

Confident Monte Carlo: Rigorous Analysis of Guessing Curves for Probabilistic Password Models

anonymous submission

Abstract—In password security a defender would like to identify and warn users with weak passwords. Similarly, the defender may also want to predict what fraction of passwords would be cracked within B guesses as the attacker’s guessing budget B varies from small (online attacker) to large (offline attacker). Towards each of these goals the defender would like to quickly estimate the guessing number for each user password pwd assuming that the attacker uses a password cracking model M i.e., how many password guesses will the attacker check before s/he cracks each user password pwd . Dell’Amico and Filippone [8] developed an efficient Monte Carlo algorithm to estimate the guessing number of a given password pwd — naïve brute-force enumeration can be prohibitively expensive when the guessing number is very large. While Dell’Amico and Filippone proved that their estimator is unbiased there is no guarantee that the Monte Carlo estimates are accurate nor does the method provide confidence ranges on the estimated guessing number or even indicate if/when there is a higher degree of uncertainty.

Our contributions are as follows: First, we identify theoretical examples where, with high probability, Monte Carlo Strength estimation produces highly inaccurate estimates of individual guessing numbers as well as the entire guessing curve. Second, we introduce Confident Monte Carlo Strength Estimation as an extension of Dell’Amico and Filippone [8]. Given a password our estimator generates an upper and lower bound with the guarantee that, except with probability δ , the true guessing number lies within the given confidence range. Our techniques can also be used to characterize the attacker’s guessing curve. In particular, given a probabilistic password cracking model M we can generate high confidence upper and lower bounds on the fraction of passwords that the attacker will crack as the guessing budget B varies.

Index Terms—Monte Carlo Estimation, Password Cracking Models, Concentration bounds

1. Introduction

In addition to their offensive uses, Probabilistic Password Cracking Models have many defensive applications. One defensive application is to use the password cracking model to estimate the strength of a user’s password during account registration so that we can warn users who attempt to register with a weak password that would be easy for an attacker to guess. For this application we would like

to quickly determine how many guesses an attacker using password cracking model M would need to check before s/he cracks a particular user’s password pwd . One way to determine the guessing number of a particular password pwd would simply be to enumerate all possible password guesses, ordered according to their probability under the model M , keeping track of the number of incorrect guesses which appear before the correct password pwd appears. However, naïve brute-force enumeration is often prohibitively expensive especially when the guessing number is very large e.g., $> 10^{15}$.

Dell’Amico and Filippone [8] developed a Monte Carlo algorithm to efficiently output an estimate of the guessing number of a given password pwd without resorting to brute-force enumeration. Their technique applies generically to any probabilistic password model M under the assumption that (1) the model M defines a distribution over passwords and we can efficiently sample from this distribution, and (2) given a particular password pwd we can quickly compute the probability that this password is generated by M . Dell’Amico and Filippone proved that their Monte Carlo estimator is unbiased i.e., the expected value of the estimate is equal to the actual guessing number. However, there is no absolute guarantee that the estimate is accurate. Their Monte Carlo estimator does not provide any statistical confidence intervals for the range of possible guessing numbers or even indicate if/when there is a high degree of uncertainty about the true guessing number.

In other defensive applications we may want to estimate the attacker’s entire guessing curve. How many consecutive incorrect login attempts should be allowed before we lock-down an account? Will doubling the cost of the password hash function significantly reduce the fraction of passwords that an offline attacker will crack?

Formally, let $\lambda_{M,B,D}$ denote the fraction of passwords in a dataset D which would be cracked within the first B guesses generated by model M . Similarly, let $\lambda_{M,B}$ denote the probability that a fresh password sampled from an (unknown) distribution \mathcal{P} over user passwords would be cracked within the first B guesses output by the model M . Characterizing the entire guessing curve $\lambda_{M,B,D}$ (or $\lambda_{M,B}$) as B ranges from small (online attack) to very large (offline attack) can help a defender set password policies. Thus, we would like to generate high confidence upper/lower bounds on the guessing curves $\lambda_{M,B,D}$ and/or $\lambda_{M,B}$.

However, when B is sufficiently large we cannot efficiently compute $\lambda_{M,B,D}$ since that would require enumer-

ation of the top B passwords in the distribution M . As a heuristic approach we could use Monte Carlo strength estimation [8] to quickly estimate the guessing number of each password in D and then compute the fraction of passwords whose estimated guessing number is below B . While this heuristic approach has become popular in the password literature (e.g., [2], [11], [14], [16]), it could yield poor estimates of $\lambda_{M,B,D}$ if/when the estimated guessing numbers are inaccurate. Indeed, in our empirical analysis we identify several instances where the heuristically estimated guessing curve is highly inaccurate. The bottom line is that there is no absolute guarantee that our estimate $\lambda_{M,B,D}$ is accurate nor does this heuristic approach provide any statistical confidence intervals on the attacker’s guessing curve.

In many settings we would like to characterize the guessing curve $\lambda_{M,B}$ using a relatively small number of samples from our unknown password distribution. For example, suppose that we conduct a user study to determine whether or not a particular password policy intervention, e.g., requiring users to select passwords with upper and lower case letters¹, effectively strengthens the distribution over user passwords. We can view the data collected from the user study as samples S_1 (control group) and S_2 (intervention group) from two different (unknown) password distributions \mathcal{P}_1 and \mathcal{P}_2 respectively. We would like to use these samples to draw statistical comparisons about the guessing curves before/after the policy intervention i.e., does the policy reduce the fraction of user passwords cracked by an offline attacker making $B = 10^{15}$ guesses per user. Observe that if we assume that the dataset S_i was sampled iid from distribution \mathcal{P}_i that we have $\lambda_{M,B} = \mathbb{E}[\lambda_{M,B,S_i}]$. We could use the popular heuristic estimate for λ_{M,B,S_i} (described above) to estimate the guessing curves $\lambda_{M,B}$ for the two different distributions. However, we now have an additional source of estimation error due to sampling of the datasets S_1 and S_2 — in addition to the error from Monte Carlo strength estimation. One challenge is that the number of samples collected from the user study will typically be constrained by research budgets making it harder to ensure that the sampling error is small.

1.1. Contributions

First, we provide theoretical examples of models M and passwords pwd where regular Monte Carlo will (1) dramatically overestimate the true guessing number of pwd with high probability, and (2) underestimate the guessing number by a factor of ≈ 2 with probability at least 0.3. We argue that such issues are inherent to *any* blackbox method for Monte Carlo strength estimation. We also consider the

1. There are many password interventions that one could consider. We could require passwords to include special characters or numbers (password composition policies). We could display a password strength meter during registration. We could warn users about the risks of picking weak passwords before registration. We could require users to participate in training activities or watch an instructional video about picking good passwords before registration.

popular heuristic of using regular Monte Carlo strength estimation to estimate the attacker’s guessing curve. We provide a dramatic (theoretical) example of a model M and a password distribution \mathcal{P} such that (1) an attacker following model M will actually crack 0% of passwords within B guesses i.e., $\lambda_{M,B} = 0$ and $\lambda_{M,B,D} = 0$, but (2) the popular heuristic using regular Monte Carlo Strength estimation incorrectly predicts that an attacker will crack 100% of user passwords sampled from \mathcal{P} within B guesses!

Second, we introduce several rigorous statistical techniques to upper/lower bound the guessing number of a password and show how these techniques can be extended to upper/lower bound an attacker’s entire guessing curve. On a technical note our upper/lower bounds are derived using concentration inequalities such as Hoeffding and Chernoff as well as a strategic application of Markov’s inequality. We call our new toolkit Confident Monte Carlo. In particular, given a password cracking model M , a particular password pwd and a confidence parameter δ Confident Monte Carlo will output an upper bound U and a lower bound L on the (unknown) true guessing number $G(pwd)$ with the guarantee that $\Pr[L \leq G(pwd)] \geq 1 - \delta$ and $\Pr[U \geq G(pwd)] \geq 1 - \delta$. Confident Monte Carlo works under the exact same generic assumptions as regular Monte Carlo strength estimation². Thus, whenever we can apply regular Monte Carlo strength estimation to estimate guessing numbers for a password cracking model M we can also apply our Confident Monte Carlo techniques to obtain confidence bounds for the estimated guessing number. Empirical analysis shows that our upper/lower bounds on the guessing number are usually quite close. Typically, we find that the estimates generated by regular Monte Carlo lie within our confidence range, but we also find that for many rare passwords the estimated guessing number generated by regular Monte Carlo is *demonstrably* inaccurate e.g., the estimate lies below our *lower bound*.

Third, we develop rigorous statistical techniques to upper/lower bound the attacker’s entire guessing curve. In particular, given a dataset D we can generate curves $\lambda_{M,B,D}^{ub}$ (resp. $\lambda_{M,B,D}^{lb}$) such that with probability at least $1 - \delta$ for all guessing budgets B we have $\lambda_{M,B,D} \leq \lambda_{M,B,D}^{ub}$ (resp. $\lambda_{M,B,D} \geq \lambda_{M,B,D}^{lb}$).

If we assume that our dataset D was sampled iid from the (unknown) password distribution \mathcal{P} then we can apply McDiarmid’s inequality to argue that (whp) $|\lambda_{M,B,D} - \lambda_{M,B}| \leq \epsilon$. This allows us to confidently upper/lower bound the guessing curve $\lambda_{M,B}$ for our (unknown) user password distribution \mathcal{P} using only a dataset D sampled from this distribution. This observation also provides a rigorous statistical framework for many natural tasks in password research (1) analyzing the impact of a password policy on the guessing curve given samples S_1 (resp. S_2) from the unknown distribution \mathcal{P}_1 (resp. \mathcal{P}_2) representing the

2. We assume that (1) the model M describes a distribution \mathcal{M} over passwords and we can efficiently sample from this distribution, and (2) given a particular password pwd we can quickly compute the probability that this password is generated by M .

distribution of user passwords before (resp. after) the policy intervention, (2) comparing the performance of different password cracking models M and M' against an unknown password distribution \mathcal{P} .

Finally, we evaluate Confident Monte Carlo empirically using several large breached password datasets. We find that our upper/lower bounds on the guessing curve are typically very close and thus tightly bound the true values $\lambda_{M,B,D}$ or $\lambda_{M,B}$. We compare our upper/lower bounds with the popular heuristic estimate using Monte Carlo strength estimation. We find the heuristic estimates for $\lambda_{M,B}$ tend to be accurate when B is small, but as the guessing budget B increases the estimates are *provably inaccurate* i.e., lies above our upper bound.

1.2. Related Work

Password Guessing Models. Offline password attacks have been a concern since the Unix system was devised [15]. Many sophisticated probabilistic password models have been proposed to generate password guesses for an online attacker such as Probabilistic Context-Free Grammars [10], [21], [23], Markov models [5], [6], [12], [19], and neural networks [14]. Each of these probabilistic password models are compatible with regular Monte-Carlo Strength estimation [8] and our Confident Monte Carlo techniques. Thus, we can apply our statistical techniques to derive high confidence upper/lower bounds on guessing numbers for each of these models. By contrast, heuristic (rule-based) tools such as Hashcat [1] and John the Ripper [9] are not compatible with Monte Carlo Strength Estimation. Liu et al. [11] developed tools to estimate guessing numbers for Hashcat and John the Ripper without resorting naïve brute-force enumeration.

Password Strength Estimation. number of any given password outputted by a probabilistic model without requiring the defender to simulate the full attack. Regular Monte Carlo strength estimation [8] has been widely used in the password research literature to understand the impact of culture/language on password strength [22], evaluate the impact of policy interventions such as password composition policies [14], [17], [19], develop password strength meters [18] and evaluate the effectiveness of key-stretching mechanisms against offline attacks [2]. However, to the best of our knowledge the problem of providing rigorous confidence intervals for the estimated guessing numbers has not been explored³. Blocki and Liu [4] recently focused on the problem of upper/lower bounding the guessing curve of a perfect knowledge attacker who *knows* the user password distribution \mathcal{P} . While this can be a useful goal, in practice it can also be useful to characterize the guessing curve of an

attacker following a state of the art password cracking model M since a real world attacker will not have perfect knowledge of the user password distribution. We also note that the upper/lower bounds of Blocki and Liu [4] for the guessing curve of a perfect knowledge attacker rapidly diverge even for moderately large values of B e.g., $B = 10^7$. By contrast, we are able to obtain relatively tight upper/lower bounds on the guessing curve of an attacker using model M even when the guessing budget B is very large e.g., $B = 10^{24}$.

2. Background

Probabilistic Password Guessing Model. In this work we assume that our attacker is untargeted and that the attacker uses a Probabilistic Password Model M to crack passwords. To apply regular Monte Carlo or Confident Monte Carlo we make several assumptions about the model M . First, we assume that the model M implicitly defines a distribution \mathcal{M} over passwords and that M allows us to efficiently sample from the distribution \mathcal{M} . Second, given an arbitrary password pwd we assume that we can efficiently compute $p_{pwd}^M \doteq \Pr_{x \leftarrow \mathcal{M}}[x = pwd]$ i.e., likelihood of the password pwd according to our distribution \mathcal{M} . We note that these assumptions hold for most sophisticated password cracking models such as Probabilistic Context-Free Grammars [10], [21], [23], Markov models [5], [6], [12], [19], and neural networks [14].

It will be convenient to let pwd_1, pwd_2, \dots denote the list of passwords in the support of our distribution \mathcal{M} and to let $p_i^M \doteq p_{pwd_i}^M$ denote the probability of password i . It will also be convenient to assume that these passwords are ordered such that $p_1^M \geq p_2^M \geq \dots$ i.e., so that an attacker using model M would check guesses in the order pwd_1, pwd_2, \dots . We let $G(pwd)$ denote the number of guesses that an attacker, following model M , would need to attempt in order to crack the password pwd i.e., $G(pwd_i) = i$.

Given a probability value $q \in [0, 1]$ we would like to define $G(q)$ as the hypothetical guessing number for a password pwd with probability $p_{pwd}^M = q$. However, if there are multiple passwords in \mathcal{M} with probability exactly q there will be multiple different values of the guessing number. To avoid ambiguity we instead define an exclusive bound $G^{\text{EX}}(q) := |\{i : p_i^M > q\}|$ to count the number of passwords with probability *strictly greater than* and an inclusive bound $G^{\text{IN}}(q) := |\{i : p_i^M \geq q\}|$ to count the number of passwords with probabilities *greater than or equal to* q . Observe that for a password pwd with probability p_{pwd}^M we have $G^{\text{EX}}(q) + 1 \leq G(pwd) \leq G^{\text{IN}}(q)$. It will sometimes be convenient to write $G^{\text{EX}}(pwd) \doteq G^{\text{EX}}(p_{pwd}^M)$ or $G^{\text{IN}}(pwd) \doteq G^{\text{IN}}(p_{pwd}^M)$.

Regular Monte Carlo Estimation. To compute $G(pwd)$ (or $G^{\text{EX}}(q)$ or $G^{\text{IN}}(q)$) exactly a defender would need to enumerate all possible passwords in M whose probability is above a given threshold (p_{pwd}^M). This can be prohibitively expensive for the defender when the guessing number is large. Thus, Dell’Amico and Filippone [8] developed a

3. Melicher et al. [14] mention that with at least one million samples typically they observe “95% confidence intervals of less than 10% of the value of the guess-number estimate” and “passwords for which the error exceeded 10% tended to be guessed only after more than 10^{18} guesses” in their experiments. However, the paper does not contain any details about these claimed confidence intervals or the methodology by which they were derived. We reached out the the authors to provide clarification, but received no response.

Monte Carlo algorithm to efficiently estimate $G(y)$. More accurately, for any probability value q their algorithm produces an unbiased estimate of $G^{\text{EX}}(q)$ — we have $G(pwd) = G^{\text{EX}}(q) + 1$ in whenever there is a unique password pwd with probability $p_{pwd}^M = q$. This regular Monte Carlo algorithm works as follows: (1) draw k iid samples from the distribution \mathcal{M} i.e., $S \leftarrow \mathcal{M}^k$, (2) output the estimate $\hat{G}_S^{\text{EX}}(q) = \frac{1}{k} \sum_{x \in S, p_x^M > q} \frac{1}{p_x^M}$. In practice, Dell’Amico and Filippone [8] proposed that one could draw the sample S ahead of time and use this sample to obtain out strength estimate $\hat{G}_S^{\text{EX}}(pwd)$ for multiple different passwords.

Dell’Amico and Filippone [8] proved that for any probability parameter $q \in [0, 1]$ the expectation of the estimation is equal to its true value, i.e. $\mathbb{E}(\hat{G}_S^{\text{EX}}(q)) = G^{\text{EX}}(q)$. They also argued that the variance $\text{Var}(\hat{G}_S^{\text{EX}}(q)) = \frac{1}{k} (\sum_{i: p_i^M > q} \frac{1}{p_i^M} - G^{\text{EX}}(q))^2$ converges to 0 as the sample size k gets to infinite. However, in practice the sample size k is finite and can be very small compared to $\frac{1}{q}$.

Similarly, as a trivial extension of [8] one can define $\hat{G}_S^{\text{IN}}(q) := \frac{1}{k} \sum_{x \in S, p_x^M \geq q} \frac{1}{p_x^M}$ as an unbiased estimate for our inclusive term $G^{\text{IN}}(q)$.

Password Guessing Curve. Given a dataset D^4 we let $\lambda_{M,B,D} \doteq \{x \in D : G(x) \leq B\}$ denote the fraction of passwords in D that would be cracked within B guesses by an untargeted attacker following model M . Similarly, we define $\lambda_{M,B} := \Pr_{y \leftarrow \mathcal{P}}[G(y) \leq B]$ to be the probability that a random password y sampled from \mathcal{P} would be cracked within y guesses. We will typically assume that the user password distribution \mathcal{P} is unknown, but that we are given a dataset D consisting of iid samples from \mathcal{P} . In this case we have $\lambda_{M,B} = \mathbb{E}[\lambda_{M,B,D}]$ where the randomness is taken over the selection of D from the unknown distribution \mathcal{P} .

For an online attacker B is usually small since an authentication service can lock out the user account after several failed login attempts. For an offline attacker B can be much larger since with the stolen (salted) cryptographic hash of the user’s password an offline attacker can check as many passwords as s/he wants by comparing the (salted) cryptographic hash with the hashes of the top B guesses pwd_1, \dots, pwd_B . An offline attacker is limited only by the resources s/he is willing to invest cracking and by the cost of repeatedly evaluating the password hash function.

3. Limitations of Regular Monte Carlo Strength Estimation

In this section we discuss the limitations of the regular Monte Carlo strength estimation [8]. We provide an example of password models where regular Monte Carlo will dramatically underestimate the guessing number with high probability, and another example of a password model where Monte Carlo will overestimate the guessing number by a factor of 2 with probability at least 0.3. Specifically, for any sample size k we define a model for which the estimation is

4. D may contain duplicated passwords since different users might select the same password.

inaccurate with a significant probability. We also provide an example of a password model and a password distribution where the regular Monte Carlo estimation has a large error on predicting the guessing curve. We will further explain that the above issue is inherent to any blackbox method for estimation given only a finite amount of samples without exploiting specific properties of the password guessing model.

3.1. Error on Guessing Number $G(y)$

3.1.1. Underestimating the Guessing Number. We first provide an example of a password model where Monte Carlo will (whp) dramatically underestimate the true guessing number.

The Model/Distribution \mathcal{M} : Consider a model M which induces a distribution \mathcal{M} over passwords pwd_1, pwd_2, \dots such that $\Pr_{x \leftarrow \mathcal{M}}[x = pwd_i] = 2^{-i}$.

Actual Guessing Number: The actual guessing number of each password pwd_i is $G(y) = i$.

Analysis of Monte Carlo Estimate: Suppose that we apply Monte Carlo with $k = |S| \geq 2^{10}$ iid samples $S \leftarrow \mathcal{M}^k$ from \mathcal{M} to estimate the guessing number of a password. For any $i > 2 \log(k)$ with high probability $(1 - 2^{-i+1})^k > 0.99$ our sample set S will contain *no* passwords with probability equal to or less than 2^{-i} . In this case for any $j \geq i$ our estimated guessing number for pwd_j is $\hat{G}_S^{\text{EX}}(pwd_j) = \frac{1}{k} \sum_{x \in S} \frac{1}{p_x}$ where $\frac{1}{k} \sum_{x \in S} \frac{1}{p_x} \leq \frac{2^{\log(k^2)}}{k} = k$. Thus, if $j \gg k$ Monte Carlo will *dramatically underestimate* the true guessing number with high probability.

The authors of [8] state that the variance of $\hat{G}_S^{\text{EX}}(pwd_j)$ approaches 0 as k increases. However, in the above example the variance for the j th most probable password pwd_j is $\text{Var}(\hat{G}_S^{\text{EX}}(pwd_j)) = \frac{1}{k} (\sum_{t=1}^{j-1} \frac{1}{2^{-t}} - (j-1)^2) = \frac{1}{k} (2^j - 1 - (j-1)^2)$. While the variance does decrease linearly with the sample size k it also increases exponentially with the true guessing number j .

3.1.2. Overestimating the Guessing Number. We now give an example where regular Monte Carlo estimation can *overestimate* the guessing number by a factor of ≈ 2 with non-negligible probability (e.g. 0.3).

The Model/Distribution \mathcal{M} : Fix the sample size $k \geq 5$ and consider a model M corresponding to a password distribution \mathcal{M} where there are $n_1 = 2k - 1$ passwords with probability $p_1 = \frac{1}{2k}$, $n_2 = \frac{k^2}{2} - 1$ passwords with probability $p_2 = \frac{1}{k^3}$, and $n_i = 1$ password with probability $p_i = \frac{1}{2^{i-2}k^3} < p_{i-1}$ for any $i \geq 3$. Observe that $\sum_{i \geq 1} n_i p_i = 1$.

Actual Guessing Number: The actual guessing number of the password with probability p_3 is $G(p_3) = n_1 + n_2 + 1 = 2k - 1 + \frac{k^2}{2}$.

Analysis of Monte Carlo Estimate: Suppose $|S| = k$ ($k \geq 5$) samples are generated to estimate the guessing numbers. Consider the event E that S contains exactly one password with probability p_2 and all the remaining $(k -$

1) samples have probability p_1 (i.e., none of the remaining samples in S have probability smaller than p_2). If the event E occurs then the estimated guessing number of password with probability p_3 is $\hat{G}_S^{\text{EX}}(p_3) + 1 = \frac{1}{k}((k-1)\frac{1}{p_1} + \frac{1}{p_2}) + 1 = 2k - 1 + k^2$ while the actual guessing number is $G(p_3) = n_1 + n_2 + 1 = 2k - 1 + \frac{1}{2}k^2$ i.e., we overestimate the guessing number by a factor $\frac{\hat{G}_S^{\text{EX}}(p_3) + 1}{G(p_3)} \approx 2$. The probability of the event E is $n_2 \binom{k}{1} p_2 (n_1 p_1)^{k-1} = (\frac{1}{2} - \frac{1}{k^2})(1 - \frac{1}{2k})^{k-1} > 0.3$.

