

# TIFS'17- Zipf's Law in Passwords

概述: 用户口令的实际分布一直是一个开放问题, 本文深入探讨这个问题, 基于自然语言处理和计算统计技术, 通过14个大型的真实数据集提出了两个Zipf-like model用于描述口令的分布, 其中PDF-Zipf可以描述流行口令, 而CDF-Zipf适用于整个口令数据集。接下来根据口令分布, 提出评估口令数据集强度的方法, 并且简述了Zipf model的应用场景及将会影响的领域。最后对未来工作进行展望。

优点: ①解决了安全领域的一个基本问题: 用户生成的口令是如何分布的。

②提出的观点都进行了详尽的实验论证和剖析。

③对该领域的背景、现状和未来工作都进行讨论(很详实的Conclusion)

问题: ①本文可进一步分析不同数据集的特征对最终分布的影响, 如对于CDF-Zipf Model,  $C$ 和 $s$ 与数据集特征的关系。

②第IV节的B部分, 文章设计了从完美Zipf分布中取样的实验, 共考虑33组可能影响拟合的实验, 包括Zipf分布、样本大小和LF, 应简要分析这样选取实验参数的原因。

③相对于PDF-Zipf和Clauaset'09的方案, 在拟合数据集时, CDF-Zipf的时间复杂度较大, 在实际测试中时间(14.67 hours)远超另外两种(32.40 seconds和69.39 seconds), 应简要讨论原因。由于CDF-Zipf在KS统计量 $D$ 和口令集覆盖范围方面表现较优, 应简要讨论提升拟合速度的方法。

▲口令策略会影响用户口令的选择, 若设计较好, 口令策略可以在保证usability的情况下改善口令安全。→ 口令策略设计者应如何评估策略强度? 管理员应如何选择口令策略?

→ 由于口令分布(一个网站内)可能随时间产生变化, 安全管理员应评估口令分布的强度, 并据此强度调整自己的策略。

→ 如何精确评测基于某策略的口令分布的强度引出另一基础问题: 用户生成的口令服从什么分布函数? !关于口令分布的不合实际的假设通常会引起严重的安全性和可用性问题。

研究口令分布的重大意义(其中的两个)。

①Schechter等提出了基于流行度的口令生成策略, 当一个口令超过给定的流行度阈值 $T$ 则被拒绝, 可有效抵抗漫步猜测攻击, 但 $T$ 若无法合适选择, 可能影响安全性和可用性。而当得知用户选择的口令分布, 可更合理地选取阈值。

②在衡量基于口令的认证方案抗离线字典空间的能力时, 主要考虑口令搜索空间大小, 并根据单次猜测成本来衡量敌手的成本-收益。

▲口令分布的强度评测可以依据单次攻击方法。

▲口令数据集: 14个大规模真实口令集, 来自不同服务器类型和规模的网站, 用户的地理位置、语言信仰和文化背景也有不同。另外, 对于不同数据集中的口令长度、字符构成等进行了统计, 并分析了产生这些特征的可能原因。

▲两个统计学技术

①线性回归: 拟合可能存线性关系的观测数据, 求解两个变量之间的关系(常用最小二乘法)

决定系数 $R^2 \in [0, 1]$ 衡量回归线与实际数据点拟合优劣的一个评价指标:  $R^2$ 越接近1越好。

②Kolmogorov-Smirnov检验: 针对离散数据的最流行的非参数检验方法, 通过假设检验测量样本和理论分布模型之间的“距离”。→ 实际数据的累积分布函数CDF  $F_n(x)$ 与理论分布CDF  $F(x)$ 的距离:  $D = \sup |F_n(x) - F(x)|$ ,  $D \in [0, 1]$ 表示两条CDF的最大距离,  $D$ 越小越好。

→ 由于口令分布不是固定不变, 采用Monte Carlo方法: 先使用理论模型拟合给定口令集得到口令分布函数并计算 $D$ , 接下来使用从实际口令集拟合的分布参数生成合成数据集, 使用理论模型拟合合成数据集并计算 $D'$ , 使用 $p$ -value代表 $D'$ 大于 $D$ 的比例。若 $p$ -value较大, 可以表明实际口令数据分布与理论模型无显著差异。(It is safer)

用户口令中的Zipf's law. (Zipf定律最早用较画自然语言中单词“排列 vs. 出现频率”的关系,  $f_r = \frac{C}{r^s}$ )

▲去除每个口令中最不流行的口令(次数少于3次或5次), 进行线性回归拟合, → 口令分布服从相似的规律, 对于口令集 $DS$ 一个口令的排列 $r$ 及频率 $f_r$ 满足 $f_r = \frac{C}{r^s}$ ,  $C$ 和 $s$ 是依赖 $DS$ 的常数。

取对数:  $\log f_r = \log C - s \cdot \log r$ ,  $R^2$ 很接近于1(与数据集大小及数据泄露方式有关)

Zipf理论模型模拟口令的概率分布函数PDF, 故称之为PDF-Zipf模型(注:  $s \in [0, 1]$ )



在KS检验中,  $p$ -value 比较小, 分析得出假设不一致的结论的原因

① 非流行部分的口令 ( $f_r < L_F$ ) 占整个口令的比重为 50%~92%, 构成长尾部分。

根据大数定律, 其在口令集出现的频率不能反应在口令分布中的真实概率。

通过将给定口令  $p_{wi}$  的每次观测值作为一个服从 Bernoulli 分布的随机变量, 得到相对标准误差:

→  $RSE: \frac{s}{\mu} \approx \sqrt{\frac{1}{f_{p_{wi}}}}$ , 则当  $f_{p_{wi}}$  较大时, 真实概率  $p_{p_{wi}}$  才可用经验概率  $\frac{f_{p_{wi}}}{f_{os}}$  拟合。

