# PROBABILISTIC PASSWORD MODELING: PREDICTING PASSWORDS WITH MACHINE LEARNING

JAY DESTORIES

MENTOR: ELIF YAMANGIL

## Abstract

Many systems use passwords as the primary means of authentication. As the length of a password grows, the search space of possible passwords grows exponentially. Despite this, people often fail to create unpredictable passwords. This paper will explore the problem of creating a probabilistic model for describing the distribution of passwords among the set of strings. This will help us gain insight into the relative strength of passwords as well as alternatives to existing methods of password candidate generation for password recovery tools like John the Ripper. This paper will consider methods from the field of natural language processing and evaluate their efficacy in modeling human-generated passwords.

## Contents

## 1. Introduction

Since ancient times, passwords have been used as a means for authentication. Given the rise of personal computers and the internet, passwords are now commonly used by the general public for authentication with services like social networks or online banking. The purpose of a password, and authentication in general in the scope of computer systems, is to confirm an identity. For instance, when using an online banking service you provide a password so that you can perform actions on an account, but others cannot. This sense of security relies on a password being something that you can remember, but others cannot discover. The security of a system protected by a password is dependent on that password being difficult to guess. This is helped by the fact that an increase in the length of a password exponentially increases the number of possible passwords. Despite this, people often create unoriginal and easy to guess passwords. In a 2014 leak of approximately 5 million gmail account passwords, 47779 users had used the password "123456". In fact, there are so many instances of different people using the same passwords that one of the dominant ways to guess passwords in password recovery tools like John the Ripper is to simply guess every password in a list of leaked passwords. It is apparent that there is some underlying thought process in coming up with a password that leads to people creating similar or even identical passwords. This paper will attempt to find that structure in the form of a probabilistic model. It will consider some techniques from the field of natural language processing for inferring the structure of the language. If there is a good way to model the language of passwords, it could be used to generate new password guesses that don't already exist in a wordlist, improving the performance of password recovery tools when wordlists are exhausted. It could also evaluate the strength of a password in terms of the expected number of guesses to find it given a model.

## 2. Classifying the Language of Passwords

2.1. **Regular Languages.**

2.2. **Context-Free Languages.**

2.3. **Context-Sensitive and Recursively Enumerable Languages.**

## 3. Probability Modeling

3.1. **N-Gram Language Model.**

3.2. **Hidden Markov Model.**

3.3. **Context-Free Grammar.**

## 4. Generating Passwords

4.1. **Baseline.**

4.2. **N-Gram Language Model.**

4.3. **Hidden Markov Model.**

4.4. **Context-Free Grammar.**

## 5. Applications

### 5.1. **Password Strength.**

### 5.2. **Password Recovery.**

## 6. Conclusion

## References

[1] Ghahramani, Z. (2001) An Introduction to Hidden Markov Models and Bayesian Networks. *International Journal of Pattern Recognition and Artificial Intelligence.*

[2] Rosenfeld, Ronald. Two Decades of Statistical Language Modeling: Where do we go from here? www.cs.cmu.edu/~roni/papers/survey-slm-IEEE-PROC-0004.pdf

## 7. APPENDIX

### 7.1. Language Model Evaluation.

```python
"""
    Test the performance of various password models
"""

import lmgenerator

# basic wordlist attack
def baselineGenerator(training_corpus):
    for pwd in training_corpus:
        yield pwd
    while True:
        yield ""

# See how many things in test_corpus the generator can guess with some number of
# tries
def testGenerator(gen, test_corpus, tries):
    found = 0
    test_set = set(test_corpus)
    guesses = set()
    for i in xrange(tries):
        guess = gen.next()
        if not guess in guesses:
            guesses.update([guess])
            if guess in test_set:
                found += 1
    return found

def testCorpora(training_corpus, test_corpus):
    print "First 5 training passwords: ", training_corpus[:5]
    print "First 5 test passwords: ", test_corpus[:5]

    tries = 100000
    baseline = testGenerator(baselineGenerator(training_corpus), test_corpus, tries)
    print "Baseline wordlist attack -- %d tries: %d." % (tries, baseline)
    bigramlmgen = lmgenerator.SimplePrunedBigramLMGenerator(training_corpus)
    bigramlm = testGenerator(bigramlmgen, test_corpus, tries)
    print "Bigram LM attack -- %d tries: %d." % (tries, bigramlm)


def main():
    print "###################################################################"
    print "Training corpus: rockyou"
    print "Test corpus: gmail"
    print "###################################################################"
    rockyou_nocount = open('corpora/rockyou_nocount', 'r')
    training_corpus = [pwd.rstrip() for pwd in rockyou_nocount]
    gmail_nocount = open('corpora/gmail_nocount', 'r')
    gmail_corpus = [pwd.rstrip() for pwd in gmail_nocount]
    test_corpus = gmail_corpus[:-5000]
    held_out_corpus = gmail_corpus[-5000:]
```

```
51        testCorpora ( training_corpus , test_corpus )
52
53
54  if __name__ == "__main__" :
55        main ( )
```