3.1.3. Inherent Limitations of Any Blackbox Estimate.

We now argue that *any* blackbox method for estimating the guessing number will have similar issues. In particular, for any sample size k and any blackbox strength estimation method we can find cases where the model will (whp) give us highly inaccurate strength estimates. As a concrete example consider the distribution $\mathcal{M} = 2^{-1}, 2^{-2}, 2^{-3}, \dots, 2^{-i}, \dots$ defined in Section 3.1.1. Suppose we define a similar distribution \mathcal{M}' to match \mathcal{M} on the first i passwords i.e., for any $j \leq i$ the password j has probability 2^{-j} . However, \mathcal{M}' has 2^{2i-1} passwords with probability 2^{-3i} (total probability mass 2^{-i-1}) and 2^{2i} passwords with probability 2^{-3i-1} (total probability mass 2^{-i-1}). the two distributions \mathcal{M} and \mathcal{M}' produce *very different* guessing numbers for many passwords. For example, consider a password x with probability $p_x = 2^{-3i-1}$. Under model \mathcal{M} the actual guessing number would simply be $G(x) = 3i + 1$. By contrast, under model \mathcal{M}' the actual guessing number for this password would be $G(x) \geq 2^{2i-1}$. In this example it is simply not possible to confidently generate an accurate estimate of the true guessing number. However, it is not even possible to distinguish between $k = o(2^i)$ samples drawn from \mathcal{M} and k samples drawn from \mathcal{M}' except with probability $o(1)$. To see this let E denote the event that we never sample a password y with probability $p_y < 2^{-i}$. Observe that if we condition on the event E the two models \mathcal{M} and \mathcal{M}' are equivalent. Furthermore, by union-bounds we have $\Pr[\bar{E}] \leq k2^{-i} = o(1)$. If the event E occurs while we were trying to estimate the guessing number for our password x (with $p_x = 2^{-3i-1}$) then it will not be possible to accurately approximate $G(x)$. In this case could be particularly problematic if our strength estimator outputs an estimated guessing number for x without providing confidence ranges or providing any indication that there is a high degree of uncertainty.

3.2. Error on Guessing Curve

Another application of Monte Carlo strength estimation is generating (an approximation of) the attacker's guessing curve. For example, we might like to estimate the probability $\lambda_{M,B}$ that an attacker can crack a random user's password (sampled from potentially unknown distribution \mathcal{P}) within B guesses using model M or, given a dataset D of user passwords we might want to estimate $\lambda_{M,B,D}$ — the fraction of passwords in D that the attacker using model M can crack within B guesses per user account. The standard (heuristic) way to approximate $\lambda_{M,B,D}$ and $\lambda_{M,B}$ is to

simply compute $\hat{\lambda}_{M,B,D} := \frac{1}{|D|} |\{y \in D : \hat{G}^{\text{EX}}(y) \leq B\}|$. This heuristic approach has been widely adopted in the password literature (e.g., [2], [11], [14], [16]). However, for any guessing budget B and sampling parameter k we can provide an example of a model M and a distribution \mathcal{P} over user passwords such that: (1) An attacker who makes B guesses will never crack any user password sampled from \mathcal{P} i.e., $\lambda_{M,B} = 0$, and (2) with high probability $\hat{\lambda}_{M,B,D} = 1$ i.e., the widely adopted heuristic using regular Monte Carlo strength estimation predicts that the attacker cracks 100% of passwords within B guesses.

Model Distribution \mathcal{M} : Consider a distribution \mathcal{M} generated by a cracking model M where there are $k^2 - 1$ passwords $pwd_1, \dots, pwd_{k^2-1}$ with probability $\frac{1}{k^2}$. The remaining k^2 passwords $pwd_{k^2}, \dots, pwd_{2k^2-1}$ in the support of our distribution satisfy the following properties: (1) $\sum_{i=k^2}^{2k^2-1} \Pr_{x \leftarrow \mathcal{M}}[x = pwd_i] = \frac{1}{k^2}$ and (2) $\Pr_{x \leftarrow \mathcal{M}}[x = pwd_i] > \Pr_{x \leftarrow \mathcal{M}}[x = pwd_{i+1}]$ for each $k^2 - 1 \leq i < 2k^2 - 1$ i.e., passwords are listed in descending order of probability. Let $B = k^2 + 1$, let $F = \{pwd_1, \dots, pwd_B\}$ denote the most popular B passwords and let $L = \{pwd_{B+1}, \dots, pwd_{k^2-1}\}$ denote the remaining passwords.

Actual Password Distribution \mathcal{P} : For the actual user password distribution we can consider any distribution \mathcal{P} with support L i.e., $\Pr_{x \leftarrow \mathcal{P}}[x \in F] = 0$.

Actual Guessing Curve: Notice that the attacker's first B guesses will all be from the set F . By construction, the support of our user password distribution \mathcal{P} does not include any password from F . Thus, for any dataset D of passwords sampled from \mathcal{P} we will have $\lambda_{M,B,D} = 0$ and we also have $\lambda_{M,B} = \mathbb{E}[\lambda_{M,B,D}] = 0$.

Analysis of Monte Carlo Curve: Given a sample set S with size k from \mathcal{M} , except with probability at most $\frac{1}{k}$, all k samples are from $pwd_1, \dots, pwd_{k^2-1}$. In this case the estimated guessing number for *any* password pwd_i with $k^2 - 1 \leq i \leq 2k^2 - 1$ is $\hat{G}_S^{\text{EX}}(pwd_i) + 1 = 1 + \frac{1}{k} \sum_{x \in S} k^2 = 1 + k^2 \leq B$. Thus, with probability at least $1 - \frac{1}{k}$, we have $\hat{\lambda}_{M,B,D} = 1$ meaning that this widely adopted heuristic incorrectly predicts that the attacker will crack 100% of passwords in our dataset D .

4. Theoretical Bounds on Guessing Number

In this section, we present two techniques to rigorously bound the actual guessing number $G(y)$ of an arbitrary password y using iid password samples randomly selected from a model distribution \mathcal{M} generated by a password cracking model M . In particular, we can ensure that the actual guessing number $G(y)$ is sandwiched between our upper/lower bounds with high confidence (e.g., 99%) allowing us to quantify the uncertainty of guessing number estimates.

4.1. Upper/Lower Bounds Using Hoeffding's Inequality

In this section, we prove an upper bound and a lower bound on the actual guessing number $G(q)$ of a password

with probability q using Hoeffding's inequality. The formal statement is shown below:

Theorem 1. *Given a set S with k iid password samples sampled from distribution \mathcal{M} , for any probability $q \in [0, 1]$ and any parameter $\epsilon \geq 0$, we have:*

$$\begin{aligned}\Pr[G(q) \leq \hat{G}_S^{\text{IN}}(q) + \epsilon/q] &\geq 1 - \exp(-2k\epsilon^2) \\ \Pr[G(q) \geq \hat{G}_S^{\text{EX}}(q) + 1 - \epsilon/q] &\geq 1 - \exp(-2k\epsilon^2)\end{aligned}$$

where the randomness is taken over the sample set $S \leftarrow \mathcal{M}^k$.

Due to space limitations the formal proof is deferred to Appendix D.1. Briefly, we observe that $\hat{G}_S^{\text{EX}}(q)$ can be viewed as the sum of k independent random variables $X_1^{\text{EX}}, \dots, X_k^{\text{EX}}$ where $X_i^{\text{EX}} = 0$ if the i th password z from our sample has probability $p_z^M > q$ and $X_i^{\text{EX}} = 1/p_z^M$ otherwise. Since the random variable are independent we can apply Hoeffding's inequality to upper bound the probability that our estimate $\hat{G}_S^{\text{EX}}(q)$ is significantly smaller than its expectation $\mathbb{E}[\hat{G}_S^{\text{EX}}(q)] = G_S^{\text{EX}}(q)$ where $G(q) \geq G_S^{\text{EX}}(q) + 1$. Similarly, we can apply Hoeffding's inequality to bound $\hat{G}_S^{\text{IN}}(q)$. In particular, for $\phi \in \{\text{EX}, \text{IN}\}$ we have $|\hat{G}_S^\phi(q) - G^\phi(q)| \leq \epsilon/q$ with any $\epsilon \geq 0$ as below:

$$\Pr[G^\phi(q) \leq \hat{G}_S^\phi(q) + \epsilon/q] \geq 1 - \exp(-2k\epsilon^2) \quad (1)$$

$$\Pr[G^\phi(q) \geq \hat{G}_S^\phi(q) - \epsilon/q] \geq 1 - \exp(-2k\epsilon^2) \quad (2)$$

Note that given a password y the upper/lower bounds in Theorem 1 differ by an additive factor of ϵ/p_y^M . Thus, we can obtain tight upper/lower bounds $\epsilon/p_y^M \ll G(y)$ although the bounds can diverge as p_y^M grows small i.e., when the password is particularly rare.

4.2. A Tighter Lower Bound For Rare Passwords

The upper and lower bounds in Section 4.1 will diverge when the password is rare. In the worst case, the lower bound will be useless if $\hat{G}_S^{\text{EX}}(q) + 1 - \epsilon/q$ becomes a negative value. Is it possible to derive tighter bounds on the guessing numbers of rare passwords? In this section we further tighten the lower bound for rare passwords by applying Markov's inequality and taking the median estimates. In particular, for any password probability $q \in [0, 1]$ we define $\hat{G}_{\mathbb{S}, \text{med}}^\phi(q) = \text{median}(\{\hat{G}_{S_i}^\phi(q)\}_{1 \leq i \leq n})$ where $\mathbb{S} = \{S_1, S_2, \dots, S_n\}$ and each S_i contains k independent samples from our model. Fixing any parameter $\delta < \frac{1}{2}$ we show that (whp) $\delta \cdot \hat{G}_{\mathbb{S}, \text{med}}^{\text{EX}}(q)$ lower bounds the true guessing number as stated in the following theorem:

Theorem 2. *For any password probability $q \in [0, 1]$, and any parameters $0 \leq \delta \leq \frac{1}{2}$ and $0 \leq \epsilon \leq \frac{1}{2} - \delta$,*

$$\Pr[G(q) \geq \delta \cdot \hat{G}_{\mathbb{S}, \text{med}}^{\text{EX}}(q) + 1] \geq 1 - \exp(-2n\epsilon^2)$$

where the randomness is taken over n sets of k Monte Carlo samples $\mathbb{S} = \{S_1, \dots, S_n\}$ from model M .

Theorem 2 often allows us to tighten the lower bounds on the guessing number of rare passwords. Due to space limitations the formal proof is deferred to Appendix D.2. Intuitively, we can define an indicator random variable X_i^ϕ for $\phi \in \{\text{EX}, \text{IN}\}$ such that $X_i^\phi = 1$ if and only if $G^\phi(q) \geq \delta \cdot \hat{G}_{S_i}^\phi(q)$ i.e., if and only if the i th sample set S_i overestimates $G^\phi(q)$ by at most a factor of $1/\delta$. As long as $\sum_{i=1}^n X_i^\phi \geq n/2$ we are guaranteed that the median estimate is not too bad and that $G(q) \geq \delta \cdot \hat{G}_{\mathbb{S}, \text{med}}^{\text{EX}}(q) + 1$. Thus, it suffices to upper bound the probability that $\sum_{i=1}^n X_i^\phi \geq n/2$. We first apply Markov's inequality to show that $\Pr[X_i^\phi = 1] \geq 1 - \delta$. Finally, since each X_i^ϕ is independent we can apply concentration bounds to show that $\sum_{i=1}^n X_i^\phi \geq n/2$ with high probability. In particular, by Chernoff bounds for any $0 \leq \epsilon \leq \frac{1}{2} - \delta$ we have

$$\begin{aligned}\Pr[\hat{G}_{\mathbb{S}, \text{med}}^\phi(q) \leq \frac{1}{\delta} G^\phi(q)] &\geq \Pr[\sum_{i=1}^n X_i^\phi \geq \frac{n}{2}] \\ &\geq \Pr[\sum_{i=1}^n X_i^\phi \geq n(1 - \delta - \epsilon)] \geq 1 - \exp(-2n\epsilon^2).\end{aligned} \quad (3)$$

5. Theoretical Bounds on Fraction of Cracked Passwords

So far our focus has been on upper/lower bounding the guessing number of a particular user password against a cracking model M . However, in some defensive applications our goal will be to upper/lower bound the attacker's entire guessing curve e.g., to determine whether or not a password policy intervention resulted in a password distribution that is harder for the adversary to crack. More formally, in this section we develop techniques to upper/lower bound the curves $\lambda_{M, B, D}$ and $\lambda_{M, B}$ as the guessing budget B varies from small to large. Recall that $\lambda_{M, B, D}$ denotes the fraction of passwords in dataset D that would be cracked within B guesses, and that $\lambda_{M, B}$ denotes the probability that randomly sampled password would be cracked within B guesses.

Given a dataset D of independent samples from an *unknown* password distribution, our first observation is that the expected value of $\lambda_{M, B, D}$ (over the random selection of D) is simply $\lambda_{M, B}$ and, if D is large enough, the random variable $\lambda_{M, B, D}$ is tightly concentrated around its mean — see Theorem 3. Given this result our main task will be to develop high confidence upper/lower bounds on $\lambda_{M, B, D}$ which will immediately yield high-confidence upper/lower bound for $\lambda_{M, B}$ as a corollary. Thus, we will focus primarily on bounding $\lambda_{M, B, D}$ in the remainder of this section. Theorem 3 follows directly from McDiarmid's inequality [13]. The formal proof is deferred to Appendix E.

Theorem 3. *For any guessing number $B \geq 0$ and any $0 \leq$*

$\epsilon \leq 1$, we have:

$$\begin{aligned} \Pr[\lambda_{M,B} \geq \lambda_{M,B,D} - \epsilon] &\geq 1 - \exp(-2|D|\epsilon), \quad \text{and} \\ \Pr[\lambda_{M,B} \leq \lambda_{M,B,D} + \epsilon] &\geq 1 - \exp(-2|D|\epsilon) \end{aligned}$$

where the randomness is taken over the sample set $D \leftarrow \mathcal{P}^{|D|}$.

5.1. The General Framework

In this section, we propose a generalized framework for converting confident upper/lower bounds on guessing numbers into confident upper and lower bounds on $\lambda_{M,B,D}$ (and by extension $\lambda_{M,B}$). Suppose that $G(q)$ denotes the guessing number for a password pwd whose probability (according to our model) is q — for simplicity of exposition let us first suppose that there is only one such password with probability exactly q . Although $B = G(q)$ is unknown we observe that it is still possible to compute the quantity $\lambda_{M,B,D}$ i.e., by computing the fraction of passwords in D (i.e. $pw \in D$) whose probability is $p_{pw}^M \geq q$. Unfortunately, this is still not sufficient to plot the curve $\lambda_{M,B,D}$ since we do not actually know the value of B . However, if we are given upper/lower bounds $L \leq B \leq U$ then we can use the value of $\lambda_{M,B,D}$ as an upper bound for $\lambda_{M,L,D}$ and as a lower bound for $\lambda_{M,U,D}$ since we know that $\lambda_{M,L,D} \leq \lambda_{M,B,D} \leq \lambda_{M,U,D}$.

Our key idea is to pick a sequence q_1, q_2, \dots, q_ℓ of ℓ probability mesh points and obtain upper (resp. lower) bounds U_1, \dots, U_ℓ (resp. L_1, \dots, L_ℓ) on the corresponding guessing numbers. Intuitively, as long as all of our upper (resp. lower) bounds are valid we can use them to lower (resp. upper) bound the guessing curve $\lambda_{M,B,D}$ at multiple points $B \in \{U_1, \dots, U_\ell\}$ (resp. $B \in \{L_1, \dots, L_\ell\}$).

The formal framework is slightly complicated by the fact that we will occasionally have multiple passwords with the same probability q according to model M . However, we can deal with this concern using by lower (resp. upper) bounding the exclusive (resp. inclusive) guessing numbers. Recall that for any $\phi = \{\text{EX}, \text{IN}\}$ equations (1), (2) and (3) provide high confidence upper and lower bounds on $G^\phi(q)$. In general, for any probability $0 \leq q \leq 1$ we define $\text{UB}_{G^\phi, |S|}(q, S)$ (resp. $\text{LB}_{G^\phi, |S|}(q, S)$) to be an arbitrary upper (resp. lower) bound of $G^\phi(q)$ with error rate $\text{ERR}(\text{UB}_{G^\phi, |S|})$ (resp. $\text{ERR}(\text{LB}_{G^\phi, |S|})$), i.e., with randomness taken over the selection of sample set S from model M we have:

$$\Pr[G^\phi(q) \geq \text{LB}_{G^\phi, |S|}(q, S)] \geq 1 - \text{ERR}(\text{LB}_{G^\phi, |S|}) \quad (4)$$

$$\Pr[G^\phi(q) \leq \text{UB}_{G^\phi, |S|}(q, S)] \geq 1 - \text{ERR}(\text{UB}_{G^\phi, |S|}) \quad (5)$$

Formally, we define $\hat{\lambda}_{M,B,D,S}^{ub}$ and $\hat{\lambda}_{M,B,D,S}^{lb}$ as below:

$$\hat{\lambda}_{M,B,D,S}^{ub} := \min_{1 \leq i \leq \ell, B \leq \text{LB}_{G^{\text{IN}}, |S|}(q_i, S)} (\lambda_{M, G^{\text{IN}}(q_i), D}), \quad (6)$$

$$\hat{\lambda}_{M,B,D,S}^{lb} := \max_{1 \leq i \leq \ell, B \geq \text{UB}_{G^{\text{EX}}, |S|}(q_i, S)} (\lambda_{M, G^{\text{EX}}(q_i), D}). \quad (7)$$

For completeness, we set our upper (resp. lower) bound $\hat{\lambda}_{M,B,D,S}^{ub} = 1$ (resp. $\hat{\lambda}_{M,B,D,S}^{lb} = 0$) if $B > \max_{1 \leq i \leq \ell} \{\text{LB}_{G^{\text{IN}}, |S|}(q_i, S)\}$ (resp. $B < \min_{1 \leq i \leq \ell} \{\text{UB}_{G^{\text{EX}}, |S|}(q_i, S)\}$).

Theorem 4 shows that $\hat{\lambda}_{M,B,D,S}^{ub}$ (resp. $\hat{\lambda}_{M,B,D,S}^{lb}$) upper (resp. lower) bounds the value of $\lambda_{M,B,D}$ with high confidence. Intuitively, the error terms $\ell \cdot \text{ERR}(\text{LB}_{G^{\text{IN}}, |S|})$ and $\ell \cdot \text{ERR}(\text{UB}_{G^{\text{EX}}, |S|})$ are obtained by taking union bounds.

Theorem 4. *Given a password dataset D containing iid samples from distribution \mathcal{P} and a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$, we have:*

$$\Pr[\forall B \geq 0, \lambda_{M,B,D} \leq \hat{\lambda}_{M,B,D,S}^{ub}] \geq 1 - \ell \cdot \text{ERR}(\text{LB}_{G^{\text{IN}}, |S|})$$

$$\Pr[\forall B \geq 0, \lambda_{M,B,D} \geq \hat{\lambda}_{M,B,D,S}^{lb}] \geq 1 - \ell \cdot \text{ERR}(\text{UB}_{G^{\text{EX}}, |S|})$$

where the randomness is taken over the sample set S from model M .

If we assume that our dataset D is sampled iid from our unknown password distribution we can apply Theorem 3 to upper/lower bound $\lambda_{M,B}$ as an immediate corollary of Theorem 4 — see Corollary 5.

Corollary 5. *Given a password distribution \mathcal{P} and a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$, for any guessing number $B > 0$ and any parameters $0 \leq \epsilon \leq 1$, we have:*

$$\begin{aligned} \Pr[\lambda_{M,B} \leq \hat{\lambda}_{M,B,D,S}^{ub} + \epsilon] \\ \geq 1 - \ell \cdot \text{ERR}(\text{LB}_{G^{\text{IN}}, |S|}) - \exp(-2|D|\epsilon^2) \end{aligned}$$

$$\begin{aligned} \Pr[\lambda_{M,B} \geq \hat{\lambda}_{M,B,D,S}^{lb} - \epsilon] \\ \geq 1 - \ell \cdot \text{ERR}(\text{UB}_{G^{\text{EX}}, |S|}) - \exp(-2|D|\epsilon^2) \end{aligned}$$

where the randomness is taken over the sample set S from model M and the dataset $D \leftarrow \mathcal{P}^{|D|}$.

The formal proofs can be found in Appendix E.1.

5.2. Concrete Bounds on Guessing Curves

We now derive concrete upper and lower bounds on the attacker's guessing curves ($\lambda_{M,B,D}$ and $\lambda_{M,B}$) by applying our general framework from Section 5.1 with our concrete upper/lower bounds on the guessing numbers $G^{\text{EX}}(q)$ and $G^{\text{IN}}(q)$ from Section 4. Due to space limitations, we only present the concrete bounds on $\lambda_{M,B}$ in this section and defer formal proofs and the formal statements of concrete bounds on $\lambda_{M,B,D}$ to Appendix E.2.

First Concrete Upper/Lower Bound. We first apply Theorem 4 to the upper and lower bounds on $G^{\text{EX}}(q)$ and $G^{\text{IN}}(q)$ from equations (1) and (2). In particular, we define $\hat{\lambda}_{M,B,D,S,\epsilon}^{ub1}$ and $\hat{\lambda}_{M,B,D,S,\epsilon}^{lb1}$ as the concrete upper/lower bounds on $\lambda_{M,B,D}$:

$$\hat{\lambda}_{M,B,D,S,\epsilon}^{ub1} := \min_{1 \leq i \leq \ell, B \leq \hat{G}_S^{\text{IN}}(q_i) - \epsilon/q_i} (\lambda_{M, G^{\text{IN}}(q_i), D}),$$

$$\hat{\lambda}_{M,B,D,S,\epsilon}^{lb1} := \max_{1 \leq i \leq \ell, B \geq \hat{G}_S^{\text{EX}}(q_i) + \epsilon/q_i} (\lambda_{M, G^{\text{EX}}(q_i), D}).$$

For completeness, we set $\hat{\lambda}_{M,B,D,S,\epsilon}^{ub1} = 1$ (resp. $\hat{\lambda}_{M,B,D,S,\epsilon}^{lb1} = 0$) if $B > \max_{1 \leq i \leq \ell} \{\hat{G}_S^{\text{IN}}(q_i) - \epsilon/q_i\}$ (resp. $B < \min_{1 \leq i \leq \ell} \{\hat{G}_S^{\text{EX}}(q_i) + \epsilon/q_i\}$).

