# TECHNICAL REPORT ON MACHINE LEARNING EXPERIMENTS FOR THE MÖBIUS FUNCTION

DAVID LOWRY-DUDA
LAST UPDATED: 2024.10.19

ABSTRACT. Last week, I was at the Mathematics and Machine Learning program at Harvard's Center of Mathematical Sciences and Applications. The underlying topic was on number theory and I've been studying various number theoretic problems from a machine learning perspective.

This is a technical report, including details related to actually running the code and analyzing the results.

## CONTENTS

## 1. INTRODUCTION

I've been computing several experiments related to estimating the Mobius function $\mu(n)$. Previous machine learning experiments on studying $\mu(n)$ have used neural networks or classifiers. Francois Charton made an integer sequence to integer sequence transformer-based translator, [Int2Int] (available at Int2Int), and I thought it would be fun to see if this works any different.

Initially, I sought to get Int2Int to work. I describe aspects of that and how to run it in various ways here.

I'm splitting my description into two parts: a general report [DLD-General] and a technical report. This is the technical report. This includes many details related to actually running and analyzing the code.

---

## 2. Int2Int

Int2Int can be found at [Int2Int](). It is possible to run Int2Int using a CPU, but it's much slower and probably not worth trying.

Francois and Edgar Costa (and to a lesser extent, me) have tried ot make Int2Int as self-contained as possible. It's certainly easier now than it was a few weeks ago — it might be possible for the Reader to experiment with https://github.com/f-charton/Int2Int/ using only the README there.

2.1. **Training from data files.** By default, Int2Int expects to be able to generate valid inputs and outputs on the fly. We wanted to use and experiment with data that is nontrivial to compute (such as the Möbius function or data associated to elliptic curves). For that, we've added the ability to train from data files.

> It would be fair to say that running Int2Int with default settings is so easy that the hardest part to get up and running is creating the data files. And this is only as hard as the data is to generate and store.

To train from data files, the file needs ot have a particular format. Recall that Int2Int fundamentall reads and outputs sequences of integers. Each integer is encoded as `s ad ... a0`, where `s` is either `+` or `-` and is the sign, and `ad` through `a0` are the digits in a given base (which defaults to 1000). For example, the number 12345 is encoded as `+ 12 345`.

An array of $n$ integers is encoded as `Vn z1 ... zn`, where the `n` in `Vn` is the actual number. For example, the array $(1, 1234, 1234567)$ is encoded as `V3 + 1 + 1 234 + 1 234 567`.

A datafile should have the input given as an array of the appropriate length, followed by a tab character `\t`, followed by the output. As an even more technical note, the output can be specified by a range of values instead of as an integer or integer array; this is useful with $\mu(n)$ since it can only take 3 values. This has to do with the symbol table that Int2Int uses, and the fact that it uses [cross-entropy loss]() to measure performance.

A complete datafile could be the following.

```
1  V5 + 1 + 2 + 1 + 3 + 4 + 5\t+ 1
2  V5 + 0 + 2 + 1 + 3 + 1 + 5\t+ 0
```

This data file has the spec `int[5]:int`. If we wanted more than a single `int` as output, we would have to use `Vn` appropriately. There is an additional datatype called `range` (with python-like semantics). In practice, if we know the output is a single constrained integer, there is a minor boost from using `range` instead of `int`.

2.2. **Datafile Generation Scripts.** I generated most of my datafiles using a script that closely looked like the following. This is a `sagemath` script. That is, it's mostly python, but it has inbuilt commands `primes` and `moebius` that I take for granted.

```python
1  primes_100 = list(primes(542))  # generate list of 100 primes
2
3  def encode_integer(val, base=1000, digit_sep=" "):
4      if val == 0:
5          return '+ 0'
6      sgn = '+' if val >= 0 else '-'
7      val = abs(val)
8      r = []
9      while val > 0:
10         r.append(str(val % base))
11         val = val//base
12     r.append(sgn)
13     r.reverse()
14     return digit_sep.join(r)
15
16 # Each line has an input, a tab, and an output.
17 def make_line(n):
18     return make_input(n) + "\t" + make_output(n) + "\n"
19
20 def make_input(n):
21     ret = []
22     count = len(primes_100)
23     ret.append(f"V{2*count}")
24     for p in primes_100:
25         ret.append(encode_integer(n % p))  # feed in n mod p
26         ret.append(encode_integer(p))       # followed by p
27     return ' '.join(ret)
28
29 def make_output(n):
30     return str(moebius(n))
```

What's left is to determine what family of $n$ to input. For performing regression-like tasks, Francois noted that using a log-distribution tends to work best. This is more like a classification task as the output is one of $0$, $-1$, or $1$. Thus I generically uniformly sampled integers $n$ up to some large bound like $10^{13}$ without repetition. To do this, I generate random integers in the range and check to make sure that I don't generate the same one twice.

```python
1  import random
2
3  seen = set()
4  with open("mu_modp_and_p.txt", "w", encoding="utf8") as
   ↪  outfile:
5      while len(seen) < 10**7:
6          n = random.randint(2, 10**13)
7          if n in seen:
```

```
8                  continue
9            seen.add(n)
10           outfile.write(make_line(n))
```

Note that this creates $10^7$ lines, each having approximately $200 \cdot 3 \sim 1000$ characters. The resulting file will be approximately 10GB. Adjust the parameters appropriately!

The slow part is computing $\mu(n)$ for random integers. Generating random numbers (including the $10^{-6}$ chance of hitting a previously seen number) and writing to the file is fast; computing $\mu(n)$ for a random 12 digit number can be slow-ish.