7.2. **N-Gram Language Model Implementation.**

```
1   from math import log , exp
2   import random
3
4   start_token = "<S>"
5   end_token = "</S>"
6
7   def Preprocess ( corpus ):
8         return [[ start_token ] + [ token for token in pwd] + [ end_token ] for pwd in corpus]
9
10  class BigramLM :
11        def __init__ ( self ):
12              self . bigram_counts = {}
13              self . unigram_counts = {}
14
15        def Train ( self , training_corpus ):
16              training_set = Preprocess ( training_corpus )
17              for pwd in training_set :
18                    for i in xrange ( len (pwd) − 1 ):
19                          token = pwd [ i ]
20                          next_token = pwd [ i + 1]
21                          if not token in self . unigram_counts :
22                                self . unigram_counts [ token ] = 0
23                          if not token in self . bigram_counts :
24                                self . bigram_counts [ token ] = {}
25                          if not next_token in self . bigram_counts [ token ] :
26                                self . bigram_counts [ token ][ next_token ] = 0
27                          self . unigram_counts [ token ] += 1
28                          self . bigram_counts [ token ][ next_token ] += 1
29
30        def GenerateSample ( self ):
31              sample = [ start_token ]
32              while not sample [ −1] == end_token :
33                    selector = random . uniform ( 0 , self . unigram_counts [ sample [ −1]])
34                    sum_bc = 0
35                    for bigram in self . bigram_counts [ sample [ −1]]:
36                          sum_bc += self . bigram_counts [ sample [ −1]][ bigram ]
37                          if sum_bc > selector :
38                                sample . append ( bigram )
39                                break
40              return ''. join ( sample [1: −1])
41
42        # gets the ( unsmoothed ) probability of a string given the bigramlm
43  #     def StringLogProbability ( self , string ):
44
45
46  def BigramLMGenerator ( training_corpus ):
```

```
47        lm = BigramLM()
48        lm.Train(training_corpus)
49        while True:
50            yield lm.GenerateSample()
51
52    def SimplePrunedBigramLMGenerator(training_corpus):
53        tries = set()
54        gen = BigramLMGenerator(training_corpus)
55        while True:
56            pwd = gen.next()
57            if not pwd in tries:
58                tries.update([pwd])
59                yield pwd
```