By Theorem 4 it follows that with probability at least $\ell \cdot \exp(-2k\epsilon^2)$ that $\hat{\lambda}_{M,B,D,S,\epsilon}^{ub1}$ (resp. $\hat{\lambda}_{M,B,D,S,\epsilon}^{lb1}$) is an upper (resp. lower) bound on $\lambda_{M,B,D}$ for every $B \geq 1$ where the randomness depends only on the selection of k samples $|S| = k$ from our model M — see Theorem 12 in Appendix E.2.1 for the formal statement. If we additionally assume that our dataset D was sampled iid from our unknown password distribution \mathcal{P} then we can apply Theorem 3 (or Corollary 5) and use $\hat{\lambda}_{M,B,D,S,\epsilon}^{ub1}$ (resp. $\hat{\lambda}_{M,B,D,S,\epsilon}^{lb1}$) to upper (resp. lower) bound $\lambda_{M,B}$ for every $B \geq 1$. To ensure that we obtain high confidence bounds we include a small additional slack term (ϵ_2) to account for sampling error selecting our dataset D — see Theorem 6.

Theorem 6. *Given a password distribution \mathcal{P} and a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$, for any guessing number $B > 0$ and any parameters $0 \leq \epsilon_1, \epsilon_2 \leq 1$, we have:*

$$\begin{aligned} \Pr \left[\lambda_{M,B} \leq \hat{\lambda}_{M,B,D,S,\epsilon_1}^{ub1} + \epsilon_2 \right] &\geq 1 - \alpha \\ \Pr \left[\lambda_{M,B} \geq \hat{\lambda}_{M,B,D,S,\epsilon_1}^{lb1} - \epsilon_2 \right] &\geq 1 - \alpha \end{aligned}$$

where $\alpha = \ell \cdot \exp(-2k\epsilon_1^2) - \exp(-2|D|\epsilon_2^2)$ and the randomness is taken over the sample set $S \leftarrow \mathcal{M}^k$ with size k and the dataset $D \leftarrow \mathcal{P}^{|D|}$.

As long as the number of mesh points $\ell = |Q|$ is not too large and the sample size k (and $|D|$) is not too small, both the upper and lower bounds will hold with high probability. **Second Upper Bound.** Recall that in equation (3) we derived a second guessing number lower bound using Markov's inequality and concentration bounds. This lower bound can be effective for rare passwords. We can use this lower bound on the guessing number to derive a second upper bound on the attacker's guessing curve. In particular, we define $\hat{\lambda}_{M,B,D,S,\delta}^{ub2}$ as another upper bound on $\lambda_{M,B,D}$:

$$\hat{\lambda}_{M,B,D,S,\delta}^{ub2} := \min_{1 \leq i \leq \ell, B \leq \delta \cdot \hat{G}_{S,\text{med}}^{\text{IN}}(q_i)} (\lambda_{M,G^{\text{IN}}(q_i),D})$$

As before we set $\hat{\lambda}_{M,B,D,S,\delta}^{ub2} = 1$ whenever $B > \max_{1 \leq i \leq \ell} \delta \cdot \hat{G}_{S,\text{med}}^{\text{IN}}(q_i)$. Applying Theorem 4 we can conclude that with high probability we have $\lambda_{M,B,D} \leq \hat{\lambda}_{M,B,D,S,\delta}^{ub2}$ for all $B \geq 1$ — see Theorem 13 in Appendix E.2.2 for the formal statement. If we additionally assume that the dataset D was sampled iid from our (unknown) password distribution \mathcal{P} we can additionally upper bound $\lambda_{M,B}$ as in Theorem 7.

Theorem 7. *Given a password distribution \mathcal{P} and a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$, for any guessing number $B > 0$ and any parameters $0 < \delta \leq \frac{1}{2}$, $0 \leq \epsilon_1 \leq \frac{1}{2} - \delta$, $0 \leq \epsilon_2 \leq 1$, we have:*

$$\Pr \left[\lambda_{M,B} \leq \hat{\lambda}_{M,B,D,S,\delta}^{ub2} + \epsilon_2 \right] \geq 1 - \alpha$$

where $\alpha = \ell \cdot \exp(-2n\epsilon_1^2) - \exp(-2|D|\epsilon_2^2)$ and the randomness is taken over n Monte Carlo sample sets $\mathbb{S} = \{S_1, S_2, \dots, S_n\}$ each of which contains k samples from model M and the dataset $D \leftarrow \mathcal{P}^{|D|}$.

5.3. A Trivial Upper Bound for Large Guessing Number

The previous section presents two upper bounds and one lower bound on $\lambda_{M,B,D}$ using a series of mesh points q_1, \dots, q_ℓ . However, when B gets large (i.e., $B > \max(\hat{G}_{\text{med},S}(q_\ell) - \epsilon/q_\ell, \delta \cdot \hat{G}_{\text{med},S}(q_\ell))$) we will run out of mesh points q_1, \dots, q_ℓ and the two upper bounds will immediately jump to 1.

Note that if a password y is never outputted by a password cracking model M (i.e., $p_y^M = 0$), then the attacker using M will not be able to successfully guess this password in a dataset D . Denote $\hat{\lambda}_{M,D}^{ub3} := \frac{1}{|D|} |y \in D : p_y^M > 0|$ be the percentage of passwords that will be eventually guessed by model M . Trivially, we have $\lambda_{M,B,D} \leq \hat{\lambda}_{M,D}^{ub3}$ for all finite guessing budgets $B \geq 0$ — see Theorem 14 in Appendix E.3 for the formal theorem statement. We can also obtain the following upper bound on $\lambda_{M,B}$:

Theorem 8. *Given a password distribution \mathcal{P} and a password cracking model M , for any parameters $0 \leq \epsilon \leq 1$, we have:*

$$\Pr \left[\forall B \geq 0, \lambda_{M,B} \leq \hat{\lambda}_{M,D}^{ub3} + \epsilon \right] \geq 1 - \exp(-2|D|\epsilon^2)$$

where the randomness is taken over the dataset $D \leftarrow \mathcal{P}^{|D|}$.

Theorem 8 is derived by the observation that $\hat{\lambda}_{M,D}^{ub3}$ is concentrated on its expectation $\sum_{y \in \mathcal{P}} p_y^M$ which is the total probability mass of passwords in distribution \mathcal{P} that will be guessed with non-zero probability in model M , and $\lambda_{M,B} \leq \lambda_{M,\infty} = \sum_{y \in \mathcal{P}} p_y^M$. The formal proof is in Appendix E.3.

We remark that for the neural network models we consider we have $p_{pw}^M > 0$ for every password $pw \in D$ in our datasets. Thus, the upper bound $\hat{\lambda}_{M,D}^{ub3} = 1$ becomes trivial. However, as we will show in Section 6, for some other probabilistic models such as Markov Models and PCFG, over 20% and 40% of passwords in some of the datasets we consider had $p_{pw}^M = 0$ indicating that these passwords will never be guessed by these particular cracking models. In these cases, our trivial upper bound $\hat{\lambda}_{M,D}^{ub3}$ can yield tighter bounds for large guessing number B .

5.4. Password Composition Policies

Some organizations impose restrictions (password composition policies) on the passwords that user's are allowed to select e.g., users may be required to include numbers, special symbols and/or capital letters. Even if our model M was trained entirely on passwords that are consistent with the policy \mathbb{C} it is still possible that some of the guesses generated by the model will be inconsistent with \mathbb{C} . In this case a trivial optimization for the attacker would be

to simply filter out inconsistent guesses since they cannot appear in our password dataset D or in the support of our user password distribution \mathcal{P} . Intuitively, our bounds work by sampling $|S| = k$ passwords from our model M and then filtering to obtain $S' \subseteq S$ the subset of passwords which are consistent with our policy. We can then show that $\frac{1}{k} \sum_{z \in S', p_z^M > q} \frac{1}{p_z^M}$ is an unbiased estimate of the updated (exclusive) guessing number after filtering. As before we can apply concentration bounds to argue that (whp) the actual guessing value will be close to our estimate.

We show how our statistical techniques can be extended to provide confident upper/lower bounds on the guessing numbers *after* this filtering step. We can then apply our general framework from Section 5.1 to upper/lower bound the attacker’s updated guessing curves $\lambda_{M,B}^C$ and $\lambda_{M,B,D}^C$ after filtering out password guesses that are inconsistent with \mathcal{C} . See Appendix B for formal claims (Theorems 9, 10, 20) about our high confidence upper and lower bounds on $\lambda_{M,B}^C$. See Appendix F and G for all formal proofs and all the remaining theorems of bounding guessing numbers and $\lambda_{M,B,D}^C$.

6. Empirical Experiments

In this section we apply our statistical techniques to upper/lower bound the guessing numbers for user passwords and to upper/lower bound the attacker’s guessing curve $\lambda_{M,B,D}$ and $\lambda_{M,B}$ as the guessing budget B varies from small to large ⁵. To apply our statistical techniques we need to fix a password cracking model M and a password dataset D .

Password Cracking Models We consider 3 generative probabilistic models: Transformer neural network [20], 4-gram Markov model [6] and PCFG [23] (probabilistic context free grammar), each representing a different category of password cracking models. For Markov Models and PCFGs we use the same implementations as [7]. We chose Transformer as the representative of neural network because Transformer is the state-of-the-art machine learning model in learning sequential data. It is faster in training and sampling than RNN [14]; also, it was more efficient in guessing strong passwords in our local tests. The structure and hyperparameters in training Transformer model can be found in the Appendix A.1. We did not expend significant effort optimizing our password cracking models as our primary focus is demonstrating how our statistical techniques can be applied to obtain tight upper/lower bounds on guessing numbers and the attacker’s guessing curve $\lambda_{M,B}$.

Password Datasets We consider six breached password datasets in our experiments: Bfield (0.54m user accounts), Brazzers (0.93m), Clixsense (2.2m), CSDN (6.4m), Neopets (68.3m), 000webhost (153m). When we analyze the guessing curve $\lambda_{M,D}$ we will assume that each dataset D represents $|D|$ independent samples from an (unknown) probability distribution over user passwords — this unknown

password distribution may be different at different sites. We remark that when analyzing the guessing curve $\lambda_{M,B,D}$ we do not need to make any assumptions about how the dataset D was sampled. However, we argue that the independent samples assumption is reasonable for the datasets we consider. Blocki and Liu [4] developed a linear programming technique which can detect when a dataset is blatantly inconsistent with our assumption about iid samples, e.g., if many user registered for two accounts with the same password or if a large fraction of the dataset was duplicated. Their technique rejects the LinkedIn frequency corpus, which contained far more passwords than unique e-mail addresses. By contrast, all of the datasets that we consider passed the consistency checks from [4].

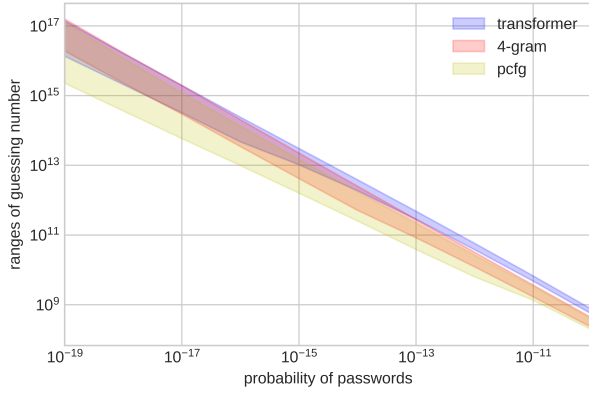
Ethical Considerations The usage of password datasets which contain stolen passwords that were subsequently leaked on the internet raises important ethical considerations. Our usage of the datasets does not pose any additional risk to users as the datasets are already publicly available. We do not crack any new passwords as part of our analysis nor do we attempt to deanonymize the datasets by linking passwords to particular user accounts. Furthermore, we believe that the statistical techniques developed in this paper may benefit users by helping defenders to pick informed password policies.

6.1. Bounding the Guessing Number with Confidence

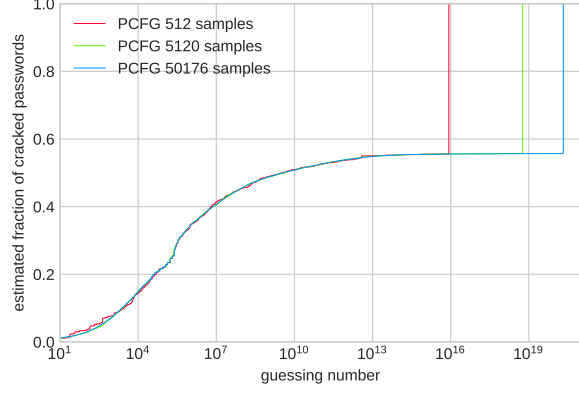
We begin by using Theorems 1 and 2 to upper and lower bound the guessing numbers for specific passwords. When applying Theorem 1 we set the number of samples $k = 206848$ ⁶ and $\epsilon = 0.005$ to obtain confidence $> 99\%$ that each upper/lower bound are correct. Similarly, when applying Theorem 2 we set $n = 186$ and $k = 5120$, $\delta = 0.333$ and $\epsilon = 0.167$ to obtain confidence 99% that each upper bound is correct. Since we have two separate lower bounds on the guessing number we will take the maximum of the lower bounds obtained from Theorems 1 and 2 — union bounds imply that the maximum lower bound will be correct with probability at least 98% . Figure 1a plots our upper/lower bounds for the guessing number as the probability q ranges from 10^{-19} to 10^{-12} . We consider three models: PCFG, Markov and Transformer each trained on the Bfield dataset (due to space limitations we omit similar plots for models trained on other datasets). As we can see, the distance between the upper/lower bounds increases as the password probability decreases. For example, consider the PCFG model, for passwords with probability $q = 10^{-19}$ our upper/lower bound on the guessing number range from 2.3×10^{15} to 1.3×10^{17} . By contrast, when $q = 10^{-11}$ the range is $[1.9 \times 10^8, 3.9 \times 10^8]$. Estimating the guessing numbers for strong passwords is particularly error prone. Thus, it is important to consider the confi-

5. Our code is publicly available in an anonymous Github repository <https://github.com/ConfidentMonteCarlo/ConfidentMonteCarlo.git> and constantly maintained for better modularization

6. Throughout the experiments we set number of samples a multiple of 512 since sample generation using transformer is computed by GPU in parallel with batch size 512



(a) Ranges of Guessing Number vs Password Probability



(b) Regular Monte Carlo Estimation with Varying Sample sizes

Figure 1: Limitation of Regular Monte Carlo Estimation

dence range of guessing numbers when measuring password strength/resistance to brute-force guessing attacks.

Limitation of Regular Monte Carlo Strength Estimation

Figure 1b shows what happens if we apply regular Monte Carlo estimation as a heuristic to estimate the attacker’s guessing curve. We train our PCFG model using the bfield dataset withholding 25,000 passwords D_{test} for testing. For the purpose of comparison we run regular Monte Carlo strength estimation with three different sampling parameters $k \in \{512, 5120, 50176\}$. The figure plots the guessing budget B vs. the estimated fraction of cracked passwords i.e., the fraction of passwords in D_{test} whose estimated guessing number is less than B . As we noted previously our guessing number estimates become less and less certain as q decreases so intuitively we might expect the estimated guessing curves to be less accurate when the guessing budget B is large. Indeed this is what we observe. Each of the PCFG guessing curves suddenly spikes to 100%, and these sudden spikes occur at different points for different values of our sampling parameter $k \in \{512, 5120, 50176\}$ since sample size affects the strongest password being sampled⁷. This motivates the need to derive upper/lower bounds on the attacker’s guessing curve which hold with high probability.

6.2. Confident Guessing Curves

We now turn our attention to the problem of upper/lower bounding the attacker’s guessing curve using Theorem 6, 7,

7. This is due to the definition of regular Monte Carlo. If the password y is strictly stronger than any password in the set S of $|S| = k$ passwords sampled from our model M then Monte Carlo Strength Estimation will output $\hat{G}(y) := \frac{1}{k} \sum_{x: p_x^M \geq p_y^M} \frac{1}{p_x^M}$ as the estimated guessing number for y — this observation holds even if the true guessing number for password y should be $G(y) = \infty$ i.e., if y cannot be generated by the model M . Thus, for $B \geq \frac{1}{k} \sum_{x \in S} \frac{1}{p_x^M}$ we will have $\hat{G}(y) \leq B$ for all passwords y sampled from our unknown password distribution. Thus, the “spike” in our plot occurs when the guessing budget is $B = \frac{1}{k} \sum_{x \in S} \frac{1}{p_x^M}$. We expect that $\frac{1}{k} \sum_{x \in S} \frac{1}{p_x^M}$ is grows with the parameter k .

8, 12 13 and 14.

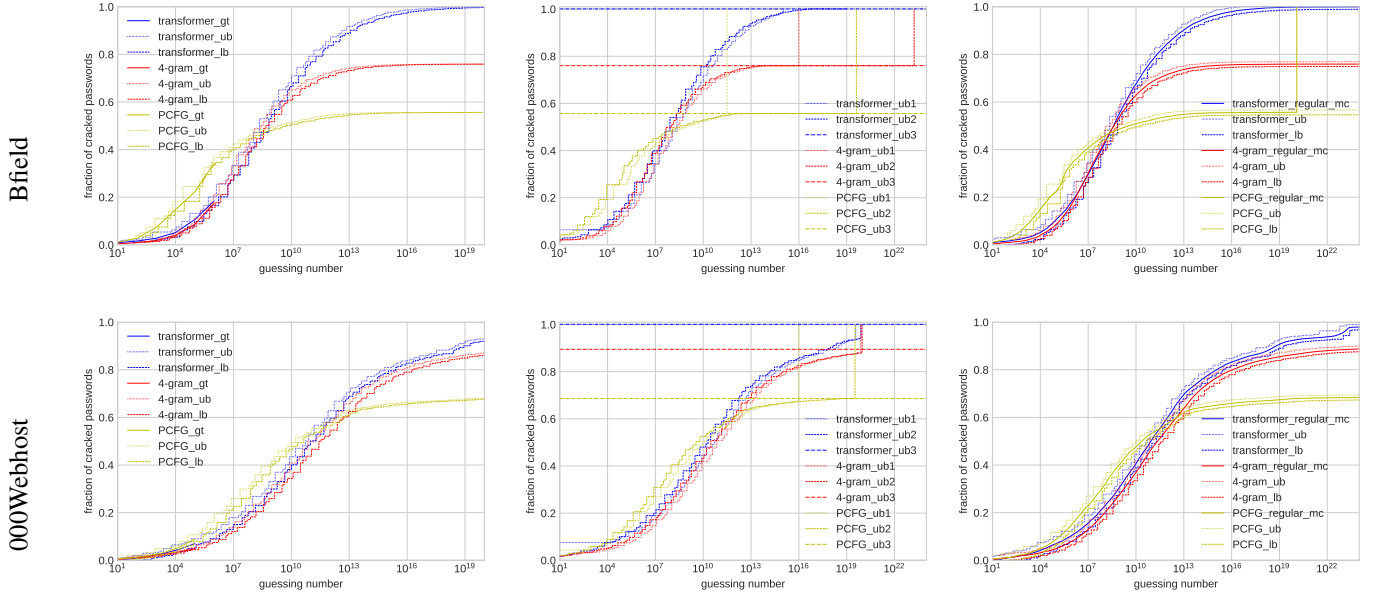
Experimental Setup: For each password dataset $D_{original}$ we first perform train-test split to obtain D_{train} and D_{test} with $|D_{test}| = 25,000$. All 3 probabilistic models are trained with D_{train} , then we use D_{test} to upper/lower bound the attacker’s guessing curve.

To apply Theorem 6 and 7 we need to define a set of probability mesh points $Q = \{q_1, \dots, q_\ell\}$. In particular, we fix probability mesh points to be $j \cdot 10^{-4-i}$ for $i \in [1, 25]$ and $j \in \{0.25, 0.5, 0.75, 1\}$ for a total of $\ell = 25 \cdot 4 = 100$ mesh points. In Theorem 6 and 7 there are two sources of confidence loss. The term $\ell \cdot \exp(2k\epsilon_1^2)$ (resp. $\ell \cdot \exp(-2n\epsilon_1^2)$) in Theorem 6 (resp. 7) upper bounds the guessing number error associated with *any* point in our probability mesh. The term $\exp(-2|D_{test}|\epsilon_2^2)$ accounts for confidence loss due to sampling error from our unknown password distribution.

The parameter setting of Theorem 12 13 and 14 is identical with that of Theorem 6 7 8, respectively.

In our experiments we instantiate the parameters k , ϵ_1, ϵ_2 , n and $|D_{test}|$ to ensure that the total probability of failure for each bound is at most 0.01. More specifically, we set the number of samples $k = 206848$, $|D_{test}| = 25,000$ and $\epsilon_2 = 0.01$ in both Theorem 6 and 7. We set $\epsilon_1 = 0.005$ in Theorem 6, in Theorem 7 we set $\epsilon_1 = 0.167$, $n = 186$ and $\delta = 0.333$. Similarly, in Theorem 8 we set $\epsilon = 0.0096$ and $|D_{test}| = 25,000$ to ensure that, except with probability 0.01, $\hat{\lambda}_{M,D}^{ub3} + \epsilon$ is an upper bound on $\lambda_{M,B}$ for all guessing budgets $B \geq 0$.

Because we have multiple techniques to generate upper bounds it will often make sense to consider the best upper/lower bound. Thus, we define $\hat{\lambda}_{M,B}^{ub} = \min \left\{ \hat{\lambda}_{M,B,D,S,\epsilon_1}^{ub1} + \epsilon_2, \hat{\lambda}_{M,B,D,S,\delta}^{ub2} + \epsilon_2, \hat{\lambda}_{M,D}^{ub3} + \epsilon, 1 \right\}$ and $\hat{\lambda}_{M,B,D}^{ub} = \min \left\{ \hat{\lambda}_{M,B,D,S,\epsilon_1}^{ub1}, \hat{\lambda}_{M,B,D,S,\delta}^{ub2}, \hat{\lambda}_{M,D}^{ub3}, 1 \right\}$ which are best upper bounds for $\lambda_{M,B}$ and $\lambda_{M,B,D}$, respectively. Applying union bounds the probability that the curve $\hat{\lambda}_{M,B}^{ub}$ (resp. $\hat{\lambda}_{M,B,D}^{ub}$) is not a valid upper bound for $\lambda_{M,B}$ (resp. $\lambda_{M,B,D}$) is at most 0.03. For notional



(a) Upper/Lower Bounds on $\lambda_{M,B,D}$

(b) 3 Upper Bounds on $\lambda_{M,B}$

(c) Upper/Lower Bounds on $\lambda_{M,B}$

Figure 2: Confident Bounds

convenience we also define $\hat{\lambda}_{M,B}^{lb} = \hat{\lambda}_{M,B,D,S,\epsilon_1}^{lb1} - \epsilon_2$ and $\hat{\lambda}_{M,B,D}^{lb} = \hat{\lambda}_{M,B,D,S,\epsilon_1}^{lb1}$. In figure legends we use modelname_ub to denote $\hat{\lambda}_{M,B}^{ub}$ or $\hat{\lambda}_{M,B,D}^{ub}$ contingent on the figure caption. The same case applies to lower bound legend modelname_lb.

