 occurs. There is $\mathbf{r} \in \mathcal{R}$, such that $|\mathbf{r} \cap x| > \varepsilon n$ and $\mathbf{r} \cap N = \emptyset$. The required probability is at least the probability that for this specific \mathbf{r} , we have $X = |\mathbf{r} \cap T| \geq \frac{\varepsilon n}{2}$. The variable X is a sum of $t = |\mathbf{r} \cap x| \geq \varepsilon n$ random independent 0/1 variables, each one has probability m/n to be one. Setting $\mu = \mathbb{E}[X] = tm/n \geq \varepsilon m$ and $\xi = 1/2$, we have by Chernoff’s inequality that

$$\mathbb{P}[|\mathbf{r} \cap T| \leq \varepsilon m/2] \leq \mathbb{P}[X < (1 - \xi)\mu] \leq \exp(-\mu\xi^2/2) = \exp(-\varepsilon m/8) < 1/2,$$

if $\varepsilon m \geq 8$. Thus, for $\mathbf{r} \in \mathcal{E}_1$, we have $\mathbb{P}[\mathcal{E}_2] / \mathbb{P}[\mathcal{E}_1] \geq \mathbb{P}[|\mathbf{r} \cap T| \geq \frac{\varepsilon m}{2}] = 1 - \mathbb{P}[|\mathbf{r} \cap T| < \frac{\varepsilon m}{2}] \geq \frac{1}{2}$. ■

Claim 39.2.2 implies that to bound the probability of \mathcal{E}_1 , it is enough to bound the probability of \mathcal{E}_2 . Let

$$\mathcal{E}'_2 = \left\{ \exists \mathbf{r} \in \mathcal{R} \mid \mathbf{r} \cap N = \emptyset \text{ and } |\mathbf{r} \cap T| \geq \frac{\varepsilon m}{2} \right\}.$$

Clearly, $\mathcal{E}_2 \subseteq \mathcal{E}'_2$. Thus, bounding the probability of \mathcal{E}'_2 is enough to prove **Theorem 39.2.1**. Note, however, that a shocking thing happened! We no longer have \mathbf{x} participating in our event. Namely, we turned bounding an event that depends on a global quantity (i.e., the ground set \mathbf{x}) into bounding a quantity that depends only on a local quantity/experiment (involving only N and T). This is the crucial idea in this proof.

39.2.1.2. Using double sampling to finish the proof

Claim 39.2.3. $\mathbb{P}[\mathcal{E}_2] \leq \mathbb{P}[\mathcal{E}'_2] \leq 2^{-\varepsilon m/2} G_\delta(2m)$.

Proof: We fix the content of $\mathbf{R} = N \cup T$. The range space $(\mathbf{R}, \mathcal{R}_{|\mathbf{R}})$ has $\mathcal{G}_\delta(|\mathbf{R}|)$ ranges. Fix a range \mathbf{r} in this range space. Let $T = \mathbf{r} \cap \mathbf{R}$. If $b = |T| < \varepsilon m/2$ then the \mathcal{E}'_2 can not happened. Otherwise, the probability that \mathbf{r} is a bad range is $\mathbb{P}\left[T \subseteq T \text{ and } T \cap N = \emptyset \mid T \subseteq \mathbf{R}\right] \leq \frac{1}{2^b}$, by **Lemma 39.1.2**. In particular, by the union bound over all ranges, we have $\mathbb{P}[\mathcal{E}'_2 \mid \mathbf{R}] \leq 2^{-\varepsilon m/2} \mathcal{G}_\delta(|\mathbf{R}|)$. As such, we have

$$\mathbb{P}[\mathcal{E}'_2] = \sum_{\mathbf{R}} \mathbb{P}[\mathcal{E}'_2 \mid \mathbf{R}] \mathbb{P}[\mathbf{R}] \leq \sum_{\mathbf{R}} 2^{-\varepsilon m/2} \mathcal{G}_\delta(|\mathbf{R}|) \mathbb{P}[\mathbf{R}] \leq 2^{-\varepsilon m/2} \mathbb{E}\left[\mathcal{G}_\delta(|\mathbf{R}|)\right] \leq 2^{-\varepsilon m/2} G_\delta(2m).$$