### 7.3. Hidden Markov Model Implementation.

```
1    from math import log, exp, log1p
2    import random
3    from memoize import memoize
4
5    start_token = "<S>"
6    end_token = "</S>"
7    wildcard_token = "<*>"
8
9    # reduce floating point imprecision in adding probabilities in log space
10   def SumLogProbs(lps):
11       # ln(e^lp1 + e^lp2) == ln(e^lp2 (e^(lp1 - lp2) + 1)) = ln(e^(lp1 - lp2) + 1) + lp2
12       def adderhelper(lp1, lp2):
13           return log1p(exp(lp1 - lp2)) + lp2 if lp2 > lp1 else log1p(exp(lp2 - lp1)) + lp1
14       return reduce(adderhelper, lps)
15
16
17   def Preprocess(corpus):
18       return [[start_token] + [token for token in pwd] + [end_token] for pwd in corpus]
19
20   # gets a count-length array of random probabilities summing to s
21   def RandomPartition(count, s):
22       if count is 1:
23           return [s]
24       split_prob = (random.random() * .4 + .2) * s
25       split_count = len(count) / 2
26       return RandomPartition(split_count, split_prob) + \
27               RandomPartition(count - split_count, s - split_prob)
28
29   # gets an array of log probabilities [p1, p2, ...] where e^p1 + e^p2 + ... = 1
30   def RandomLogProbs(count):
31       total = 4000000000
32       partition = RandomPartition(count, total)
33       return [log(p) - log(total) for p in partition]
34
35
36   class BigramHMM:
37       def __init__(self, vocabulary, state_count):
38           self.o_vocabulary = set(vocabulary)
```

```python
39          self.states = range(state_count)
40          self.start_probability = {state: prob for (state, prob) in zip(self.states, Rand
41          self.transition_probability = {state1: {state2: prob for (state2, prob) in (self
42          self.end_probability = {state: prob for (state, prob) in zip(self.states, Random
43          self.emission_probability = {state: {symbol: prob for (symbol, prob) in zip(voca
44
45      @memoize
46      def ForwardMatrix(pwd):
47          bp = [{state: None for state in self.states} for c in pwd]
48
49          # initialization
50          bp[0] = {state: self.start_probability[state] + self.emission_probability[state]
51
52          # recursion
53          for i in xrange(1, len(pwd)):
54              bp[i] = {state: SumLogProbs(map(lambda p: bp[i - 1][p] + self.transition_pro
55
56          return bp
57
58
59      @memoize
60      def BackwardMatrix(pwd):
61          bp = [{state: None for state in self.states} for c in pwd]
62
63          # initialization
64          bp[len(pwd) - 1] = {state: self.end_probability[state] for state in self.states}
65
66          # recursion
67          for i in reversed(xrange(0, len(pwd) - 1)):
68              bp[i] = {state: SumLogProbs(map(lambda n: bp[i + 1][n] + self.transition_pro
69
70          return bp
71
72
73      @memoize
74      def ForwardProbability(step, state, pwd):
75          matrix = self.ForwardMatrix(pwd)
76          if state == wildcard_token:
77              return SumLogProbs(matrix[step].values())
78          return matrix[step][state]
79          """
80          if step is 0:
81              return self.start_probability[end_state]
82
83          bp = [{state: None for state in self.states} for c in pwd]
84
85          # initialization
86          bp[0] = {state: self.start_probability[state] + self.emission_probability[state]
87
88          # recursion
89          for i in xrange(1, step - 1):
90              bp[i] = {state: sum(map(lambda p: bp[i - 1][p] + self.transition_probability
91
```

```python
92              # termination
93              if end_state == wildcard_token:
94                  return sum(map(lambda state: sum(map(lambda p: bp[step - 1][p] + self.transi
95
96              return sum(map(lambda p: bp[step - 1][p] + self.transition_probability[p][end_st
97          """
98
99          @memoize
100         def BackwardProbability(step, state, pwd):
101             matrix = self.BackwardMatrix(pwd)
102             return matrix[step][state]
103             """
104             last_step = len(pwd) - 1
105             if step == last_step:
106                 return self.end_probability[start_state]
107
108             bp = [{state: None for state in self.states} for c in pwd]
109
110             # initialization
111             bp[last_step] = {state: self.end_probability[state] for state in self.states}
112
113             # recursion
114             for i in reversed(xrange(step + 1, last_step - 1)):
115                 bp[i] = {state: sum(map(lambda n: bp[i + 1][n] + self.transition_probability
116
117             # termination
118             return sum(map(lambda n: bp[step + 1][n] + self.transition_probability[start_sta
119         """
120
121         @memoize
122         def TimeStateProbability(step, state, pwd):
123             return self.ForwardProbability(step, state, pwd) + \
124                    self.BackwardProbability(step, state, pwd) - \
125                    self.ForwardProbability(len(pwd) - 1, wildcard_token, pwd)
126
127         @memoize
128         def StateTransitionProbability(step, state1, state2, pwd):
129             return self.ForwardProbability(step, state1, pwd) + \
130                    self.BackwardProbability(step + 1, state2, pwd) + \
131                    self.transition_probability[state1][state2] + \
132                    self.emission_probability[state2][pwd[step + 1]] - \
133                    self.ForwardProbability(len(pwd) - 1, wildcard_token, pwd)
134
135         def ForwardBackward():
136             # for now assume convergence in constant number of iterations
137             for i in xrange(10):
138                 # expectation
139
140                 # maximization
141
142                 # reset memos
143                 self.ForwardMatrix.reset()
144                 self.BackwardMatrix.reset()
```

```
145                    self . ForwardProbability . reset ()
146                    self . BackwardProbability . reset ()
147                    self . TimeStateProbability . reset ()
148                    self . StateTransitionProbability . reset ()
```

## 7.4. Context-Free Grammar Implementation.