Upper/Lower Bounds on $\lambda_{M,B,D}$. Figure 4a plots our best upper bound $\hat{\lambda}_{M,B,D}^{ub}$ (resp. best lower bound $\hat{\lambda}_{M,B,D}^{lb}$) as the guessing budget B varies for each password model M . Due to space limitations we only show the plots for the Bfield and 000webhost datasets in the main body. The remaining plots can be found in the Appendix C — see Figure 4. We additionally plot the ground truth⁸ $\lambda_{M,B,D}$ for $B \leq 10^6$ — denoted by modelname_gt in figure legends.

Discussion. Consider Figure 4a (Bfield) as an example. We find the the upper and lower bounds are reasonably close to each other. Furthermore, when $B \leq 10^6$ we note that $\lambda_{M,B,D}$ (the ground truth of fraction of cracked passwords against dataset D) is sandwiched between our upper/lower bounds for all three models M . We can also use our confident guessing curves to draw rigorous statistical comparisons between the three cracking models. In general, if we find that $\hat{\lambda}_{M_1,B,D}^{lb} > \hat{\lambda}_{M_2,B,D}^{ub}$ then this supports the hypothesis that model M_1 outperforms model M_2 with guessing budget B against dataset D . Notice that whenever $B \geq 10^{11}$ the lower bound for Transformer is strictly higher than the upper bound of 4-gram Markov model. This supports the hypothesis that Transformers outperforms 4-gram for larger guessing budgets.

8. We generated a dictionary of the top 1 million popular passwords in each model-defined distribution \mathcal{M} by brute-force and use the dictionary to crack passwords in D_{test} .

Comparing Our Upper Bounds on $\lambda_{M,B}$. For the sake of comparison Figure 2b plots our three upper bounds on $\lambda_{M,B}$ separately. In the legend modelname_ub1 denotes the upper bound $\min\{1, \hat{\lambda}_{M,B,D,S,\epsilon_1}^{ub1} + \epsilon_2\}$, modelname_ub2 denotes the upper bound $\min\{1, \hat{\lambda}_{M,B,D,S,\delta}^{ub2} + \epsilon_2\}$ and modelname_ub3 denotes the upper bound $\min\{1, \hat{\lambda}_{M,D}^{ub3} + \epsilon\}$. In Figure 2b we occasionally observe spiking behavior for the first two upperbounds (ub1 and ub2). We note that for an upper bound this behavior is not problematic for two reasons. First, even if an upper bound spikes to 1.0 the upper bound continues to be accurate i.e., the interpretation is simply that the current statistical approach does not rule out the *possibility* that the attackers cracks 100% of passwords. By contrast, for regular Monte Carlo strength estimation if the guessing curve spikes to 1 this represents a (possibly incorrect) *prediction* that the attacker will crack 100% of passwords. Second, in all of the plots from Figure 2b we find that the first two upper bounds ub1 and ub2 approach the third upper bound ub3 (a straight line) *before* we observe the spiking behavior. Thus, we expect to obtain reasonably tight upper bounds by considering best of the three bounds i.e., $\hat{\lambda}_{M,B}^{ub}$.

Confident Bounds on $\lambda_{M,B}$. Figure 2c compares our best upper/lower bounds $\hat{\lambda}_{M,B}^{ub}$ and $\hat{\lambda}_{M,B}^{lb}$ with the guessing curve obtained from regular Monte Carlo Strength estimation — denoted by $\lambda_{M,B}^{MC}$ ⁹. The number of samples used to compute

9. Note that the regular Monte Carlo estimate $\lambda_{M,B}^{MC}$ will also depend on the test dataset D (sampled from the unknown password distribution) and samples S (sampled from the model M) in addition to the model M and guessing budget B . We omit S and D from the subscript to simplify notation.

$\lambda_{M,B}^{MC}$ is 10240, a multiple of 512 that is closest to sample size 10000 which is recommended in [8].

Discussion. If $\lambda_{M,B}^{MC} < \hat{\lambda}_{M,B}^{lb}$ (resp. $\lambda_{M,B}^{MC} > \hat{\lambda}_{M,B}^{ub}$) then we can confidently conclude that regular Monte Carlo Strength Estimation is underestimating (resp. overestimating) the true guessing curve. On the other hand if $\hat{\lambda}_{M,B}^{lb} \leq \lambda_{M,B}^{MC} \leq \hat{\lambda}_{M,B}^{ub}$ then the estimate $\lambda_{M,B}^{MC}$ is *plausibly accurate*. Furthermore, if we have $\hat{\lambda}_{M,B}^{lb} \leq \lambda_{M,B}^{MC} \leq \hat{\lambda}_{M,B}^{ub}$ and the difference $\hat{\lambda}_{M,B}^{ub} - \hat{\lambda}_{M,B}^{lb}$ is sufficiently small (e.g., < 0.05) then we can confidently conclude that $\lambda_{M,B}^{MC}$ is an accurate estimation.

As an example, consider Figure 2c (Bfield) using PCFG as a our probabilistic password model. When $B \approx 1.2 \times 10^{20}$, $\lambda_{M,B}^{MC}$ suddenly jumps to 1 whereas $\hat{\lambda}_{M,B}^{ub} = 0.56$. Thus, we can confidently conclude that regular Monte Carlo Strength estimation significantly overestimates the fraction of passwords cracked by PCFG when $B \geq 1.2 \times 10^{20}$. For $B < 1.2 \times 10^{20}$ we have $\hat{\lambda}_{M,B}^{lb} \leq \lambda_{M,B}^{MC} \leq \hat{\lambda}_{M,B}^{ub}$. Moreover, in the range of $[2.8 \times 10^4, 1.7 \times 10^5]$ we have $\hat{\lambda}_{M,B}^{ub} - \hat{\lambda}_{M,B}^{lb} > 0.05$ so the estimate $\lambda_{M,B}^{MC}$ is *plausibly accurate*. For $B < 2.8 \times 10^4$ or $1.7 \times 10^5 \leq B \leq 1.2 \times 10^{20}$ we have $\hat{\lambda}_{M,B}^{ub} - \hat{\lambda}_{M,B}^{lb} < 0.05$ allowing us to confidently conclude that $\lambda_{M,B}^{MC}$ is an accurate estimate for $\lambda_{M,B}$.

We can also apply our results to compare the distributions from different datasets. For example, we compare Bfield and 000webhost by fixing the password probabilistic model M to be transformer and guessing budget to be $B = 10^{12}$, then we consider bounds for fraction of cracked passwords. We have $\hat{\lambda}_{M,B}^{lb} = 0.84$ for the Bfield distribution and $\hat{\lambda}_{M,B}^{ub} = 0.63$ for the 000webhost distribution. Thus, we can confidently conclude that the 000webhost distribution is more resistant to attacks by an attacker using the Transformer Cracking Model with guessing budget $B = 10^{12}$.

6.3. Password Composition Policies

In this subsection we apply our statistical techniques from Section 5.4 to examine the impact of password composition policies (PCPs) on the attacker’s guessing curve. Specifically, we consider 3 password composition policies: 1) *1class8*, passwords with length no less than 8; 2) *1class16*, passwords with length no less than 16; 3) *3class8*, passwords with length no less than 8 and containing at least 3 categories of characters out of 4 (lower case, upper case, digits, special characters). Applying a composition policy \mathbb{C} will alter the password distribution \mathcal{P} over user passwords. Of course the hope is that passwords from the new distribution $\mathcal{P}^{\mathbb{C}}$ will be harder for an attacker to crack. **Normalized Probability Model.** One immediate challenge is that we need to obtain a test dataset sampled from the new (unknown) distribution $\mathcal{P}^{\mathbb{C}}$ (after applying our restriction \mathbb{C}) before we can apply our statistical techniques. However, each of our dataset(s) D was sampled from the original distribution \mathcal{P} . The normalized probability model [3] is a heuristic assumption about the way a restriction \mathbb{C} impacts the password distribution i.e., the probability of sampling the password x under $\mathcal{P}^{\mathbb{C}}$ is assumed to be the conditional

probability of selecting x conditioning on the event that the sampled password y is allowed by our composition policy. More formally, for any disallowed password $x \notin \mathbb{C}$ we have $\Pr_{y \sim \mathcal{P}^{\mathbb{C}}}[y = x] = 0$ and for any allowed password we have $\Pr_{y \sim \mathcal{P}^{\mathbb{C}}}[y = x] = \Pr_{y \sim \mathcal{P}}[y = x : y \in \mathbb{C}]$. Under the normalized probabilities assumption if we have a dataset D sampled iid from \mathcal{P} then the dataset $D^{\mathbb{C}}$, obtained by filtering D to remove any passwords that are not consistent with \mathbb{C} , can be viewed as iid samples from $\mathcal{P}^{\mathbb{C}}$. In this subsection we will make the heuristic normalized probabilities assumption so that we can obtain iid samples from $\mathcal{P}^{\mathbb{C}}$ for our analysis.

Filtering Training Data, Filter Guesses or Both We consider a scenario where the attacker obtains a (leaked) password dataset A which can be used to train a password model M . The attacker then runs an offline attack against passwords sampled from $\mathcal{P}^{\mathbb{C}}$ i.e., following the composition policy \mathbb{C} . Before training the model M the attacker might filter the dataset A to remove passwords that do not follow our policy \mathbb{C} . Alternatively, the attacker might choose to train M on the entire dataset (e.g., to maximize the number of training examples) since we can always filter out guesses M generates when they are inconsistent with \mathbb{C} . Thus, we obtain four possible strategies for the attacker: *generic-unfiltered*, *generic-filtered*, *specific-unfiltered* and *specific-filtered*. Here *generic* (resp. *specific*) indicates that we train the model on the full dataset M (resp. the subset of A consistent with \mathbb{C}) and *filtered* (resp. *unfiltered*) indicates that we filter (resp. do not filter) the password guessed output by our model M to remove guesses that are inconsistent with \mathbb{C} .

In Figure 3a we adopt the *generic-filtered* approach and train our transformer model M trained on the entire Clixsense dataset. We upper and lower bound the attacker’s guessing curve against three groups of CSDN passwords. Observe that for $B \leq 10^{18}$ the upper bound for 1class16 passwords lies below the lower bound for 3class8 passwords and the upper bound for 3class8 passwords lies below the lower bound for 1class8 passwords. Thus, under our heuristic normalized probabilities assumption, we can conclude that the 1class16 PCP (resp. 3class8 PCP) provides the strongest (resp. second strongest) resistance to an offline attacker using our transformer model M . In Figure 3b we fix the PCP to be 1class8 and the password model to be transformer. We compare the efficiency of the four different attacker strategies. In this example, we can see that the attacker benefits by only training on Clixsense passwords that are consistent with the 1class8 policy. For example, in order to crack 40% the 1class8 passwords, a model trained only on Clixsense passwords that are consistent with the 1class8 policy requires 4×10^7 guesses *at most*. By contrast, a generic model trained on the entire Clixsense dataset requires *at least* 10^9 to achieve the same success rate. However, the benefit of filtering out inconsistent password guesses generated by our model is insignificant. This is because almost all of the passwords generated by our model are already consistent with the 1class8 PCP — whether or not the training data was generic or specific to the 1class8

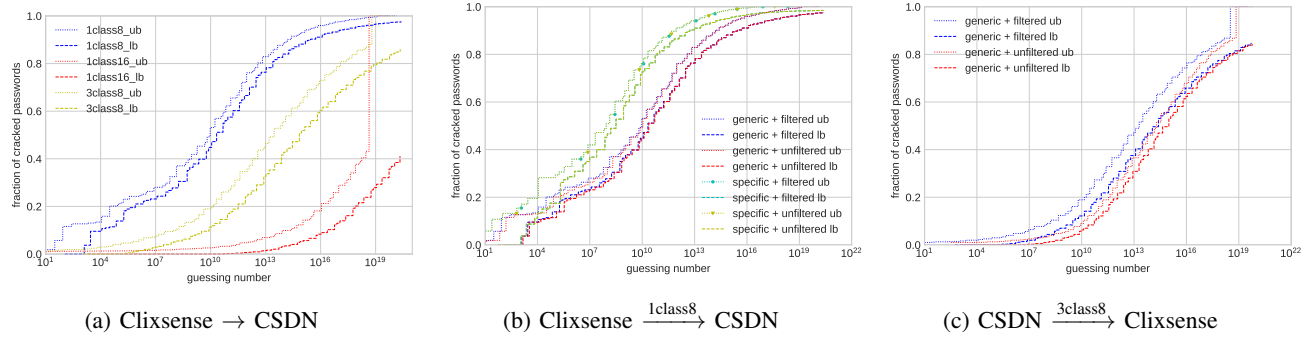


Figure 3: Confident Guessing Curves When Considering Password Composition Policies

PCP.

In Figure 3c we fix a strong password policy 3class8. We consider a generic transformer model trained on the entire CSDN dataset¹⁰ and attack 3class8 passwords from the Clixsense dataset. We compare the performance of the model with and without filtering. We observe that with a filtered dictionary the attacker is able to crack more passwords. We can be confident in the prior conclusion as long as $5 \cdot 10^7 < B < 10^{15}$ because the upper bound for our unfiltered curve lies below the lower bound for our filtered curve.

7. Conclusion

In this paper, we provided theoretical and empirical evidence that regular Monte Carlo will sometimes yield inaccurate estimations of the guessing number of a password and of the attacker’s guessing curve. We extend the regular Monte Carlo method by developing rigorous statistical techniques to confidently upper/lower bound the guessing number of a password. We also showed how to use our Confident Monte Carlo framework to provide high confidence upper/lower bounds on the attacker’s guessing curve. Our rigorous statistical framework allows us to evaluating the impact of a password policy interventions (e.g., password composition policies) on password strength, rigorously compare the performance of different cracking models and characterize the resistance of a password dataset/distribution to an offline password attacker.

References

- [1] Hashcat: advanced password recovery. <https://hashcat.net/hashcat/>.
- [2] Wenjie Bai and Jeremiah Blocki. Dahash: distribution aware tuning of password hashing costs. In *International Conference on Financial Cryptography and Data Security*, pages 382–405. Springer, 2021.
- [3] Jeremiah Blocki, Saranga Komanduri, Ariel Procaccia, and Or Sheffet. Optimizing password composition policies. In *Proceedings of the fourteenth ACM conference on Electronic commerce*, pages 105–122. ACM, 2013.
- [4] Jeremiah Blocki and Peiyuan Liu. Towards a rigorous statistical analysis of empirical password datasets. In *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2023.
- [5] Claude Castelluccia, Abdelberi Chaabane, Markus Dürmuth, and Daniele Perito. When privacy meets security: Leveraging personal information for password cracking. *arXiv preprint arXiv:1304.6584*, 2013.
- [6] Claude Castelluccia, Markus Dürmuth, and Daniele Perito. Adaptive password-strength meters from Markov models. In *NDSS 2012*. The Internet Society, February 2012.
- [7] Matteo Dell’Amico. Implementation of monte carlo estimation, 2017. <https://github.com/matteodellamico/montecarlopwd>.
- [8] Matteo Dell’Amico and Maurizio Filippone. Monte Carlo strength evaluation: Fast and reliable password checking. In Indrajit Ray, Ninghui Li, and Christopher Kruegel, editors, *ACM CCS 2015*, pages 158–169. ACM Press, October 2015.
- [9] Solar Designer. John the ripper password cracker, 2006.
- [10] Patrick Gage Kelley, Saranga Komanduri, Michelle L. Mazurek, Richard Shay, Timothy Vidas, Lujio Bauer, Nicolas Christin, Lorrie Faith Cranor, and Julio Lopez. Guess again (and again and again): Measuring password strength by simulating password-cracking algorithms. In *2012 IEEE Symposium on Security and Privacy*, pages 523–537. IEEE Computer Society Press, May 2012.
- [11] Enze Liu, Amanda Nakanishi, Maximilian Golla, David Cash, and Blase Ur. Reasoning analytically about password-cracking software. In *2019 IEEE Symposium on Security and Privacy*, pages 380–397. IEEE Computer Society Press, May 2019.
- [12] Jerry Ma, Weining Yang, Min Luo, and Ninghui Li. A study of probabilistic password models. In *2014 IEEE Symposium on Security and Privacy*, pages 689–704. IEEE Computer Society Press, May 2014.
- [13] Colin McDiarmid et al. On the method of bounded differences. *Surveys in combinatorics*, 141(1):148–188, 1989.
- [14] William Melicher, Blase Ur, Sean M. Segreti, Saranga Komanduri, Lujio Bauer, Nicolas Christin, and Lorrie Faith Cranor. Fast, lean, and accurate: Modeling password guessability using neural networks. In Thorsten Holz and Stefan Savage, editors, *USENIX Security 2016*, pages 175–191. USENIX Association, August 2016.
- [15] Robert Morris and Ken Thompson. Password security: A case history. *Communications of the ACM*, 22(11):594–597, 1979.

10. If we remove passwords from CSDN that are inconsistent with the 3class8 policy then the remaining training set is quite a bit smaller. Thus, we found that the specific-filtered strategy does not outperform the generic-filtered strategy.

- [16] Dario Pasquini, Marco Cianfriglia, Giuseppe Ateniese, and Massimo Bernaschi. Reducing bias in modeling real-world password strength via deep learning and dynamic dictionaries. In Michael Bailey and Rachel Greenstadt, editors, *USENIX Security 2021*, pages 821–838. USENIX Association, August 2021.
- [17] Richard Shay, Saranga Komanduri, Adam L. Durity, Phillip (Seoung) Huh, Michelle L. Mazurek, Sean M. Segreti, Blase Ur, Lujo Bauer, Nicolas Christin, and Lorrie Faith Cranor. Designing password policies for strength and usability. *ACM Trans. Inf. Syst. Secur.*, 18(4), may 2016.
- [18] Blase Ur, Felicia Alfieri, Maung Aung, Lujo Bauer, Nicolas Christin, Jessica Colnago, Lorrie Faith Cranor, Henry Dixon, Pardis Emami Naeini, Hana Habib, et al. Design and evaluation of a data-driven password meter. In *Proceedings of the 2017 chi conference on human factors in computing systems*, pages 3775–3786, 2017.
- [19] Blase Ur, Sean M. Segreti, Lujo Bauer, Nicolas Christin, Lorrie Faith Cranor, Saranga Komanduri, Darya Kurilova, Michelle L. Mazurek, William Melicher, and Richard Shay. Measuring real-world accuracies and biases in modeling password guessability. In Jaeyeon Jung and Thorsten Holz, editors, *USENIX Security 2015*, pages 463–481. USENIX Association, August 2015.
- [20] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [21] Rafael Veras, Christopher Collins, and Julie Thorpe. On semantic patterns of passwords and their security impact. In *NDSS 2014*. The Internet Society, February 2014.
- [22] Ding Wang, Ping Wang, Debiao He, and Yuan Tian. Birthday, name and bifacial-security: Understanding passwords of chinese web users. In Nadia Heninger and Patrick Traynor, editors, *USENIX Security 2019*, pages 1537–1555. USENIX Association, August 2019.
- [23] Matt Weir, Sudhir Aggarwal, Breno de Medeiros, and Bill Glodek. Password cracking using probabilistic context-free grammars. In *2009 IEEE Symposium on Security and Privacy*, pages 391–405. IEEE Computer Society Press, May 2009.

Appendix A. Experiment Details

A.1. Transformer neural network

Our Transformer neural network is composed of a classic encoder-decoder structure. The encoder contains 16 identical layers stacked upon each other. Each layer has 2 sub-layers. The first is a multi-head self-attention mechanism with 16 heads, and the second is a simple, position-wise fully connected feed-forward network. The decoder is a basic linear layer. In addition, we set the embedding size to be 128 and the size of hidden layers to be 1024.

A.2. Data Preprocessing

For Transformer neural network, we restrict the alphabet set to be the union of 96 printable ASCII characters and 2 special characters — start/ \perp denoting the start/end of a password text string, respectively. Also, we fix the maximum length of passwords to be 16. Thus, passwords in D_{train} composed of non-ASCII characters or having length larger than 16 are filtered out, and the remaining passwords are prepend/append with start/ \perp . Passwords in D_{test} are untouched in evaluation, passwords composed of characters not

in the alphabet set are assigned probability 0, which implies not being cracked under any circumstances. For 4-gram Markov model, we adopted implementation from [7] which only performs prepending/appending without any smoothing techniques. We did not apply any smoothing techniques to make sure the probabilities of all allowable strings add up to 1, i.e., the model strictly defines a probability distribution. How heuristics like smoothing techniques affect rigorous analysis of password probabilistic models remains an open question.

Appendix B. Bounding Guessing Curve Under Password Composition Policies

For any password pwd we use $\text{Allowed}^{\mathbb{C}}(pwd)$ to describe the password policy \mathbb{C} an attacker would follow when making guesses. We set $\text{Allowed}^{\mathbb{C}}(pwd) = 1$ if and only if password pwd satisfies the policy \mathbb{C} (e.g. passwords must be at least 8 characters long). Similar to $G^{\text{EX}}(\cdot)$ and $G^{\text{IN}}(\cdot)$, we define $G^{\mathbb{C},\text{EX}}(q) = |\{z \in \mathcal{M} : p_z^M > q \wedge \text{Allowed}^{\mathbb{C}}(z) = 1\}|$ and $G^{\mathbb{C},\text{IN}}(q) = |\{z \in \mathcal{M} : p_z^M \geq q \wedge \text{Allowed}^{\mathbb{C}}(z) = 1\}|$ such that $G^{\mathbb{C},\text{EX}}(q) + 1$ and $G^{\mathbb{C},\text{IN}}(q)$ are the smallest and largest possible guessing number of a password with probability q in M under policy \mathbb{C} , i.e., $G^{\mathbb{C},\text{EX}}(q) + 1 \leq G^{\mathbb{C}}(q) \leq G^{\mathbb{C},\text{IN}}(q)$. Given a set S of k iid samples randomly selected from model distribution \mathcal{M} , we let $\hat{G}_S^{\mathbb{C},\text{EX}}(q) = \frac{1}{k} \sum_{z \in S, p_z^M > q, \text{Allowed}^{\mathbb{C}}(z)=1} \frac{1}{p_z^M}$ be the regular Monte Carlo estimate of $G^{\mathbb{C},\text{EX}}(q)$, and $\hat{G}_S^{\mathbb{C},\text{IN}}(q) = \frac{1}{k} \sum_{z \in S, p_z^M \geq q, \text{Allowed}^{\mathbb{C}}(z)=1} \frac{1}{p_z^M}$ be the regular Monte Carlo estimate of $G^{\mathbb{C},\text{IN}}(q)$. For $\phi = \{\text{EX}, \text{IN}\}$, we denote $\hat{G}_{S,med}^{\mathbb{C},\phi}(q) = \text{median}(\{\hat{G}_{S_i}^{\mathbb{C},\phi}(q)\}_{1 \leq i \leq n})$ where $\mathbb{S} = \{S_1, S_2, \dots, S_n\}$ and each S_i contains k independent samples from model M .