by **Lemma 39.1.8**. ■

Proof of THEOREM 39.2.1. By **Claim 39.2.2** and **Claim 39.2.3**, we have that $\mathbb{P}[\mathcal{E}_1] \leq 2 \cdot 2^{-\varepsilon m/2} G_\delta(2m)$. It thus remains to verify that if m satisfies **Eq. (39.2)**, then the above is smaller than φ . Which is equivalent to

$$\begin{aligned} 2 \cdot 2^{-\varepsilon m/2} G_\delta(2m) \leq \varphi &\iff 16 \cdot 2^{-\varepsilon m/2} \left(\frac{4em}{\delta}\right)^\delta \leq \varphi \iff -4 + \frac{\varepsilon m}{2} - \delta \lg\left(\frac{4em}{\delta}\right) \geq \lg \frac{1}{\varphi} \\ &\iff \left(\frac{\varepsilon m}{8} - 4 - \delta \lg \frac{4e}{\delta}\right) + \left(\frac{\varepsilon m}{8} - \lg \frac{1}{\varphi}\right) + \left(\frac{\varepsilon m}{4} - \delta \lg\left(\frac{m}{\delta}\right)\right) \geq 0 \end{aligned}$$

We remind the reader that the value of m we pick is such that $m \geq \max\left(\frac{8}{\varepsilon} \lg \frac{4}{\varphi}, \frac{16\delta}{\varepsilon} \lg \frac{16}{\varepsilon}\right)$. In particular, $m \geq 64\delta/\varepsilon$ and $-4 - \delta \lg\left(\frac{4e}{\delta}\right) \geq -4 - 4\delta \leq -8\delta \geq -\varepsilon m/8$. Similarly, by the choice of m , we have $\varepsilon m/8 \geq \lg \frac{1}{\varphi}$. As such, we need to show that $\frac{\varepsilon m}{4} \geq \delta \lg\left(\frac{m}{\delta}\right) \iff m \geq \frac{4\delta}{\varepsilon} \lg \frac{m}{\delta}$, and one can verify using some easy but tedious calculations that this holds if $m \geq \frac{16\delta}{\varepsilon} \lg \frac{16}{\varepsilon}$. ■

39.3. From previous lectures

Lemma 39.3.1. *Let X_1, \dots, X_n be n independent Bernoulli trials, where $\mathbb{P}[X_i = 1] = p_i$, and $\mathbb{P}[X_i = 0] = 1 - p_i$, for $i = 1, \dots, n$. Let $X = \sum_{i=1}^n X_i$, and $\mu = \mathbb{E}[X] = \sum_i p_i$. For $\delta \in (0, 4)$, we have*

$$\mathbb{P}[X > (1 + \delta)\mu] < \exp(-\mu\delta^2/4),$$

Theorem 39.3.2. Let X_1, \dots, X_n be n independent variables, where $\mathbb{P}[X_i = 1] = p_i$ and $\mathbb{P}[X_i = 0] = q_i = 1 - p_i$, for all i . Let $X = \sum_{i=1}^n X_i$. $\mu = \mathbb{E}[X] = \sum_i p_i$. For any $\delta > 0$, we have

$$\mathbb{P}[X > (1 + \delta)\mu] < \left(e^\delta / (1 + \delta)^{1+\delta} \right)^\mu.$$

Theorem 39.3.3. Let X_1, \dots, X_n be n independent random variables, such that $\mathbb{P}[X_i = 1] = \mathbb{P}[X_i = -1] = \frac{1}{2}$, for $i = 1, \dots, n$. Let $Y = \sum_{i=1}^n X_i$. Then, for any $\Delta > 0$, we have

$$\mathbb{P}[Y \geq \Delta] \leq \exp(-\Delta^2/2n).$$

References

- [Har11] S. Har-Peled. *Geometric approximation algorithms*. Vol. 173. Math. Surveys & Monographs. Boston, MA, USA: Amer. Math. Soc., 2011.
- [HW87] D. Haussler and E. Welzl. *ϵ -nets and simplex range queries*. *Discrete Comput. Geom.*, 2: 127–151, 1987.