⇒ 不流行的口令较大程度上会干扰拟合

② 收集到的口令集相对于根据策略产生的口令集仍然太小, 由于 Zipf 分布具有口令概率的多式递减特性, 小样本中低概率事件将压倒高概率事件。

▲ 进一步证明用户生成的口令样本服从 Zipf 定律, 从完美 Zipf 理论分布中随机抽取一些样本, 又观测拟合特性。⇒ 考虑了 3 组可能影响拟合的变量, Zipf 分布 3 种, 样本大小 8 种, 最小频率阈值  $L_F$  5 种

→ 共 120 个实验, 表明给定 Zipf 分布 (确定  $N$  和  $S$ ), 刚开始更大的  $L_F$  会导致更好的回归, 而  $L_F$  进一步增加会导致拟合恶化。

→ 找到最优拟合, 表明样本量越大, 更大比例的流行事件将用于拟合。

收集到的口令远小于样本空间, 故应先去除非流行部分的口令。

▲ CDF-Zipf Model (直接对口令的 CDF 建模)。

$F_r = C' \cdot Y^{S'}$ ,  $F_r$  为排名第  $r$  的口令的累积频率,  $C'$  和  $S'$  是取决于口令集的常数 (用最小二乘法)

由  $F_r$  为阶梯函数:  $f_r = F_r - F_{r-1} = C' \cdot Y^{S'} - C' \cdot (Y-1)^{S'} (若视为连续函数  $f_r = \frac{d(F_r)}{dY} = C' \cdot S' \cdot Y^{S'-1}$ )$

KS 检验统计量  $D$  为 0.006170~0.045874 (均值 0.018457), 比 PDF-Zipf 小。

当给定一个足够大的样本, 小且非显著差异会在统计学中认为是显著的, 故  $p$ -value 较小, 而当样本容量较小 ( $10^5$ ) 时, 多数 CDF-Zipf 模型拟合可通过 KS 检验 ( $p$ -value  $> 0.01$ )

→ KS 统计量  $D$  的最大改进量被限制在区间  $[0, 0.018457]$  内, CDF-Zipf model 可改进空间有限。

通过与 Clauset'09 的方案对比, 在 KS 统计量  $D$  方面: CDF-Zipf 和 Clauset'09 优于 PDF-Zipf,

对于百万规模以上的口令集, CDF-Zipf 通常更好, 百万级以下的, Clauset'09 表现更好。

② 从口令集覆盖范围方面, CDF-Zipf 表现最佳, PDF-Zipf 次之, Clauset'09 较差

另外, 给定中等的计算资源, 三种模型均可在可接受的时间内完成拟合。

▲ 数据集涵盖了不同的服务类型、规模、泄漏原因等, 表明了 Zipf model 的广泛适用性

Zipf' law 带来的三点启示。

1> 基于口令的密码协议: 在 CDF-Zipf model 中, 敌手优势刻画为:  $Adv_{A,p}(k) = C' \cdot Q(k)^{S'} + E(k)$ 。

2> 口令创建策略: 对基于流行度的口令生成策略的阈值做更合理的选择。

$P(\text{rank} \leq r) = C' \cdot Y^{S'}$ ,  $P(\text{rank} \leq N) = C' \cdot N^{S'} = 1 \Rightarrow N = (\frac{1}{C'})^{\frac{1}{S'}}$ ,  $N$  为独立口令的个数。

前  $N$  个独立口令的累积频率:  $P(\text{rank} \leq N) = C' \cdot (N \cdot N^{S'}) = C' \cdot (N \cdot (\frac{1}{C'})^{\frac{1}{S'}})^{S'} = C' \cdot Y^{S'} \cdot (\frac{1}{C'}) = Y^{S'}$ 。

若口令频率  $f_r$  为连续型变量,  $f_r = \frac{d(F_r)}{dY} = \frac{d(C' \cdot Y^{S'})}{dY} = C' \cdot S' \cdot Y^{S'-1} \Rightarrow T = C' \cdot S' \cdot Y^{S'-1}$

则流行度高于  $T$  的比例为:  $\eta = (\frac{T}{C' \cdot S'})^{\frac{1}{S'-1}} \cdot (C')^{\frac{1}{S'}}$

潜在受影响的用户比例:  $W_p(Y) = \eta^{S'}$ , 而真实受影响的用户比例为:  $W_a(Y) = (1 - S') \cdot Y^{S'}$

3> 口令分布强度指标  $\alpha$ -guesswork: 在两种情形下实际上是非参数化的, 在 Zipf model 下,  $\alpha$ -guesswork 的应用可大为简化

▲ 口令数据集的强度指标

对于一个聪明的攻击者来说, 给定数据集  $X$ , 使用破解结果  $\lambda_X(n)$  作为强度指标:

$\lambda_X(n) = \sum_{i=1}^n P_i(X) = \frac{1}{|DS|} \sum_{i=1}^n f_i(X)$ ,  $|DS|$  为数据集大小,  $n$  为猜测数。

由  $P_i(X) = \frac{f_i(X)}{|DS|} = F_r(X) - F_{r-1}(X) \approx C' \cdot Y^{S'} - C' \cdot (Y-1)^{S'}$ 。

⇒  $\lambda_X(n) = \sum_{i=1}^n P_i(X) = F_n(X) \approx C' \cdot n^{S'} = \lambda_X(n)$

通过真实数据集作图,  $\lambda_X(n)$  和  $\lambda_X(n)$  的图像大多重叠, 通过与其它方案对比, 基于 CDF-Zipf 的强度指标  $\lambda_X(n)$  最好。