In this section, we state the theorems that upper and lower bound the guessing curve of an attacker who follows given password policies \mathbb{C} . We define $\lambda_{M,B,D}^{\mathbb{C}} = \frac{1}{|D|} |y \in D : G^{\mathbb{C}}(y) \geq B|$ to be the guessing curve of set D against an attacker with model M under password policy \mathbb{C} , i.e., the percentage of passwords in D cracked by making the top B most probable guesses outputted by model M satisfying password policy \mathbb{C} . Define $\lambda_{M,B}^{\mathbb{C}} = \Pr_{y \leftarrow \mathcal{P}}[G^{\mathbb{C}}(y) \leq B]$ to be the probability of a password from some unknown distribution \mathcal{P} cracked by an attacker making the top B guesses of model M satisfying password policy \mathbb{C} .

First Upper/Lower Bound. We define $\hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},ub1}$ and $\hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},lb1}$ below to be the first upper and lower bounds on $\lambda_{M,B,D}^{\mathbb{C}}$:

$$\begin{aligned} \hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},ub1} &:= \min_{1 \leq i \leq \ell, B \leq \hat{G}_S^{\mathbb{C},\text{IN}}(q_i) - \epsilon/q_i} \left(\lambda_{M,G^{\mathbb{C}}(q_i),D}^{\mathbb{C}} \right), \\ \hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},lb1} &:= \max_{1 \leq i \leq \ell, B \geq \hat{G}_S^{\mathbb{C},\text{EX}}(q_i) + \epsilon/q_i} \left(\lambda_{M,G^{\mathbb{C}}(q_i),D}^{\mathbb{C}} \right). \end{aligned}$$

For completeness, we set $\hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},ub1} = 1$ (resp. $\hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},lb1} = 0$) if $B > \max_{1 \leq i \leq \ell} \{\hat{G}_S^{\mathbb{C},\text{IN}}(q_i) - \epsilon/q_i\}$ (resp.

$B < \min_{1 \leq i \leq \ell} \{\hat{G}_S^{\text{C,EX}}(q_i) + \epsilon/q_i\}$. Then our first upper and lower bounds on $\lambda_{M,B}^{\text{C}}$ is shown below:

Theorem 9. *Given a password distribution \mathcal{P} , a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$ and a password policy \mathbb{C} , for any guessing number $B \geq 0$ and any parameters $0 \leq \epsilon_1, \epsilon_2 \leq 1$, we have:*

$$\begin{aligned} \Pr \left[\lambda_{M,B}^{\text{C}} \leq \hat{\lambda}_{M,B,D,S,\epsilon_1}^{\text{C,ub1}} + \epsilon_2 \right] &\geq 1 - \alpha \\ \Pr \left[\lambda_{M,B}^{\text{C}} \geq \hat{\lambda}_{M,B,D,S,\epsilon_1}^{\text{C,lb1}} - \epsilon_2 \right] &\geq 1 - \alpha \end{aligned}$$

where $\alpha = \ell \cdot \exp(-2k\epsilon_1^2) - \exp(-2|D|\epsilon_2^2)$ and the randomness is taken over the sample set $S \leftarrow \mathcal{M}^k$ with size k and the dataset $D \leftarrow \mathcal{P}^{|D|}$.

Second Upper Bound. We define $\hat{\lambda}_{M,B,D,S,\delta}^{\text{C,ub2}}$ below to be our second bound on $\lambda_{M,B,D}^{\text{C}}$:

$$\hat{\lambda}_{M,B,D,S,\delta}^{\text{C,ub2}} := \min_{1 \leq i \leq \ell, B \leq \delta \cdot \hat{G}_{S,\text{med}}^{\text{C,IN}}(q)} \left(\lambda_{M,G^{\text{C}}(q_i),D}^{\text{C}} \right)$$

For completeness, we set $\hat{\lambda}_{M,B,D,S,\epsilon}^{\text{C,ub2}} = 1$ if $B > \max_{1 \leq i \leq \ell} \{\delta \cdot \hat{G}_{S,\text{med}}^{\text{C,IN}}(q)\}$. Then our second upper bound on $\lambda_{M,B}^{\text{C}}$ is shown below:

Theorem 10. *Given a password distribution \mathcal{P} , a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$ and a password policy \mathbb{C} , for any guessing number $B \geq 0$ and any parameters $0 < \delta \leq \frac{1}{2}, 0 \leq \epsilon_1 \leq \frac{1}{2} - \delta, 0 \leq \epsilon_2 \leq 1$, we have:*

$$\Pr \left[\lambda_{M,B}^{\text{C}} \leq \hat{\lambda}_{M,B,D,S,\delta}^{\text{C,ub2}} + \epsilon_2 \right] \geq 1 - \alpha$$

where $\alpha = \ell \cdot \exp(-2n\epsilon_1^2) - \exp(-2|D|\epsilon_2^2)$ the randomness is taken over n Monte Carlo sample sets $\mathbb{S} = \{S_1, S_2, \dots, S_n\}$ each of which contains k samples from model M and the dataset $D \leftarrow \mathcal{P}^{|D|}$.

Third Upper Bound. We denote $\hat{\lambda}_{M,D}^{\text{C,ub3}} := \frac{1}{|D|} |y \in D : p_y^M > 0 \wedge \text{Allowed}^{\text{C}}(y) = 1|$ to be the percentage of passwords that will be eventually guessed by model M restricted by policy \mathbb{C} . We have $\lambda_{M,B,D}^{\text{C}} \leq \hat{\lambda}_{M,D}^{\text{C,ub3}}$ for any guessing number $B > 0$ due to the fact that passwords with zero probability in M will never be guessed by M . Then we have the following trivial upper bound on $\lambda_{M,B}^{\text{C}}$:

Theorem 11. *Given a password distribution \mathcal{P} , a password cracking model M and a password policy \mathbb{C} , for any parameters $0 \leq \epsilon \leq 1$, we have:*

$$\Pr \left[\forall B \geq 0, \lambda_{M,B}^{\text{C}} \leq \hat{\lambda}_{M,D}^{\text{C,ub3}} + \epsilon \right] \geq 1 - \exp(-2|D|\epsilon^2)$$

where the randomness is taken over the dataset $D \leftarrow \mathcal{P}^{|D|}$.

Appendix C. Additional Plots

Figure 4 shows additional plots for experiments described in Section 6.

Appendix D.

Complete Proofs of Bounds on Guessing Number

D.1. Complete Proof of Theorem 1

Recall that the actual guessing number $G(q)$ is in between $G^{\text{EX}}(q)$ and $G^{\text{IN}}(q)$ as explained in Section 3, where $G^{\text{EX}}(q)$ (resp. $G^{\text{IN}}(q)$) denotes the guessing number that a password with probability q is the first (resp. last) guess made by an attacker using model M among all equally likely password guesses. To bound $G(q)$, we start by first bounding $G^{\text{EX}}(q)$ and $G^{\text{IN}}(q)$ using their regular Monte Carlo estimates $\hat{G}_S^{\text{EX}}(q)$ and $\hat{G}_S^{\text{IN}}(q)$ using Hoeffding's inequality in the following lemma:

Lemma 1. *Given a set S with k iid password samples sampled from distribution \mathcal{M} , for any password probability q and any parameter $\epsilon \geq 0$, we have:*

$$\begin{aligned} \Pr[G^\phi(q) \leq \hat{G}_S^\phi(q) + \frac{\epsilon}{q}] &\geq 1 - \exp(-2k\epsilon^2) \\ \Pr[G^\phi(q) \geq \hat{G}_S^\phi(q) - \frac{\epsilon}{q}] &\geq 1 - \exp(-2k\epsilon^2) \end{aligned}$$

where the randomness is taken over the sample set $S \leftarrow \mathcal{M}^k$, and $\phi = \text{EX}$ or IN .

Proof. Given a password probability q and a set S with k password samples randomly sampled from distribution \mathcal{M} , we consider k independent random variables $X_1^\phi, \dots, X_k^\phi$ where for any $i = 1, \dots, k$ and $\phi = \text{EX}, \text{IN}$ we define:

$$\begin{aligned} X_i^{\text{EX}} &:= \begin{cases} \frac{1}{p_z^M} & \text{if the } i\text{th sampled password is } z \text{ and } p_z^M > q; \\ 0 & \text{otherwise.} \end{cases} \\ X_i^{\text{IN}} &:= \begin{cases} \frac{1}{p_z^M} & \text{if the } i\text{th sampled password is } z \text{ and } p_z^M \geq q; \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

Then we have $0 \leq X_i^\phi \leq \frac{1}{q}$ and the expectation of X_i^ϕ is $G^\phi(q)$ as shown below:

$$\begin{aligned} \mathbb{E}(X_i^{\text{EX}}) &= \sum_{z \in \mathcal{M}, p_z^M > q} p_z^M \cdot \frac{1}{p_z^M} = |\{z \in \mathcal{M} : p_z^M > q\}| = G^{\text{EX}}(q) \\ \mathbb{E}(X_i^{\text{IN}}) &= \sum_{z \in \mathcal{M}, p_z^M \geq q} p_z^M \cdot \frac{1}{p_z^M} = |\{z \in \mathcal{M} : p_z^M \geq q\}| = G^{\text{IN}}(q) \end{aligned}$$

Observe that $\hat{G}_S^\phi(q) = \frac{1}{k} \sum_{i=1}^k X_i^\phi$ and $\mathbb{E}(\frac{1}{k} \sum_{i=1}^k X_i^\phi) = \frac{1}{k} \sum_{i=1}^k \mathbb{E}(X_i^\phi) = G^\phi(q)$. Using Hoeffding's inequality we

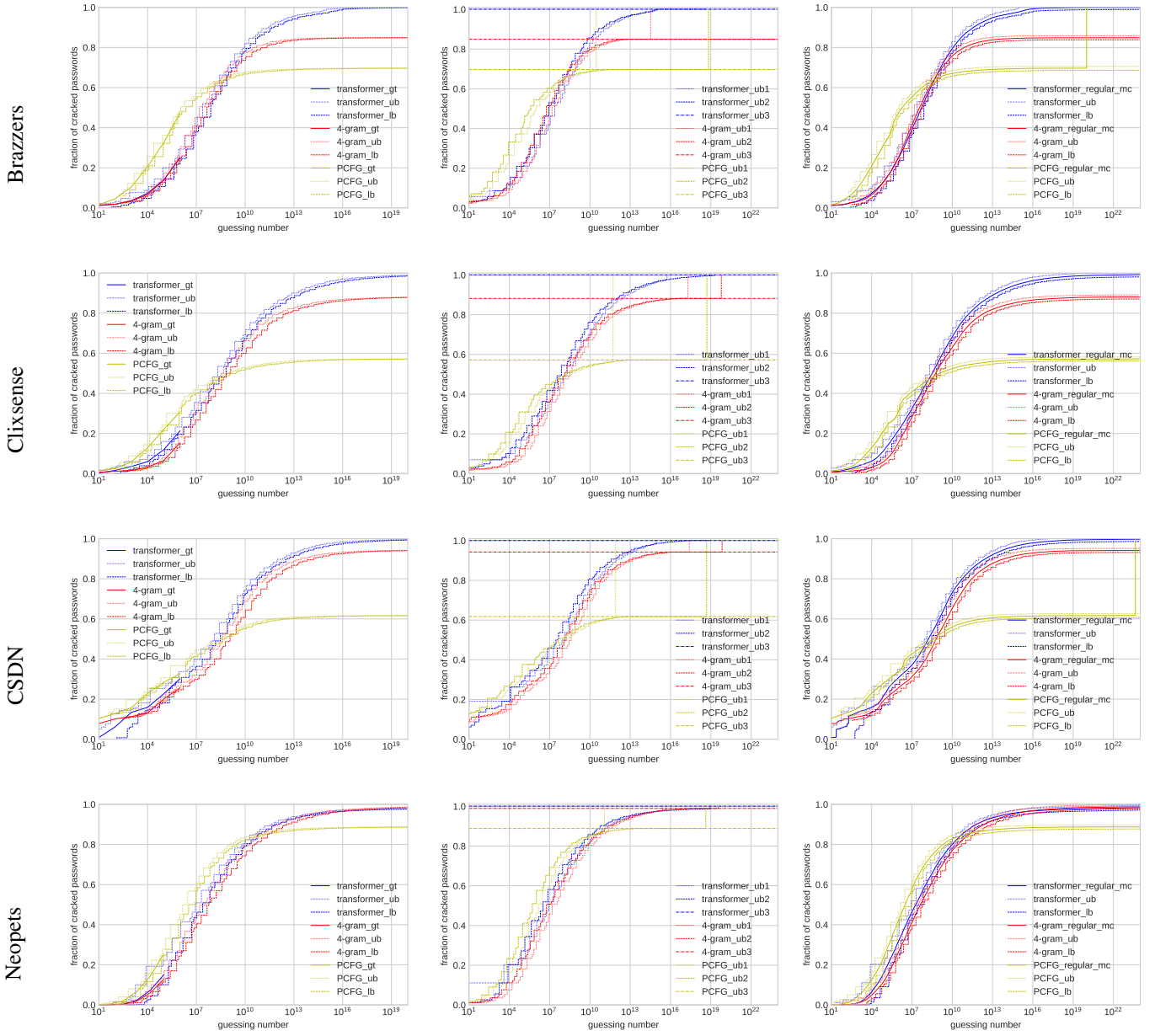


Figure 4: Confident Bounds for Additional Datasets

can bound $|\hat{G}_S^\phi(q) - G^\phi(q)|$ with any $t \geq 0$ as below:

$$\begin{aligned} \Pr[G^\phi(q) - \hat{G}_S^\phi(q) \leq t] &= \Pr\left[\mathbb{E}\left[\frac{1}{k} \sum_{i=1}^k X_i^\phi\right] - \frac{1}{k} \sum_{i=1}^k X_i^\phi \leq t\right] \\ &\geq 1 - \exp\left(\frac{-2t^2 k^2}{\sum_{i=1}^k \left(\frac{1}{q} - 0\right)^2}\right) = 1 - \exp(-2kt^2 q^2); \end{aligned}$$

$$\begin{aligned} \Pr[G^\phi(q) - \hat{G}_S^\phi(q) \geq -t] &= \Pr\left[\mathbb{E}\left[\frac{1}{k} \sum_{i=1}^k X_i^\phi\right] - \frac{1}{k} \sum_{i=1}^k X_i^\phi \geq -t\right] \\ &\geq 1 - \exp\left(\frac{-2t^2 k^2}{\sum_{i=1}^k \left(\frac{1}{q} - 0\right)^2}\right) = 1 - \exp(-2kt^2 q^2). \end{aligned}$$

Setting $t = \epsilon/q$ we have:

$$\begin{aligned} \Pr[G^\phi(q) \leq \hat{G}_S^\phi(q) + \frac{\epsilon}{q}] &\geq 1 - \exp(-2k\epsilon^2) \\ \Pr[G^\phi(q) \geq \hat{G}_S^\phi(q) - \frac{\epsilon}{q}] &\geq 1 - \exp(-2k\epsilon^2) \end{aligned}$$

□

Note that the actual guessing number $G(q)$ is between $G^{\text{EX}}(q) + 1$ and $G^{\text{IN}}(q)$, i.e., $G^{\text{EX}}(q) + 1 \leq G(q) \leq G^{\text{IN}}(q)$. Therefore, by directly applying the upper bound of $G^{\text{IN}}(q)$

and the lower bound of $G^{\text{EX}}(q)$ in Lemma 1, we can rigorously bound $G(q)$ with high confidence in Theorem 1.

Using Theorem 1 for any password y with probability p_y^M we can bound $G(y)$ with regular Monte Carlo estimates $\hat{G}_S^{\text{EX}}(y), \hat{G}_S^{\text{IN}}(y)$ as below:

$$\begin{aligned}\Pr[G(y) \leq \hat{G}_S^{\text{IN}}(y) + \frac{\epsilon}{p_y^M}] &\geq 1 - \exp(-2k\epsilon^2) \\ \Pr[G(y) \geq \hat{G}_S^{\text{EX}}(y) + 1 - \frac{\epsilon}{p_y^M}] &\geq 1 - \exp(-2k\epsilon^2)\end{aligned}$$

D.2. Complete Proof of Theorem 2

We start with lower bounding $G^\phi(q)$ by directly applying Markov's inequality:

Lemma 2. For any probability $q \in [0, 1]$ and any parameter $0 < \delta \leq 1$, we have:

$$\Pr[G^\phi(q) \geq \delta \cdot \hat{G}_S^\phi(q)] \geq 1 - \delta$$

where the randomness is taken over the selection of the sample set S from model M , and $\phi = \text{EX}$ or IN .

Proof. The proof of Lemma 1 shows that $\mathbb{E}(\hat{G}_S^\phi(q))_{S \leftarrow \mathcal{M}[S]} = G^\phi(q)$. Using Markov's inequality we have $\Pr[\hat{G}_S^\phi(q) \geq \frac{G^\phi(q)}{\delta}] \leq \delta$ for any $0 < \delta \leq 1$. Therefore, the theorem holds. \square

With probability at least $1 - \delta$, the true guessing number $G^{\text{EX}}(q)$ (resp. $G^{\text{IN}}(q)$) is off its estimation $\hat{G}_S^{\text{EX}}(q)$ (resp. $\hat{G}_S^{\text{IN}}(q)$) generated by a single execution of the regular Monte Carlo method for any password with probability q and any parameter $0 < \delta \leq 1$. Given that $G(q) \geq G^{\text{EX}}(q) + 1$, for any fixed parameter δ and ϵ Lemma 2 indicates a tighter lower bound of $G(q)$ than Theorem 1 when $q < \frac{\epsilon}{(1-\delta) \cdot \hat{G}_S^{\text{EX}}(q)}$. Note that we don't expect the confidence $1 - \delta$ to be very high (e.g. 0.99) because we want to increase δ (i.e., decrease the confidence $1 - \delta$) to tighten the lower bound $\delta \cdot \hat{G}_S^\phi(q)$ ($\phi = \text{EX}$ or IN). One example of picking δ is to set $\delta = \frac{1}{3}$. Once the bound is fixed (i.e. δ is fixed), we can increase the confidence by running the Monte Carlo sampling and estimation multiple times and outputting the median of all estimates.

Let $\mathbb{S} = \{S_1, S_2, \dots, S_n\}$ be n Monte Carlo sample sets where S_i is the sample set generated in the i th Monte Carlo simulation. For any password y , we define $\hat{G}_{\mathbb{S}, \text{med}}^\phi(y) = \text{median}(\{\hat{G}_{S_i}^\phi(y)\}_{1 \leq i \leq n})$ to be the median of n Monte Carlo estimates. Similar to $\hat{G}_S^\phi(\cdot)$, we abuse the notation $\hat{G}_{\mathbb{S}, \text{med}}^\phi(\cdot)$ and let $\hat{G}_{\mathbb{S}, \text{med}}^\phi(q)$ be the median of $\hat{G}_{S_1}^\phi(q), \dots, \hat{G}_{S_n}^\phi(q)$ for any probability $q \in [0, 1]$. Lemma 3 states that by repeating Monte Carlo sampling and estimation process n times we can improve the confidence from $1 - \delta$ to $1 - \exp(-2n\epsilon^2)$, where $0 \leq \epsilon \leq \frac{1}{2} - \delta$. We expect $1 - \exp(-2n\epsilon^2)$ to be a high confidence, e.g. 0.99.

Lemma 3. For any password probability $q \in [0, 1]$, and any parameters $0 \leq \delta \leq \frac{1}{2}$ and $0 \leq \epsilon \leq \frac{1}{2} - \delta$,

$$\Pr[G^\phi(q) \geq \delta \cdot \hat{G}_{\mathbb{S}, \text{med}}^\phi(q)] \geq 1 - \exp(-2n\epsilon^2)$$

where the randomness is taken over n sets of k Monte Carlo samples $\mathbb{S} = \{S_1, \dots, S_n\}$ from model M , and $\phi = \text{EX}$ or IN .

Proof. Let X_i^ϕ be the indicator variable corresponding to the i th execution of the regular Monte Carlo estimation with randomly selected sample set S_i . $X_i^\phi = 1$ if and only if the i th estimation $\hat{G}_{S_i}^\phi(q) \leq \frac{G^\phi(q)}{\delta}$; otherwise, $X_i^\phi = 0$. Lemma 2 indicates that $\mathbb{E}(X_i^\phi) \geq 1 - \delta$. Using Chernoff bound we have:

$$\Pr[\sum_{i=1}^n X_i^\phi \leq \frac{n}{2}] \leq \Pr[\sum_{i=1}^n X_i^\phi \leq n(1 - \delta - \epsilon)] \leq \exp(-2n\epsilon^2)$$

where we set $\delta \leq \frac{1}{2}$ and $\epsilon \leq \frac{1}{2} - \delta$. Then we have:

$$\Pr[\hat{G}_{\mathbb{S}, \text{med}}^\phi(q) \leq \frac{G^\phi(q)}{\delta}] \geq \Pr[\sum_{i=1}^n I_i \geq \frac{n}{2}] \geq 1 - \exp(-2n\epsilon^2)$$

\square

Recall that $G(q) \geq G^{\text{EX}}(q) + 1$. By directly applying the lower bound of $G^{\text{EX}}(q)$ in Lemma 3 on $G(q)$ Theorem 2 is proved.

Using Theorem 2 for any password y with probability p_y^M we can bound $G(y)$ with the median Monte Carlo estimate $\hat{G}_{\mathbb{S}, \text{med}}^{\text{EX}}(y)$ as $\Pr[G(y) \geq \delta \cdot \hat{G}_{\mathbb{S}, \text{med}}^{\text{EX}}(y) + 1] \geq 1 - \exp(-2n\epsilon^2)$.

Appendix E.

Complete Proofs of Bounds on Guessing Curve

In this section we provide the complete proofs of bounding $\lambda_{M,B,D}$ and $\lambda_{M,B}$. We start by providing the formal proof of Theorem 3.

Reminder of Theorem 3. For any guessing number $B \geq 0$ and any $0 \leq \epsilon \leq 1$, we have:

$$\begin{aligned}\Pr[\lambda_{M,B} \geq \lambda_{M,B,D} - \epsilon] &\geq 1 - \exp(-2|D|\epsilon), \quad \text{and} \\ \Pr[\lambda_{M,B} \leq \lambda_{M,B,D} + \epsilon] &\geq 1 - \exp(-2|D|\epsilon)\end{aligned}$$

where the randomness is taken over the sample set $D \leftarrow \mathcal{P}^{|D|}$.

Proof of Theorem 3. Consider $\lambda_{M,B,D}$ to be a function of the $|D|$ samples in D . For any two sets with the same number of iid samples from \mathcal{P} $D = \{d_1, \dots, d_i, \dots, d_{|D|}\}$ and $D' = \{d_1, \dots, d'_i, \dots, d_N\}$ that only differs on the one sample d_i and d'_i , the difference of $\lambda_{M,B,D}$ and $\lambda_{M,B,D'}$ is at most $1/|D|$, i.e., $|\lambda_{M,B,D} - \lambda_{M,B,D'}| \leq 1/|D|$. Therefore, using McDiarmid's inequality [13] we have:

$$\begin{aligned}\Pr[\lambda_{M,B} \geq \lambda_{M,B,D} - \epsilon] &\geq 1 - \exp(-2|D|\epsilon) \\ \Pr[\lambda_{M,B} \leq \lambda_{M,B,D} + \epsilon] &\geq 1 - \exp(-2|D|\epsilon)\end{aligned}$$

Theorem 3 shows that $\lambda_{M,B,D}$ is tightly concentrated around $\lambda_{M,B}$. Given this result our main task will be to develop high confidence upper/lower bounds on $\lambda_{M,B,D}$ which will immediately yield high-confidence upper/lower bound for $\lambda_{M,B}$ as a corollary.

E.1. Complete Proof of The General Framework

E.1.1. Step 1: Bounding with A Single Point. Recall that for any $\phi = \{\text{EX}, \text{IN}\}$ Lemma 1 and Lemma 3 prove several upper and lower bounds on $G^\phi(q)$ with high confidence. In general, for any probability $0 \leq q \leq 1$ we define $\text{UB}_{G^\phi, |S|}(q, S)$ (resp. $\text{LB}_{G^\phi, |S|}(q, S)$) to be an arbitrary upper (resp. lower) bound of $G^\phi(q)$ with error rate $\text{ERR}(\text{UB}_{G^\phi, |S|})$ (resp. $\text{ERR}(\text{LB}_{G^\phi, |S|})$), i.e., with randomness taken over the selection of sample set S from model M we have:

$$\begin{aligned} \Pr[G^\phi(q) \geq \text{LB}_{G^\phi, |S|}(q, S)] &\geq 1 - \text{ERR}(\text{LB}_{G^\phi, |S|}) \\ \Pr[G^\phi(q) \leq \text{UB}_{G^\phi, |S|}(q, S)] &\geq 1 - \text{ERR}(\text{UB}_{G^\phi, |S|}) \end{aligned}$$

We define BAD1_q^ϕ to be the bad event that $\text{LB}_{G^\phi, |S|}(q, S) > G^\phi(q)$ is not an acceptable lower bound given a set S of random samples from model M , and define BAD2_q^ϕ to be the bad event that $\text{UB}_{G^\phi, |S|}(q, S) < G^\phi(q)$ is no an acceptable upper bound. Then we have $\Pr[\text{BAD1}_q^\phi] \leq \text{ERR}(\text{LB}_{G^\phi, |S|})$ and $\Pr[\text{BAD2}_q^\phi] \leq \text{ERR}(\text{UB}_{G^\phi, |S|})$.

Observe that for any cracking model M and password dataset D , the percentage of cracked passwords $\lambda_{M, B, D}$ in D is monotonically increasing for guessing number B since an attacker can always crack more passwords or at least the same amount of passwords in D as the attacker makes more guesses, i.e., for any $0 \leq B_1 \leq B_2$ we have $\lambda_{M, B_1, D} \leq \lambda_{M, B_2, D}$. Therefore, when the bad event BAD1_q^ϕ does not happen, we have $\lambda_{M, B, D} \leq \lambda_{M, \text{LB}_{G^\phi, |S|}(q, S), D} \leq \lambda_{M, G^\phi(q), D}$ for any $B \leq \text{LB}_{G^\phi, |S|}(q, S)$; similarly, when the bad event BAD2_q^ϕ does not happen, we have $\lambda_{M, B, D} \geq \lambda_{M, \text{UB}_{G^\phi, |S|}(q, S), D} \geq \lambda_{M, G^\phi(q), D}$ for any $B \geq \text{UB}_{G^\phi, |S|}(q, S)$. The statement is formally described in the following lemma:

Lemma 4. Given a password dataset D containing iid samples from distribution \mathcal{P} , for any probability $0 \leq q \leq 1$ and parameter $0 \leq \epsilon \leq 1$ we have:

$$\begin{aligned} \Pr \left[\forall B \leq \text{LB}_{G^\phi, |S|}(q, S), \lambda_{M, B, D} \leq \lambda_{M, G^\phi(q), D} \mid \overline{\text{BAD1}_q^\phi} \right] &= 1 \\ \Pr \left[\forall B \geq \text{UB}_{G^\phi, |S|}(q, S), \lambda_{M, B, D} \geq \lambda_{M, G^\phi(q), D} \mid \overline{\text{BAD2}_q^\phi} \right] &= 1 \end{aligned}$$

where the randomness is taken over the sample set $S \leftarrow \mathcal{M}^{|S|}$, and $\phi = \text{EX}$ or IN .

Although $G^\phi(q)$ is unknown to us, $\lambda_{M, G^\phi(q), D}$ in fact can be easily computed by counting the number of passwords in D with probability outputted by model M greater than q (resp. greater than or equal to q), i.e., $\lambda_{M, G^{\text{EX}}(q), D} = \frac{|\{y \in D: p_y^M > q\}|}{|D|}$ (resp. $\lambda_{M, G^{\text{IN}}(q), D} = \frac{|\{y \in D: p_y^M \geq q\}|}{|D|}$), since p_y^M can be outputted by M for $y \in D$. Therefore, as shown in Lemma 4 we are able to use $\lambda_{M, G^\phi(q), D}$ as an upper (resp. lower) bound of $\lambda_{M, B, D}$ on some special point $B = \text{LB}_{G^\phi, |S|}(q, S)$ (resp. $B = \text{UB}_{G^\phi, |S|}(q, S)$) with an error rate that the bad event BAD1_q^ϕ (resp. BAD2_q^ϕ) happens. To guarantee the bounds hold with high confidence, we

consider the error rates $\text{ERR}(\text{LB}_{G^\phi, |S|})$ and $\text{ERR}(\text{UB}_{G^\phi, |S|})$ to be small.

E.1.2. Step 2: Bounding with Multiple Mesh Points.

Given a single probability point q , Lemma 4 allows us to bound $\lambda_{M, B, D}$ using a single value $\lambda_{M, G^\phi(q), D}$ with high confidence: for any guessing number $B \leq \text{LB}_{G^\phi, |S|}(q, S)$, $\lambda_{M, G^\phi(q), D}$ is an *upper* bound on $\lambda_{M, B, D}$; for any $B \geq \text{UB}_{G^\phi, |S|}(q, S)$, $\lambda_{M, G^\phi(q), D}$ is an *lower* bound on $\lambda_{M, B, D}$. However, this upper (resp. lower) bound can be significantly less tight as B gets much smaller than $\text{LB}_{G^\phi, |S|}(q, S)$ (resp. much larger than $\text{UB}_{G^\phi, |S|}(q, S)$), and even doesn't exist when $B > \text{LB}_{G^\phi, |S|}(q, S)$ (resp. $B < \text{UB}_{G^\phi, |S|}(q, S)$). To get complete and tighter upper and lower bounds for *all* guessing number $B \geq 0$, we need more probability points to generate multiple upper/lower bound values. In this section, we tighten the bounds and present complete upper and lower bounds on $\lambda_{M, B, D}$ (as well as $\lambda_{M, B}$) for any guessing number $B \geq 0$ by generating multiple upper/lower bound mesh points and union bounding the error rates among all mesh points.

Given a sequence of ℓ probability mesh points $Q = \{q_1, q_2, \dots, q_\ell\}$, Lemma 4 allow us to generate multiple upper (resp. lower) bounds $\lambda_{M, G^\phi(q_i), D}$ on the actual guessing curve $\lambda_{M, B, D}$ with high confidence for $B \leq \text{LB}_{G^\phi, |S|}(q_i, S)$ (resp. $B \geq \text{UB}_{G^\phi, |S|}(q_i, S)$). We denote $\hat{\lambda}_{M, B, D, S}^{ub}$ and $\hat{\lambda}_{M, B, D, S}^{lb}$ to be the tightest upper and lower bounds among them. Formally, we define $\hat{\lambda}_{M, B, D, S}^{ub}$ and $\hat{\lambda}_{M, B, D, S}^{lb}$ below as we state in equations (6) and (7):

$$\begin{aligned} \hat{\lambda}_{M, B, D, S}^{ub} &:= \min_{1 \leq i \leq \ell, B \leq \text{LB}_{G^{\text{IN}}, |S|}(q_i, S)} (\lambda_{M, G^{\text{IN}}(q_i), D}), \\ \hat{\lambda}_{M, B, D, S}^{lb} &:= \max_{1 \leq i \leq \ell, B \geq \text{UB}_{G^{\text{EX}}, |S|}(q_i, S)} (\lambda_{M, G^{\text{EX}}(q_i), D}). \end{aligned}$$

We set $\hat{\lambda}_{M, B, D, S}^{ub} = 1$ (resp. $\hat{\lambda}_{M, B, D, S}^{lb} = 0$) when no such i exists that satisfy the conditions, i.e., $B > \max_{1 \leq i \leq \ell} \{\text{LB}_{G^{\text{IN}}, |S|}(q_i, S)\}$ (resp. $B < \min_{1 \leq i \leq \ell} \{\text{UB}_{G^{\text{EX}}, |S|}(q_i, S)\}$). Note that we omit $\lambda_{M, G^{\text{EX}}(q_i), D}$ (resp. $\lambda_{M, G^{\text{IN}}(q_i), D}$) in the best upper (resp. lower) bound formula, since $\lambda_{M, G^{\text{IN}}(q_i), D}$ and $\lambda_{M, G^{\text{EX}}(q_i), D}$ are almost identical in most cases where not many distinct passwords are with the same probability q_i outputted by model M . Although including them in the lower (resp. upper) bound may slightly tighten the bound in some cases, this will also double the error rate. We can always tighten the bound more effectively by increasing the number of mesh points q_1, \dots, q_ℓ .

To guarantee the tightest upper and lower bounds on $\lambda_{M, B, D}$ hold (i.e., $\hat{\lambda}_{M, B, D, S}^{lb} \leq \lambda_{M, B, D} \leq \hat{\lambda}_{M, B, D, S}^{ub}$), none of the bad events $\text{BAD1}_{q_i}^{\text{IN}}, \text{BAD2}_{q_i}^{\text{EX}}$ for all $i = 1, 2, \dots, \ell$ should happen. Define $\text{BAD1}_Q = \bigvee_{i=1}^{\ell} \text{BAD1}_{q_i}^{\text{IN}}$ be the event that at least one of $\text{BAD1}_{q_i}^{\text{IN}}$ for $i = 1, \dots, \ell$ happens, and similarly define $\text{BAD2}_Q = \bigvee_{i=1}^{\ell} \text{BAD2}_{q_i}^{\text{EX}}$. Lemma 5 bounds the probabilities that the bad events $\text{BAD1}_Q, \text{BAD2}_Q$ happen.

Lemma 5. Given a sequence of ℓ probability mesh points $Q = \{q_1, q_2, \dots, q_\ell\}$, for any $1 \leq \epsilon \leq 1$ we have

$$\begin{aligned}\Pr[\text{BAD1}_Q] &\leq \ell \cdot \text{ERR}(\text{LB}_{G^{\text{IN}}, |S|}) \\ \Pr[\text{BAD2}_Q] &\leq \ell \cdot \text{ERR}(\text{UB}_{G^{\text{EX}}, |S|}).\end{aligned}$$

Proof. Recall that for any $1 \leq i \leq \ell$, $\Pr[\text{BAD1}_{q_i}^{\text{IN}}] \leq \text{ERR}(\text{LB}_{G^{\text{IN}}, |S|})$ and $\Pr[\text{BAD2}_{q_i}^{\text{EX}}] \leq \text{ERR}(\text{UB}_{G^{\text{EX}}, |S|})$. By applying union bounds the lemma holds. \square

As long as the bad events BAD1_Q and BAD2_Q happen with small probabilities, we can claim that $\hat{\lambda}_{M,B,D,S}^{ub}$ and $\hat{\lambda}_{M,B,D,S}^{lb}$ are the upper and lower bounds of $\lambda_{M,B,D}$ with high confidence as stated in Theorem 4.

Reminder of Theorem 4. Given a password dataset D containing iid samples from distribution \mathcal{P} and a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$, we have:

$$\begin{aligned}\Pr[\forall B \geq 0, \lambda_{M,B,D} \leq \hat{\lambda}_{M,B,D,S}^{ub}] &\geq 1 - \ell \cdot \text{ERR}(\text{LB}_{G^{\text{IN}}, |S|}) \\ \Pr[\forall B \geq 0, \lambda_{M,B,D} \geq \hat{\lambda}_{M,B,D,S}^{lb}] &\geq 1 - \ell \cdot \text{ERR}(\text{UB}_{G^{\text{EX}}, |S|})\end{aligned}$$

where the randomness is taken over the sample set S from model M .

Proof of Theorem 4. Lemma 4 shows that when the bad event $\text{BAD1}_{q_i}^{\text{IN}}$ (resp. $\text{BAD2}_{q_i}^{\text{EX}}$) does not occur, $\lambda_{M,G^{\text{IN}}(q_i),D}$ (resp. $\lambda_{M,G^{\text{EX}}(q_i),D}$) is an upper (resp. lower) bound of $\lambda_{M,B,D}$ for any $B \leq \text{LB}_{G^{\text{IN}}, |S|}(q, S)$ (resp. $B \geq \text{UB}_{G^{\text{EX}}, |S|}(q, S)$). Therefore, when all bad events do not occur, i.e., BAD1_Q (resp. BAD2_Q) does not occur, all the upper (resp. lower) bounds hold, indicating that $\lambda_{M,B,D} \leq \hat{\lambda}_{M,B,D,S}^{ub}$ (resp. $\lambda_{M,B,D} \geq \hat{\lambda}_{M,B,D,S}^{lb}$):

$$\begin{aligned}\Pr[\forall B \geq 0, \lambda_{M,B,D} \leq \hat{\lambda}_{M,B,D,S}^{ub} \mid \overline{\text{BAD1}_Q^\phi}] &= 1 \\ \Pr[\forall B \geq 0, \lambda_{M,B,D} \geq \hat{\lambda}_{M,B,D,S}^{lb} \mid \overline{\text{BAD2}_Q^\phi}] &= 1\end{aligned}$$

Then applying Lemma 5 we have:

$$\begin{aligned}\Pr[\forall B \geq 0, \lambda_{M,B,D} \leq \hat{\lambda}_{M,B,D,S}^{ub}] &\geq \Pr[\forall B \geq 0, \lambda_{M,B,D} \leq \hat{\lambda}_{M,B,D,S}^{ub} \mid \overline{\text{BAD1}_Q}] \Pr[\overline{\text{BAD1}_Q}] \\ &= \Pr[\overline{\text{BAD1}_Q}] \geq 1 - \ell \cdot \text{ERR}(\text{LB}_{G^{\text{IN}}, |S|}) \\ \Pr[\forall B \geq 0, \lambda_{M,B,D} \geq \hat{\lambda}_{M,B,D,S}^{lb}] &\geq \Pr[\forall B \geq 0, \lambda_{M,B,D} \geq \hat{\lambda}_{M,B,D,S}^{lb} \mid \overline{\text{BAD2}_Q}] \Pr[\overline{\text{BAD2}_Q}] \\ &= \Pr[\overline{\text{BAD2}_Q}] \geq 1 - \ell \cdot \text{ERR}(\text{UB}_{G^{\text{EX}}, |S|})\end{aligned}$$

Corollary 5 is an immediate corollary of Theorem 4 by applying Theorem 3.

E.2. Concrete Bounds on Guessing Curves

In this section, we apply the bounds on $\hat{G}^\phi(q)$ in Section 4 to the general proof in Section 5.1 and obtain concrete upper and lower bounds on $\lambda_{M,B,D}$ and $\lambda_{M,B}$.

E.2.1. Upper/Lower Bounds on Guessing Curves Using Hoeffding's Inequality. Lemma 1 proves upper and lower bounds on $G^\phi(q)$ using Hoeffding's Inequality. Applying Lemma 1 we set $\text{LB}_{G^{\text{IN}}, |S|}(q, S) = \hat{G}_S^{\text{IN}}(q) - \epsilon/q$, $\text{UB}_{G^{\text{EX}}, |S|}(q, S) = \hat{G}_S^{\text{EX}}(q) + \epsilon/q$, and $\text{ERR}(\text{LB}_{G^{\text{IN}}, |S|}) = \text{ERR}(\text{UB}_{G^{\text{EX}}, |S|}) = \exp(-2k\epsilon^2)$. Denote $\hat{\lambda}_{M,B,D,S}^{ub} = \hat{\lambda}_{M,B,D,S,\epsilon}^{ub1}$ and $\hat{\lambda}_{M,B,D,S}^{lb} = \hat{\lambda}_{M,B,D,S,\epsilon}^{lb1}$ where:

$$\begin{aligned}\hat{\lambda}_{M,B,D,S,\epsilon}^{ub1} &:= \min_{1 \leq i \leq \ell, B \leq \hat{G}_S^{\text{IN}}(q_i) - \epsilon/q_i} (\lambda_{M,G^{\text{IN}}(q_i),D}), \\ \hat{\lambda}_{M,B,D,S,\epsilon}^{lb1} &:= \max_{1 \leq i \leq \ell, B \geq \hat{G}_S^{\text{EX}}(q_i) + \epsilon/q_i} (\lambda_{M,G^{\text{EX}}(q_i),D}).\end{aligned}$$

Especially, we set $\hat{\lambda}_{M,B,D,S,\epsilon}^{ub1} = 1$ (resp. $\hat{\lambda}_{M,B,D,S,\epsilon}^{lb1} = 0$) when no such i exist that satisfy the conditions above, i.e., $B > \max_{1 \leq i \leq \ell} \{\hat{G}_S^{\text{IN}}(q_i) - \epsilon/q_i\}$ (resp. $B < \min_{1 \leq i \leq \ell} \{\hat{G}_S^{\text{EX}}(q_i) + \epsilon/q_i\}$).

Applying Theorem 4 we have:

Theorem 12. Given a password dataset D containing iid samples from distribution \mathcal{P} and a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$, for any parameter $0 \leq \epsilon \leq 1$, we have:

$$\begin{aligned}\Pr[\forall B \geq 0, \lambda_{M,B,D} \leq \hat{\lambda}_{M,B,D,S,\epsilon}^{ub1}] &\geq 1 - \ell \cdot \exp(-2k\epsilon^2) \\ \Pr[\forall B \geq 0, \lambda_{M,B,D} \geq \hat{\lambda}_{M,B,D,S,\epsilon}^{lb1}] &\geq 1 - \ell \cdot \exp(-2k\epsilon^2)\end{aligned}$$

where the randomness is taken over the sample set $S \leftarrow \mathcal{M}^k$ with size k .

Applying Theorem 3 (or Corollary 5) we have Theorem 6 that upper and lower bound $\lambda_{M,B}$ as stated in Section 5.2.

E.2.2. A Tighter Upper Bound for Small Mesh Points q_i . Lemma 3 tighten the lower bound on $G^\phi(q)$ for small q by taking the median estimate of n regular Monte Carlo estimation results. Applying Lemma 3 we set $\text{LB}_{G^{\text{IN}}, |S|}(q, S) = \delta \cdot \hat{G}_{S,\text{med}}^\phi(q)$, and $\text{ERR}(\text{LB}_{G^{\text{IN}}, |S|}) = \exp(-2n\epsilon^2)$. Denote $\hat{\lambda}_{M,B,D,S}^{ub} = \hat{\lambda}_{M,B,D,S,\delta}^{ub2}$ where:

$$\hat{\lambda}_{M,B,D,S,\delta}^{ub2} := \min_{1 \leq i \leq \ell, B \leq \delta \cdot \hat{G}_{S,\text{med}}^{\text{IN}}(q_i)} (\lambda_{M,G^{\text{IN}}(q_i),D})$$

Especially, we set $\hat{\lambda}_{M,B,D,S,\delta}^{ub2} = 1$ when no such i exists satisfying the condition $B \leq \delta \cdot \hat{G}_{S,\text{med}}^{\text{IN}}(q_i)$, i.e., $B > \max_{1 \leq i \leq \ell} \{\delta \cdot \hat{G}_{S,\text{med}}^{\text{IN}}(q_i)\}$. Applying Theorem 4 we have:

Theorem 13. Given a password dataset D containing iid samples from distribution \mathcal{P} and a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$, for any parameters $0 < \delta \leq \frac{1}{2}$, $0 \leq \epsilon \leq \frac{1}{2} - \delta$, we have:

$$\Pr[\forall B \geq 0, \lambda_{M,B,D} \leq \hat{\lambda}_{M,B,D,S,\delta}^{ub2}] \geq 1 - \ell \cdot \exp(-2n\epsilon^2)$$

where the randomness is taken over n Monte Carlo sample sets $\mathbb{S} = \{S_1, S_2, \dots, S_n\}$ each of which contains k samples from model M .

Applying Theorem 3 (or Corollary 5) we have Theorem 7 that upper bounds $\lambda_{M,B}$ as stated in Section 5.2.

E.3. A Trivial Upper Bound for Large Guessing Number

Recall that we define $\hat{\lambda}_{M,D}^{ub3} = \frac{1}{|D|} |y \in D : p_y^M > 0|$ be the percentage of passwords that will be eventually guessed by model M in Section 5.3. Then we have the following trivial upper bounds on $\lambda_{M,B,D}$ and $\lambda_{M,B}$:

Theorem 14. *Given a password cracking model M , for any guessing number $B \geq 0$, $\lambda_{M,B,D} \leq \hat{\lambda}_{M,D}^{ub3}$.*

Reminder of Theorem 8. *Given a password distribution \mathcal{P} and a password cracking model M , for any parameters $0 \leq \epsilon \leq 1$, we have:*

$$\Pr \left[\forall B \geq 0, \lambda_{M,B} \leq \hat{\lambda}_{M,D}^{ub3} + \epsilon \right] \geq 1 - \exp(-2|D|\epsilon^2)$$

where the randomness is taken over the dataset $D \leftarrow \mathcal{P}^{|D|}$. *Proof of Theorem 8.* Let $\sum_{y \in \mathcal{P}} p_y^M$ be the total probability mass of passwords in distribution \mathcal{P} that will be guessed with non-zero probability in model M . Then we have $\lambda_{M,B} \leq \sum_{y \in \mathcal{P}} p_y^M$ for any $B \geq 0$, since passwords with zero probability in M will never be guessed. Let $X_1, \dots, X_{|D|}$ be $|D|$ random variables where $X_i = 1$ if the i th sample in D has non-zero probability in M , and $X_i = 0$ otherwise. Then $\sum_{i=1}^{|D|} X_i = |D| \hat{\lambda}_{M,D}^{ub3}$ is the number of passwords in D that will eventually be guessed by model M . Note that $\hat{\lambda}_{M,D}^{ub3} = \frac{1}{|D|} \mathbb{E}(\sum_{i=1}^{|D|} X_i) = \sum_{y \in \mathcal{P}} p_y^M$. Using Chernoff bound, for any $0 \leq \epsilon \leq 1$ we have:

$$\Pr \left[\sum_{y \in \mathcal{P}} p_y^M \leq \hat{\lambda}_{M,D}^{ub3} + \epsilon \right] \geq 1 - \exp(-2|D|\epsilon^2).$$

Since $\forall B \geq 0, \lambda_{M,B} \leq \sum_{y \in \mathcal{P}} p_y^M$, we have:

$$\Pr \left[\forall B \geq 0, \lambda_{M,B} \leq \hat{\lambda}_{M,D}^{ub3} + \epsilon \right] \geq 1 - \exp(-2|D|\epsilon^2).$$

Appendix F. Bounds Under Password Composition Policies

In this section, we provide a complete description about our bounds guessing numbers and guessing curves under password composition policies.

F.1. Bounds on Guessing Number under Policies

For any password pwd we use $\text{Allowed}^{\mathbb{C}}(pwd)$ to describe the password policy \mathbb{C} an attacker would follow when making guesses. We set $\text{Allowed}^{\mathbb{C}}(pwd) = 1$ if and only if password pwd satisfies the policy \mathbb{C} (e.g. passwords must be at least 8 characters long). Similar to the actual guessing number $G(y)$ without applying any policies, we define $G^{\mathbb{C}}(y)$ to be the actual guessing number of password y cracked by an attacker who only makes guesses that satisfy the policy \mathbb{C} . We abuse the notation $G^{\mathbb{C}}(\cdot)$ and

let $G^{\mathbb{C}}(q)$ be the actual guessing number of a password with probability q in model M under the policy \mathbb{C} , for any $q \in [0, 1]$. Similar to $G^{\text{EX}}(\cdot)$ and $G^{\text{IN}}(\cdot)$, we define $G^{\mathbb{C}, \text{EX}}(q) = |\{z \in \mathcal{M} : p_z^M > q \wedge \text{Allowed}^{\mathbb{C}}(z) = 1\}|$ and $G^{\mathbb{C}, \text{IN}}(q) = |\{z \in \mathcal{M} : p_z^M \geq q \wedge \text{Allowed}^{\mathbb{C}}(z) = 1\}|$ to be the smallest and largest possible guessing number of a password with probability q in M , i.e., $G^{\mathbb{C}, \text{EX}}(q) + 1 \leq G^{\mathbb{C}}(q) \leq G^{\mathbb{C}, \text{IN}}(q)$. Given a set S of k Monte Carlo samples from model distribution \mathcal{M} , let $\hat{G}_S^{\mathbb{C}, \text{EX}}(q) = \frac{1}{k} \sum_{z \in S, p_z^M > q, \text{Allowed}^{\mathbb{C}}(z)=1} \frac{1}{p_z^M}$ be the estimate of $G^{\mathbb{C}, \text{EX}}(q)$, and $\hat{G}_S^{\mathbb{C}, \text{IN}}(q) = \frac{1}{k} \sum_{z \in S, p_z^M \geq q, \text{Allowed}^{\mathbb{C}}(z)=1} \frac{1}{p_z^M}$ be the estimate of $G^{\mathbb{C}, \text{IN}}(q)$.

In this section, we propose two theorems (Theorems 15 and 16) to upper and lower bound $G^{\mathbb{C}}(q)$ using the same techniques we proposed in Section 4 with small changes in the proofs to take the password policy \mathbb{C} into account.

F.1.1. First Upper/Lower Bound. To prove our first upper/lower bound in Theorem 15, we start by the following lemma:

Lemma 6. *Given a set S with k iid password samples sampled from distribution \mathcal{M} and a password policy \mathbb{C} , for any probability $0 \leq q \leq 1$ and any parameter $\epsilon \geq 0$, we have:*

$$\begin{aligned} \Pr[G^{\mathbb{C}, \phi}(q) \leq \hat{G}_S^{\mathbb{C}, \phi}(q) + \epsilon/q] &\geq 1 - \exp(-2k\epsilon^2) \\ \Pr[G^{\mathbb{C}, \phi}(q) \geq \hat{G}_S^{\mathbb{C}, \phi}(q) - \epsilon/q] &\geq 1 - \exp(-2k\epsilon^2) \end{aligned}$$

where the randomness is taken over the sample set $S \leftarrow \mathcal{M}^k$, and $\phi = \text{EX}$ or IN .

Proof. Given a password probability q and a set S with k password samples randomly sampled from distribution \mathcal{M} , we consider k independent random variables $X_1^\phi, \dots, X_k^\phi$ where for any $i = 1, \dots, k$ and $\phi = \text{EX}, \text{IN}$ we define:

$$\begin{aligned} X_i^{\text{EX}} &:= \begin{cases} \frac{1}{p_z^M} & \text{if the } i\text{th sampled password is } z \text{ and } p_z^M > q \\ & \text{and } \text{Allowed}^{\mathbb{C}}(z) = 1; \\ 0 & \text{otherwise.} \end{cases} \\ X_i^{\text{IN}} &:= \begin{cases} \frac{1}{p_z^M} & \text{if the } i\text{th sampled password is } z \text{ and } p_z^M \geq q \\ & \text{and } \text{Allowed}^{\mathbb{C}}(z) = 1; \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

Then we have $0 \leq X_i^\phi \leq \frac{1}{q}$ and the expectation of X_i^ϕ is $G^{\mathbb{C}, \phi}(q)$ as shown below:

$$\begin{aligned} \mathbb{E}(X_i^{\text{EX}}) &= \sum_{z \in \mathcal{M}, p_z^M > q} p_z^M \cdot \frac{1}{p_z^M} \\ &= |\{z \in \mathcal{M} : p_z^M > q \wedge \text{Allowed}^{\mathbb{C}}(z) = 1\}| = G^{\mathbb{C}, \text{EX}}(q) \\ \mathbb{E}(X_i^{\text{IN}}) &= \sum_{z \in \mathcal{M}, p_z^M \geq q} p_z^M \cdot \frac{1}{p_z^M} \\ &= |\{z \in \mathcal{M} : p_z^M \geq q \wedge \text{Allowed}^{\mathbb{C}}(z) = 1\}| = G^{\mathbb{C}, \text{IN}}(q) \end{aligned}$$

Observe that $\hat{G}_S^{\mathbb{C},\phi}(q) = \frac{1}{k} \sum_{i=1}^k X_i^\phi$ and $\mathbb{E}(\frac{1}{k} \sum_{i=1}^k X_i^\phi) = \frac{1}{k} \sum_{i=1}^k \mathbb{E}(X_i^\phi) = G^{\mathbb{C},\phi}(q)$. Using Hoeffding's inequality we can bound $|\hat{G}_S^{\mathbb{C},\phi}(q) - G^{\mathbb{C},\phi}(q)|$ with any $t \geq 0$ as below:

$$\begin{aligned} & \Pr[G^{\mathbb{C},\phi}(q) - \hat{G}_S^{\mathbb{C},\phi}(q) \leq t] \\ &= \Pr[\mathbb{E}(\frac{1}{k} \sum_{i=1}^k X_i^\phi) - \frac{1}{k} \sum_{i=1}^k X_i^\phi \leq t] \\ &\geq 1 - \exp\left(\frac{-2t^2 k^2}{\sum_{i=1}^k (\frac{1}{q} - 0)^2}\right) = 1 - \exp(-2kt^2 q^2); \end{aligned}$$

$$\begin{aligned} & \Pr[G^{\mathbb{C},\phi}(q) - \hat{G}_S^{\mathbb{C},\phi}(q) \geq -t] \\ &= \Pr[\mathbb{E}(\frac{1}{k} \sum_{i=1}^k X_i^\phi) - \frac{1}{k} \sum_{i=1}^k X_i^\phi \geq -t] \\ &\geq 1 - \exp\left(\frac{-2t^2 k^2}{\sum_{i=1}^k (\frac{1}{q} - 0)^2}\right) = 1 - \exp(-2kt^2 q^2). \end{aligned}$$

Setting $t = \epsilon/q$ we have:

$$\begin{aligned} \Pr[G^{\mathbb{C},\phi}(q) \leq \hat{G}_S^{\mathbb{C},\phi}(q) + \frac{\epsilon}{q}] &\geq 1 - \exp(-2k\epsilon^2) \\ \Pr[G^{\mathbb{C},\phi}(q) \geq \hat{G}_S^{\mathbb{C},\phi}(q) - \frac{\epsilon}{q}] &\geq 1 - \exp(-2k\epsilon^2) \end{aligned}$$

□

Since $G^{\mathbb{C},\text{EX}}(q) + 1 \leq G^{\mathbb{C}}(q) \leq G^{\mathbb{C},\text{IN}}(q)$, we can rigorously bound $[G^{\mathbb{C}}(q)]$ with high confidence by directly applying the upper bound of $G^{\mathbb{C},\text{IN}}(q)$ and the lower bound of $G^{\mathbb{C},\text{EX}}(q)$ in Lemma 6:

Theorem 15. *Given a set S with k iid password samples sampled from distribution \mathcal{M} and a password policy \mathbb{C} , for any password y and any parameter $\epsilon \geq 0$, we have:*

$$\begin{aligned} \Pr[G^{\mathbb{C}}(q) \leq \hat{G}_S^{\mathbb{C},\text{IN}}(q) + \epsilon/q] &\geq 1 - \exp(-2k\epsilon^2) \\ \Pr[G^{\mathbb{C}}(q) \geq \hat{G}_S^{\mathbb{C},\text{EX}}(q) - \epsilon/q + 1] &\geq 1 - \exp(-2k\epsilon^2) \end{aligned}$$

where the randomness is taken over the sample set $S \leftarrow \mathcal{M}^k$.

F.1.2. Second Upper Bound For Rare Passwords. Recall that we propose a second lower bound on $G(q)$ in Section 4.2 by first deriving a low-confidence bound using Markov inequality and then increasing the confidence by taking the median over n Monte Carlo estimates. Define $\hat{G}_{S,\text{med}}^{\mathbb{C},\phi}(q) = \text{median}_{1 \leq i \leq n} \{\hat{G}_{S_i}^{\mathbb{C},\phi}(q)\}$ to be the median of n Monte Carlo estimates computed using n Monte Carlo sample sets $\mathbb{S} = \{S_1, S_2, \dots, S_n\}$. Observe that the statements and proofs in Lemma 2 and Lemma 3 still hold if we replace $G^\phi(q), \hat{G}_S^\phi(q), \hat{G}_{S,\text{med}}^\phi(q)$ with $G^{\mathbb{C},\phi}(q), \hat{G}_S^{\mathbb{C},\phi}(q), \hat{G}_{S,\text{med}}^{\mathbb{C},\phi}(q)$ respectively. With high confidence we can lower bound $G^{\mathbb{C},\phi}(q)$ with $\hat{G}_{S,\text{med}}^{\mathbb{C},\phi}(q)$ as shown in the following lemma:

Lemma 7. *Given a password policy \mathbb{C} , for any password probability $q \in [0, 1]$, and any parameters $0 \leq \delta \leq \frac{1}{2}$ and $0 \leq \epsilon \leq \frac{1}{2} - \delta$,*

$$\Pr[G^{\mathbb{C},\phi}(q) \geq \delta \cdot \hat{G}_{S,\text{med}}^{\mathbb{C},\phi}(q)] \geq 1 - \exp(-2n\epsilon^2)$$

where the randomness is taken over n sets of k Monte Carlo samples $\mathbb{S} = \{S_1, \dots, S_n\}$ from model M , and $\phi = \text{EX}$ or IN .

Proof. Recall that the proof of Lemma 6 shows that the expectation of a Monte Carlo estimate $\hat{G}_{S_i}^{\mathbb{C},\phi}(q)$ is the actual value $G^{\mathbb{C},\phi}(q)$, i.e., $\mathbb{E}(\hat{G}_{S_i}^{\mathbb{C},\phi}(q))_{S \leftarrow \mathcal{M}^{|S_i|}} = G^{\mathbb{C},\phi}(q)$. Using Markov's inequality we have $\Pr[\hat{G}_{S_i}^{\mathbb{C},\phi}(q) \geq \frac{G^{\mathbb{C},\phi}(q)}{\delta}] \leq \delta$ for any $0 < \delta \leq 1$.

Let X_i^ϕ be the indicator variable corresponding to the i th execution of the regular Monte Carlo estimation with randomly selected sample set S_i . $X_i^\phi = 1$ if and only if the i th estimation $\hat{G}_{S_i}^{\mathbb{C},\phi}(q) \leq \frac{G^{\mathbb{C},\phi}(q)}{\delta}$; otherwise, $X_i^\phi = 0$. Using Chernoff bound we have:

$$\Pr[\sum_{i=1}^n X_i^\phi \leq \frac{n}{2}] \leq \Pr[\sum_{i=1}^n X_i^\phi \leq n(1 - \delta - \epsilon)] \leq \exp(-2n\epsilon^2)$$

where we set $\delta \leq \frac{1}{2}$ and $\epsilon \leq \frac{1}{2} - \delta$. Then we have:

$$\begin{aligned} \Pr[\hat{G}_{S,\text{med}}^{\mathbb{C},\phi}(q) \leq \frac{G^{\mathbb{C},\phi}(q)}{\delta}] &\geq \Pr[\sum_{i=1}^n I_i \geq \frac{n}{2}] \\ &\geq 1 - \exp(-2n\epsilon^2) \end{aligned}$$

□

Recall that $G^{\mathbb{C}}(q) \geq G^{\mathbb{C},\text{EX}}(q) + 1$. By directly applying the lower bound of $G^{\mathbb{C},\text{EX}}(q)$ in Lemma 7 we can bound $G^{\mathbb{C}}(q)$ as below:

Theorem 16. *Given a password policy \mathbb{C} , for any probability $0 \leq q \leq 1$ and any parameters $0 \leq \delta \leq \frac{1}{2}$ and $0 \leq \epsilon \leq \frac{1}{2} - \delta$,*

$$\Pr[G^{\mathbb{C}}(q) \geq \delta \cdot \hat{G}_{S,\text{med}}^{\mathbb{C},\text{EX}}(q) + 1] \geq 1 - \exp(-2n\epsilon^2)$$

where the randomness is taken over n sets of k Monte Carlo samples $\mathbb{S} = \{S_1, \dots, S_n\}$ from model M .

F.2. Bounds on Guessing Curve under Policies

In this section, we upper and lower bound the guessing curve of an attacker who follows some password policies. We define $\lambda_{M,B,D}^{\mathbb{C}} = \frac{1}{|D|} | \{y \in D : G^{\mathbb{C}}(y) \geq B \} |$ to be the guessing curve of set D under password policy \mathbb{C} , i.e., the percentage of passwords in D cracked by making the top B most probable guesses outputted by model M satisfying password policy \mathbb{C} . Define $\lambda_{M,B}^{\mathbb{C}} = \Pr_{y \leftarrow \mathcal{P}}[G^{\mathbb{C}}(y) \leq B]$ to be the probability of a password from some unknown distribution \mathcal{P} cracked by an attacker making the top B guesses of model M satisfying password policy \mathbb{C} .

Section 5 and Appendix E present a general proof on bounding guessing curves and three concrete bounds on guessing curves $\lambda_{M,B}, \lambda_{M,B,D}$ without any extra password policies \mathbb{C} . Notice that all the proofs and statements still hold when we take password composition policy \mathbb{C} into consideration by simply replacing $\lambda_{M,B}, \lambda_{M,B,D}, G^\phi(q), \hat{G}_S^\phi(q), \hat{G}_{S,\text{med}}^\phi(q)$ with $\lambda_{M,B}^{\mathbb{C}}, \lambda_{M,B,D}^{\mathbb{C}}, G^{\mathbb{C},\phi}(q), \hat{G}_S^{\mathbb{C},\phi}(q), \hat{G}_{S,\text{med}}^{\mathbb{C},\phi}(q)$ respectively.

Therefore, we are able to derive a generalized proof of bounding guessing curves $\lambda_{M,B}^{\mathbb{C}}, \lambda_{M,B,D}^{\mathbb{C}}$ under password policy \mathbb{C} using the same proof idea in Section 5.1 and Appendix E.1. We leave the theorem statements (Theorem 22 and Corollary 23) and formal proof of the general framework for bounding $\lambda_{M,B}^{\mathbb{C}}$ and $\lambda_{M,B,D}^{\mathbb{C}}$ Appendix G. Next we will present three different concrete upper bounds and one concrete lower bound on $\lambda_{M,B}^{\mathbb{C}}, \lambda_{M,B,D}^{\mathbb{C}}$ similar to the bounds in Section 5.2 and 5.3.

F.2.1. First Concrete Upper/Lower Bound. First, by applying Lemma 6 to equations (8) and (9) of the general framework in Appendix G we can denote the first upper/lower bound as below:

$$\hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},ub1} := \min_{1 \leq i \leq \ell, B \leq \hat{G}_{S,\phi}^{\mathbb{C}}(q_i) - \epsilon/q} \left(\lambda_{M,G^{\mathbb{C}}(q_i),D}^{\mathbb{C}} \right),$$

$$\hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},lb1} := \max_{1 \leq i \leq \ell, B \geq \hat{G}_{S,\phi}^{\mathbb{C}}(q_i) + \epsilon/q} \left(\lambda_{M,G^{\mathbb{C}}(q_i),D}^{\mathbb{C}} \right).$$

For completeness, we set $\hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},ub1} = 1$ (resp. $\hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},lb1} = 0$) if $B > \max_{1 \leq i \leq \ell} \{\hat{G}_{S,\phi}^{\mathbb{C}}(q_i) - \epsilon/q_i\}$ (resp. $B < \min_{1 \leq i \leq \ell} \{\hat{G}_{S,\phi}^{\mathbb{C}}(q_i) + \epsilon/q_i\}$). Then by applying Lemma 6 to Theorem 22 (a generalized theorem for upper/lower bounding the guessing curve) we have the following upper and lower bounds on guessing curve $\lambda_{M,B,D}^{\mathbb{C}}$ for arbitrary guessing number B :

Theorem 17. *Given a password dataset D containing iid samples from distribution \mathcal{P} , a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$ and a password policy \mathbb{C} , for any parameter $0 \leq \epsilon \leq 1$, we have:*

$$\Pr \left[\forall B \geq 0, \lambda_{M,B,D}^{\mathbb{C}} \leq \hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},ub1} \right] \geq 1 - \ell \cdot \exp(-2k\epsilon^2)$$

$$\Pr \left[\forall B \geq 0, \lambda_{M,B,D}^{\mathbb{C}} \geq \hat{\lambda}_{M,B,D,S,\epsilon}^{\mathbb{C},lb1} \right] \geq 1 - \ell \cdot \exp(-2k\epsilon^2)$$

where the randomness is taken over the sample set $S \leftarrow \mathcal{M}^k$ with size k .

Applying Corollary 23 (a generalized corollary for upper/lower bounding the guessing curve) we can bound $\lambda_{M,B}^{\mathbb{C}}$ as stated in Theorem 9 in Appendix B.

F.2.2. Second Upper Bound. Second, by applying Lemma 7 to equation (8) in the general framework in Appendix G we denote our second upper bound as:

$$\hat{\lambda}_{M,B,D,S,\delta}^{\mathbb{C},ub2} := \min_{1 \leq i \leq \ell, B \leq \delta \cdot \hat{G}_{S,med}^{\mathbb{C}}(q_i)} \left(\lambda_{M,G^{\mathbb{C}}(q_i),D}^{\mathbb{C}} \right).$$

For completeness, we set $\hat{\lambda}_{M,B,D,S,\delta}^{\mathbb{C},ub2} = 1$ if $B > \max_{1 \leq i \leq \ell} \{\delta \cdot \hat{G}_{S,med}^{\mathbb{C}}(q_i)\}$. Using Theorem 22 we have:

Theorem 18. *Given a password dataset D containing iid samples from distribution \mathcal{P} , a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$ and a password policy \mathbb{C} , for any parameters $0 \leq \delta \leq \frac{1}{2}, 0 \leq \epsilon \leq \frac{1}{2} - \delta$, we have:*

$$\Pr[\forall B \geq 0, \lambda_{M,B,D}^{\mathbb{C}} \leq \hat{\lambda}_{M,B,D,S,\delta}^{\mathbb{C},ub2}] \geq 1 - \ell \cdot \exp(-2n\epsilon^2)$$

where the randomness is taken over n Monte Carlo sample sets $\mathbb{S} = \{S_1, S_2, \dots, S_n\}$ each of which contains k samples from model M .

Applying Corollary 23 to Lemma 7 we can bound $\lambda_{M,B}^{\mathbb{C}}$ as stated in Theorem 10 in Appendix B.

F.2.3. A Trivial Upper Bound. We denote $\hat{\lambda}_{M,D}^{\mathbb{C},ub3} := \frac{1}{|D|} |y \in D : p_y^M > 0 \wedge \text{Allowed}^{\mathbb{C}}(y) = 1|$ be the percentage of passwords that will be eventually guessed by model M , due to the fact that passwords with zero probability in M or not satisfying the policy \mathbb{C} will never be guessed. Then we have the following trivial upper bound on $\lambda_{M,B,D}^{\mathbb{C}}$:

Theorem 19. *Given a password cracking model M and a password policy \mathbb{C} , for any guessing number $B \geq 0$, $\lambda_{M,B,D}^{\mathbb{C}} \leq \hat{\lambda}_{M,D}^{\mathbb{C},ub3}$.*

Following the same proof idea of Corollary 8, we can derive our third upper bound on $\lambda_{M,B}^{\mathbb{C}}$ using Theorem 19 as below:

Theorem 20. *Given a password distribution \mathcal{P} , a password cracking model M and a password policy \mathbb{C} , for any parameters $0 \leq \epsilon \leq 1$, we have:*

$$\Pr \left[\forall B \geq 0, \lambda_{M,B}^{\mathbb{C}} \leq \hat{\lambda}_{M,D}^{\mathbb{C},ub3} + \epsilon \right] \geq 1 - \exp(-2|D|\epsilon^2)$$

where the randomness is taken over the dataset $D \leftarrow \mathcal{P}^{|D|}$.

Proof. Let $\sum_{y \in \mathcal{P}, \text{Allowed}^{\mathbb{C}}(y)=1} p_y^M$ be the total probability mass of passwords in distribution \mathcal{P} that will be guessed with non-zero probability in model M . Then we have $\lambda_{M,B}^{\mathbb{C}} \leq \sum_{y \in \mathcal{P}, \text{Allowed}^{\mathbb{C}}(y)=1} p_y^M$ for any $B \geq 0$, since passwords with zero probability in M will never be guessed. Let $X_1, \dots, X_{|D|}$ be $|D|$ random variables where $X_i = 1$ if the i th sample in D has non-zero probability in M and satisfying the policy \mathbb{C} and $X_i = 0$ otherwise. Then $\sum_{i=1}^{|D|} X_i = |D| \hat{\lambda}_{M,D}^{\mathbb{C},ub3}$ is the number of passwords in D that will eventually be guessed by model M . Note that $\hat{\lambda}_{M,D}^{\mathbb{C},ub3} = \frac{1}{|D|} \mathbb{E}(\sum_{i=1}^{|D|} X_i) = \sum_{y \in \mathcal{P}, \text{Allowed}^{\mathbb{C}}(y)=1} p_y^M$. Using Chernoff bound, for any $0 \leq \epsilon \leq 1$ we have:

$$\Pr \left[\sum_{y \in \mathcal{P}, \text{Allowed}^{\mathbb{C}}(y)=1} p_y^M \leq \hat{\lambda}_{M,D}^{\mathbb{C},ub3} + \epsilon \right] \geq 1 - \exp(-2|D|\epsilon^2).$$

Since $\forall B \geq 0, \lambda_{M,B}^{\mathbb{C}} \leq \sum_{y \in \mathcal{P}, \text{Allowed}^{\mathbb{C}}(y)=1} p_y^M$, we have:

$$\Pr \left[\forall B \geq 0, \lambda_{M,B}^{\mathbb{C}} \leq \hat{\lambda}_{M,D}^{\mathbb{C},ub3} + \epsilon \right] \geq 1 - \exp(-2|D|\epsilon^2). \quad \square$$

Appendix G.

General Framework of Bounding Guessing Curve Under Password Composition Policy

In this section we present a general framework for upper/lower bounding the curves $\lambda_{M,B,D}^{\mathbb{C}}$ and $\lambda_{M,B}^{\mathbb{C}}$ as the guessing budget B varies from small to large, which follows the same idea of proving the general framework for

bounding $\lambda_{M,B,D}$ and $\lambda_{M,B}$ without policy restriction \mathbb{C} in Section 5 and Appendix E.

Recall that for an attacker who only makes guesses that satisfy the policy \mathbb{C} , $\lambda_{M,B,D}^{\mathbb{C}}$ denotes the fraction of passwords in dataset D that would be cracked within B guesses, and that $\lambda_{M,B}^{\mathbb{C}}$ denotes the probability that randomly sampled password would be cracked within B guesses. Given a dataset D of independent samples from an *unknown* password distribution, our first observation is that the expected value of $\lambda_{M,B,D}^{\mathbb{C}}$ (over the random selection of D) is simply $\lambda_{M,B}^{\mathbb{C}}$ and, if D is large enough, the random variable $\lambda_{M,B,D}^{\mathbb{C}}$ is tightly concentrated around its mean — see Theorem 21. Given this result our main task will be to develop high confidence upper/lower bounds on $\lambda_{M,B,D}^{\mathbb{C}}$ which will immediately yield high-confidence upper/lower bound for $\lambda_{M,B}^{\mathbb{C}}$ as a corollary.

Theorem 21. *Given a password policy \mathbb{C} , for any guessing number $B \geq 0$ and any $0 \leq \epsilon \leq 1$, we have:*

$$\begin{aligned}\Pr[\lambda_{M,B}^{\mathbb{C}} \geq \lambda_{M,B,D}^{\mathbb{C}} - \epsilon] &\geq 1 - \exp(-2|D|\epsilon) \\ \Pr[\lambda_{M,B}^{\mathbb{C}} \leq \lambda_{M,B,D}^{\mathbb{C}} + \epsilon] &\geq 1 - \exp(-2|D|\epsilon)\end{aligned}$$

where the randomness is taken over the sample set $D \leftarrow \mathcal{P}^{|D|}$.

Proof. Consider $\lambda_{M,B,D}^{\mathbb{C}}$ to be a function of the $|D|$ samples in D . For any two sets with the same number of iid samples from \mathcal{P} $D = \{d_1, \dots, d_i, \dots, d_{|D|}\}$ and $D' = \{d_1, \dots, d'_i, \dots, d_N\}$ that only differs on the one sample d_i and d'_i , the difference of $\lambda_{M,B,D}^{\mathbb{C}}$ and $\lambda_{M,B,D'}^{\mathbb{C}}$ is at most $1/|D|$, i.e., $|\lambda_{M,B,D}^{\mathbb{C}} - \lambda_{M,B,D'}^{\mathbb{C}}| \leq 1/|D|$. Therefore, using McDiarmid's inequality [13] we have:

$$\begin{aligned}\Pr[\lambda_{M,B}^{\mathbb{C}} \geq \lambda_{M,B,D}^{\mathbb{C}} - \epsilon] &\geq 1 - \exp(-2|D|\epsilon) \\ \Pr[\lambda_{M,B}^{\mathbb{C}} \leq \lambda_{M,B,D}^{\mathbb{C}} + \epsilon] &\geq 1 - \exp(-2|D|\epsilon)\end{aligned}$$

□

G.1. Step 1: Bounding with A Single Point

Recall that for any $\phi = \{\text{EX}, \text{IN}\}$ Lemma 6 and Lemma 7 prove several upper and lower bounds on $G^{\mathbb{C},\phi}(q)$ with high confidence. In general, for any probability $0 \leq q \leq 1$ we define $\text{UB}_{G^{\mathbb{C},\phi,|S|}}(q, S)$ (resp. $\text{LB}_{G^{\mathbb{C},\phi,|S|}}(q, S)$) to be an arbitrary upper (resp. lower) bound of $G^{\mathbb{C},\phi}(q)$ with error rate $\text{ERR}(\text{UB}_{G^{\mathbb{C},\phi,|S|}})$ (resp. $\text{ERR}(\text{LB}_{G^{\mathbb{C},\phi,|S|}})$), i.e., with randomness taken over the selection of sample set S from model M we have:

$$\begin{aligned}\Pr[G^{\mathbb{C},\phi}(q) \geq \text{LB}_{G^{\mathbb{C},\phi,|S|}}(q, S)] &\geq 1 - \text{ERR}(\text{LB}_{G^{\mathbb{C},\phi,|S|}}) \\ \Pr[G^{\mathbb{C},\phi}(q) \leq \text{UB}_{G^{\mathbb{C},\phi,|S|}}(q, S)] &\geq 1 - \text{ERR}(\text{UB}_{G^{\mathbb{C},\phi,|S|}})\end{aligned}$$

We define BAD1_q^{ϕ} to be the bad event that $\text{LB}_{G^{\mathbb{C},\phi,|S|}}(q, S) > G^{\mathbb{C},\phi}(q)$ is not an acceptable lower bound given a set S of random samples from model M , and define BAD2_q^{ϕ} to be the bad event that $\text{UB}_{G^{\mathbb{C},\phi,|S|}}(q, S) < G^{\mathbb{C},\phi}(q)$ is no an acceptable upper

bound. Then we have $\Pr[\text{BAD1}_q^{\phi}] \leq \text{ERR}(\text{LB}_{G^{\mathbb{C},\phi,|S|}})$ and $\Pr[\text{BAD2}_q^{\phi}] \leq \text{ERR}(\text{UB}_{G^{\mathbb{C},\phi,|S|}})$.

Observe that for any cracking model M and password dataset D , the percentage of cracked passwords $\lambda_{M,B,D}^{\mathbb{C}}$ in D is monotonically increasing for guessing number B since an attacker can always crack more passwords or at least the same amount of passwords in D as the attacker makes more guesses, i.e., for any $0 \leq B_1 \leq B_2$ we have $\lambda_{M,B_1,D}^{\mathbb{C}} \leq \lambda_{M,B_2,D}^{\mathbb{C}}$. Therefore, when the bad event BAD1_q^{ϕ} does not happen, we have $\lambda_{M,B,D}^{\mathbb{C}} \leq \lambda_{M,\text{LB}_{G^{\mathbb{C},\phi,|S|}}(q,S),D}^{\mathbb{C}} \leq \lambda_{M,G^{\mathbb{C},\phi}(q),D}^{\mathbb{C}}$ for any $B \leq \text{LB}_{G^{\mathbb{C},\phi,|S|}}(q, S)$; similarly, when the bad event BAD2_q^{ϕ} does not happen, we have $\lambda_{M,B,D}^{\mathbb{C}} \geq \lambda_{M,\text{UB}_{G^{\mathbb{C},\phi,|S|}}(q,S),D}^{\mathbb{C}} \geq \lambda_{M,G^{\mathbb{C},\phi}(q),D}^{\mathbb{C}}$ for any $B \geq \text{UB}_{G^{\mathbb{C},\phi,|S|}}(q, S)$. The statement is formally described in the following lemma:

Lemma 8. *Given a password dataset D containing iid samples from distribution \mathcal{P} and a password policy \mathbb{C} , for any probability $0 \leq q \leq 1$ and parameter $0 \leq \epsilon \leq 1$ we have:*

$$\begin{aligned}\Pr[\forall B \leq \text{LB}_{G^{\mathbb{C},\phi,|S|}}(q, S), \lambda_{M,B,D}^{\mathbb{C}} \leq \lambda_{M,G^{\mathbb{C},\phi}(q),D}^{\mathbb{C}} \mid \overline{\text{BAD1}_q^{\phi}}] &= 1 \\ \Pr[\forall B \geq \text{UB}_{G^{\mathbb{C},\phi,|S|}}(q, S), \lambda_{M,B,D}^{\mathbb{C}} \geq \lambda_{M,G^{\mathbb{C},\phi}(q),D}^{\mathbb{C}} \mid \overline{\text{BAD2}_q^{\phi}}] &= 1\end{aligned}$$

where the randomness is taken over the sample set $S \leftarrow \mathcal{M}^{|S|}$, and $\phi = \text{EX}$ or IN .

Although $G^{\mathbb{C},\phi}(q)$ is unknown to us, $\lambda_{M,G^{\mathbb{C},\phi}(q),D}^{\mathbb{C}}$ in fact can be easily computed by counting the number of passwords in D satisfying rule \mathbb{C} with probability outputted by model M greater than q (resp. greater than or equal to q), i.e., $\lambda_{M,G^{\mathbb{C},\text{EX}}(q),D}^{\mathbb{C}} = \frac{1}{|D|} |y \in D : p_y^M > q \wedge \text{Allowed}^{\mathbb{C}}(y) = 1|$ (resp. $\lambda_{M,G^{\mathbb{C},\text{IN}}(q),D}^{\mathbb{C}} = \frac{1}{|D|} |y \in D : p_y^M \geq q \wedge \text{Allowed}^{\mathbb{C}}(y) = 1|$), since p_y^M can be easily outputted by M for $y \in D$. Therefore, as shown in Lemma 8 we are able to use $\lambda_{M,G^{\mathbb{C},\phi}(q),D}^{\mathbb{C}}$ as an upper (resp. lower) bound of $\lambda_{M,B,D}^{\mathbb{C}}$ on some special point $B = \text{LB}_{G^{\mathbb{C},\phi,|S|}}(q, S)$ (resp. $B = \text{UB}_{G^{\mathbb{C},\phi,|S|}}(q, S)$) with an error rate that the bad event BAD1_q^{ϕ} (resp. BAD2_q^{ϕ}) happens. To guarantee the bounds hold with high confidence, we consider the error rates $\text{ERR}(\text{LB}_{G^{\mathbb{C},\phi,|S|}})$ and $\text{ERR}(\text{UB}_{G^{\mathbb{C},\phi,|S|}})$ to be small.

G.2. Step 2: Bounding with Multiple Mesh Points

Given a single probability point q , Lemma 8 allows us to bound $\lambda_{M,B,D}^{\mathbb{C}}$ using a single value $\lambda_{M,G^{\mathbb{C},\phi}(q),D}^{\mathbb{C}}$ with high confidence: for any guessing number $B \leq \text{LB}_{G^{\mathbb{C},\phi,|S|}}(q, S)$, $\lambda_{M,G^{\mathbb{C},\phi}(q),D}^{\mathbb{C}}$ is an *upper* bound on $\lambda_{M,B,D}^{\mathbb{C}}$; for any $B \geq \text{UB}_{G^{\mathbb{C},\phi,|S|}}(q, S)$, $\lambda_{M,G^{\mathbb{C},\phi}(q),D}^{\mathbb{C}}$ is a *lower* bound on $\lambda_{M,B,D}^{\mathbb{C}}$. However, this upper (resp. lower) bound can be significantly less tight as B gets much smaller than $\text{LB}_{G^{\mathbb{C},\phi,|S|}}(q, S)$ (resp. much larger than $\text{UB}_{G^{\mathbb{C},\phi,|S|}}(q, S)$), and even doesn't exist when $B > \text{LB}_{G^{\mathbb{C},\phi,|S|}}(q, S)$ (resp.

$B < \text{UB}_{G^{\mathbb{C}}, \phi, |S|}(q, S)$). To get complete and tighter upper and lower bounds for *all* guessing number $B \geq 0$, we need more probability points to generate multiple upper/lower bound values. In this section, we tighten the bounds and present complete upper and lower bounds on $\lambda_{M,B,D}^{\mathbb{C}}$ (as well as $\lambda_{M,B}^{\mathbb{C}}$) for any guessing number $B \geq 0$ by generating multiple upper/lower bound mesh points and union bounding the error rates among all mesh points.

Given a sequence of ℓ probability mesh points $Q = \{q_1, q_2, \dots, q_\ell\}$, Lemma 8 allow us to generate multiple upper (resp. lower) bounds $\lambda_{M,G^{\mathbb{C}}, \phi(q_i), D}^{\mathbb{C}}$ on the actual guessing curve $\lambda_{M,B,D}^{\mathbb{C}}$ with high confidence for $B \leq \text{LB}_{G^{\mathbb{C}}, \phi, |S|}(q_i, S)$ (resp. $B \geq \text{UB}_{G^{\mathbb{C}}, \phi, |S|}(q_i, S)$). We define $\hat{\lambda}_{M,B,D,S}^{\mathbb{C}, ub}$ and $\hat{\lambda}_{M,B,D,S}^{\mathbb{C}, lb}$ as below to be the tightest upper and lower bounds among them:

$$\hat{\lambda}_{M,B,D,S}^{\mathbb{C}, ub} := \min_{1 \leq i \leq \ell, B \leq \text{LB}_{G^{\mathbb{C}}, \text{IN}, |S|}(q_i, S)} \left(\lambda_{M,G^{\mathbb{C}}, \text{IN}(q_i), D}^{\mathbb{C}} \right), \quad (8)$$

$$\hat{\lambda}_{M,B,D,S}^{\mathbb{C}, lb} := \max_{1 \leq i \leq \ell, B \geq \text{UB}_{G^{\mathbb{C}}, \text{EX}, |S|}(q_i, S)} \left(\lambda_{M,G^{\mathbb{C}}, \text{EX}(q_i), D}^{\mathbb{C}} \right). \quad (9)$$

We set $\hat{\lambda}_{M,B,D,S}^{\mathbb{C}, ub} = 1$ (resp. $\hat{\lambda}_{M,B,D,S}^{\mathbb{C}, lb} = 0$) when no such i exists that satisfy the conditions, i.e., $B > \max_{1 \leq i \leq \ell} \{\text{LB}_{G^{\mathbb{C}}, \text{IN}, |S|}(q_i, S)\}$ (resp. $B < \min_{1 \leq i \leq \ell} \{\text{UB}_{G^{\mathbb{C}}, \text{EX}, |S|}(q_i, S)\}$). Note that we omit $\lambda_{M,G^{\mathbb{C}}, \text{EX}(q_i), D}^{\mathbb{C}}$ (resp. $\lambda_{M,G^{\mathbb{C}}, \text{IN}(q_i), D}^{\mathbb{C}}$) in the best upper (resp. lower) bound formula, since $\lambda_{M,G^{\mathbb{C}}, \text{IN}(q_i), D}^{\mathbb{C}}$ and $\lambda_{M,G^{\mathbb{C}}, \text{EX}(q_i), D}^{\mathbb{C}}$ are almost identical in most cases where not many distinct passwords are with the same probability q_i outputted by model M . Although including them in the lower (resp. upper) bound may slightly tighten the bound in some cases, this will also double the error rate. We can always tighten the bound more effectively by increasing the number of mesh points q_1, \dots, q_ℓ .

To guarantee the tightest upper and lower bounds on $\lambda_{M,B,D}^{\mathbb{C}}$ hold (i.e., $\hat{\lambda}_{M,B,D,S}^{\mathbb{C}, lb} \leq \lambda_{M,B,D}^{\mathbb{C}} \leq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, ub}$), none of the bad events $\text{BAD1}_{q_i}^{\text{IN}}, \text{BAD2}_{q_i}^{\text{EX}}$ for all $i = 1, 2, \dots, \ell$ should happen. Define $\text{BAD1}_Q = \bigvee_{i=1}^{\ell} \text{BAD1}_{q_i}^{\text{IN}}$ be the event that at least one of $\text{BAD1}_{q_i}^{\text{IN}}$ for $i = 1, \dots, \ell$ happens, and similarly define $\text{BAD2}_Q = \bigvee_{i=1}^{\ell} \text{BAD2}_{q_i}^{\text{EX}}$. Lemma 5 bounds the probabilities that the bad events $\text{BAD1}_Q, \text{BAD2}_Q$ happen.

Lemma 9. Given a sequence of ℓ probability mesh points $Q = \{q_1, q_2, \dots, q_\ell\}$ and a password policy \mathbb{C} , for any $1 \leq \epsilon \leq 1$ we have

$$\begin{aligned} \Pr[\text{BAD1}_Q] &\leq \ell \cdot \text{ERR}(\text{LB}_{G^{\mathbb{C}}, \text{IN}, |S|}) \\ \Pr[\text{BAD2}_Q] &\leq \ell \cdot \text{ERR}(\text{UB}_{G^{\mathbb{C}}, \text{EX}, |S|}). \end{aligned}$$

Proof. Recall that for any $1 \leq i \leq \ell$, $\Pr[\text{BAD1}_{q_i}^{\text{IN}}] \leq \text{ERR}(\text{LB}_{G^{\mathbb{C}}, \text{IN}, |S|})$ and $\Pr[\text{BAD2}_{q_i}^{\text{EX}}] \leq \text{ERR}(\text{UB}_{G^{\mathbb{C}}, \text{EX}, |S|})$. By applying union bounds the lemma holds. \square

As long as the bad events BAD1_Q and BAD2_Q happen with small probabilities, we can claim that $\hat{\lambda}_{M,B,D,S}^{\mathbb{C}, ub}$

$\hat{\lambda}_{M,B,D,S}^{\mathbb{C}, lb}$ are the upper and lower bounds of $\lambda_{M,B,D}^{\mathbb{C}}$ with high confidence as stated in Theorem 22:

Theorem 22. Given a password dataset D containing iid samples from distribution \mathcal{P} , a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$ and a password policy \mathbb{C} , we have:

$$\begin{aligned} \Pr \left[\forall B \geq 0, \lambda_{M,B,D}^{\mathbb{C}} \leq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, ub} \right] &\geq 1 - \ell \cdot \text{ERR}(\text{LB}_{G^{\mathbb{C}}, \text{IN}, |S|}) \\ \Pr \left[\forall B \geq 0, \lambda_{M,B,D}^{\mathbb{C}} \geq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, lb} \right] &\geq 1 - \ell \cdot \text{ERR}(\text{UB}_{G^{\mathbb{C}}, \text{EX}, |S|}) \end{aligned}$$

where the randomness is taken over the sample set S from model M .

Proof. Lemma 8 shows that when the bad event $\text{BAD1}_{q_i}^{\text{IN}}$ (resp. $\text{BAD2}_{q_i}^{\text{EX}}$) does not occur, $\lambda_{M,G^{\mathbb{C}}, \text{IN}(q_i), D}^{\mathbb{C}}$ (resp. $\lambda_{M,G^{\mathbb{C}}, \text{EX}(q_i), D}^{\mathbb{C}}$) is an upper (resp. lower) bound of $\lambda_{M,B,D}^{\mathbb{C}}$ for any $B \leq \text{LB}_{G^{\mathbb{C}}, \text{IN}, |S|}(q_i, S)$ (resp. $B \geq \text{UB}_{G^{\mathbb{C}}, \text{EX}, |S|}(q_i, S)$). Therefore, when all bad events do not occur, i.e., BAD1_Q (resp. BAD2_Q) does not occur, all the upper (resp. lower) bounds hold, indicating that $\lambda_{M,B,D}^{\mathbb{C}} \leq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, ub}$ (resp. $\lambda_{M,B,D}^{\mathbb{C}} \geq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, lb}$):

$$\begin{aligned} \Pr \left[\forall B \geq 0, \lambda_{M,B,D}^{\mathbb{C}} \leq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, ub} \mid \overline{\text{BAD1}_Q} \right] &= 1 \\ \Pr \left[\forall B \geq 0, \lambda_{M,B,D}^{\mathbb{C}} \geq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, lb} \mid \overline{\text{BAD2}_Q} \right] &= 1 \end{aligned}$$

Then applying Lemma 9 we have:

$$\begin{aligned} \Pr \left[\forall B \geq 0, \lambda_{M,B,D}^{\mathbb{C}} \leq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, ub} \right] &\geq \Pr \left[\forall B \geq 0, \lambda_{M,B,D}^{\mathbb{C}} \leq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, ub} \mid \overline{\text{BAD1}_Q} \right] \Pr \left[\overline{\text{BAD1}_Q} \right] \\ &= \Pr \left[\overline{\text{BAD1}_Q} \right] \geq 1 - \ell \cdot \text{ERR}(\text{LB}_{G^{\mathbb{C}}, \text{IN}, |S|}) \\ \Pr \left[\forall B \geq 0, \lambda_{M,B,D}^{\mathbb{C}} \geq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, lb} \right] &\geq \Pr \left[\forall B \geq 0, \lambda_{M,B,D}^{\mathbb{C}} \geq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, lb} \mid \overline{\text{BAD2}_Q} \right] \Pr \left[\overline{\text{BAD2}_Q} \right] \\ &= \Pr \left[\overline{\text{BAD2}_Q} \right] \geq 1 - \ell \cdot \text{ERR}(\text{UB}_{G^{\mathbb{C}}, \text{EX}, |S|}) \end{aligned}$$

\square

As an immediate corollary of Theorem 22 by applying Theorem 21 we can upper/lower bound $\lambda_{M,B}^{\mathbb{C}}$ as below:

Corollary 23. Given a password distribution \mathcal{P} , a sequence of probability mesh points $Q = \{q_1, \dots, q_\ell\}$ and a password policy \mathbb{C} , for any guessing number $B > 0$ and any parameters $0 \leq \epsilon \leq 1$, we have:

$$\begin{aligned} \Pr \left[\lambda_{M,B}^{\mathbb{C}} \leq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, ub} + \epsilon \right] &\geq 1 - \ell \cdot \text{ERR}(\text{LB}_{G^{\mathbb{C}}, \text{IN}, |S|}) - \exp(-2|D|\epsilon^2) \\ \Pr \left[\lambda_{M,B}^{\mathbb{C}} \geq \hat{\lambda}_{M,B,D,S}^{\mathbb{C}, lb} - \epsilon \right] &\geq 1 - \ell \cdot \text{ERR}(\text{UB}_{G^{\mathbb{C}}, \text{EX}, |S|}) - \exp(-2|D|\epsilon^2) \end{aligned}$$

where the randomness is taken over the sample set S from model M and the dataset $D \leftarrow \mathcal{P}^{|D|}$.