But in practice, the actual slow part is training the resulting ML model. I didn't work to optimize generation of $\mu$ at all.

I note that a sieve could generate all the Möbius values up to $N$ at once. Then you could sample from these values in whatever way makes sense. Something along the following lines would work (and would remove the sagemath dependency).

```python
1   def primes_up_to(X):
2       """
3       A basic implementation of Eratosthenes.
4       """
5       arr = [True] * (X + 1)
6       arr[0] = arr[1] = False
7       primes = []
8       for p in range(X + 1):
9           if arr[p]:  # is prime
10              primes.append(p)
11              for j in range(p*p, X + 1, p):
12                  arr[j] = False
13      return primes
14
15
16  def mobius_up_to(X):
17      "Eratosthenes-like"
18      arr = [1] * (X + 1)
19      arr[0] = 0
20      ps = primes_up_to(X)
21      for p in ps:
22          for j in range(p, X + 1, p):
23              arr[j] *= -1
24          for j in range(p*p, X + 1, p*p):
25              arr[j] = 0
26      return arr
```

2.3. **Making testing and training data.** I then make testing and training data.

```python
import os
def shuffle_and_create(fname, ntrain=1900000, ntest=100000):
    "Shuffle and create test and training files"
    if not fname.endswith(".txt"):
        raise ValueError("Incorrect filename assumption.")
    name = fname[:-4]  # remove ".txt"
    print("shuffling...")
    os.system(f"shuf {name}.txt > {name}.shuf.txt")
    print("making training data...")
    os.system(f"head -n {ntrain} {name}.shuf.txt >
        {name}.txt.train")
    print("making testing data...")
    os.system(f"tail -n {ntest} {name}.shuf.txt >
        {name}.txt.test")
    print("done!")
```

2.4. **Running Int2Int.** It's now time to actually run the code. It's necessary to have a python with pytorch installed (not surprisingly) and to have Int2Int somewhere. But a generic run would look like

```
python ../Int2Int/train.py
    --num_workers 0
    --dump_path ~/scratch
    --exp_name dld_mu_modp_and_p_sqfree
    --exp_id 1
    --train_data ./mu_modp_and_p_sqfree.txt.train
    --eval_data ./mu_modp_and_p_sqfree.txt.test
    --local_gpu 1
    --epoch_size 250000
    --operation data
    --data_types "int[200]:range(-1,2)"
    --optimizer "adam,lr=0.00025"
```

This was one of the commands I used when using $(n \bmod p, p)$ for the first 100 primes (giving 200 inputs total) and output just $\mu(n)$ (one int in a prescribed range).

Most of these are straightforward. The `--optimizer` command is deceptively useful, largely because changing the initial learning rate can have large impacts on the overall performance.

When run in this way, it's almost certain that you'll need to manually stop the experiment before it has a complete run. This is because the default number of epochs to train through is very large. In practice, it's a good idea to sometimes look at the outputs or parse the logs and to see how the behavior is going.

2.5. **Log parsing and graph creation.** For log parsing and graph creation, I used code largely written by someone else (maybe Edgar Costa). This is a

pile of code, but it's just parsing the pickled logs from a set of experiments. Log writing and parsing always takes piles of not-very-hard code. This is no exception.

The relevant code is in the appendix.

## 3. REPRESENTATION

I thought using residues mod several primes was a good strategy. Other experiments have shown[1] that the base in which numbers are expressed can be very important.

Something like base 100 or base 1000 would allow for almost immediate recognition that $\mu(n) = 0$ if $25 \mid n$ or if $4 \mid n$, as these congruence classes are fixed. I'm more interested in what other sorts of mathematical structures the machine can learn.

In this case, I represented each number in base 1000, but almost never needed to use any number larger than 1000 (as the 100th prime is 541). The Chinese remainder theorem shows that this allows representation for every integer up to approximately $10^{219.67}$. This is large enough to be interesting.

```python
# pure python - uses primes_up_to defined above
import math
from functools import reduce

primes_100 = primes_up_to(1000)[:100]
print(primes_100[-1])
# 541

modulus = reduce(lambda x, y: x*y, primes_100, 1)
print(modulus)
# [...enormous...]
print(math.log(modulus)/math.log(10))
# 219.67...
```

If $n < 10^{219.67}$, then $n$ is uniquely determined by its residues mod $p$ for the first 100 primes $p$.

## 4. GUESSING $\mu(n)$ FROM $n$ mod $p$ WITHOUT USING CRT

One of the questions that came up was the following **mathematical** (not programmatic) question.

> How would you guess whether $n$ is squarefree or not given $n$ mod $p$ for lots of primes $p$?

---

[1]See *Learning the greatest common divisor: explaining transformer predictions* by François Charton. François also extracted Int2Int from the models used in this paper, more or less. Thank you François.

One way would be to perform the Chinese remainder theorem, reconstruct $n$, and then actually check. There is no known polynomial-time algorithm to check if an integer is squarefree, so this approach is generically slow.

The "default" algorithm would be to note that about 60.79 percent of numbers are squarefree. So guessing `squarefree` all the time would be right just over 60 percent of the time. I want any algorithm that does better.

The Dirichlet series for squarefree numbers that are divisible by a fixed prime $q$ is

$$\frac{1}{q^s} \prod_{\substack{p \\ p \neq q}} \left(1 + \frac{1}{p^s}\right) = \frac{1}{q^s} \frac{(1 - 1/q^s)}{(1 - 1/q^{2s})} \frac{\zeta(s)}{\zeta(2s)}, \tag{1}$$

and the series for squarefree numbers that aren't divisible by a fixed prime $q$ is the same, but without $q^{-s}$. Thus the percentage of integers that are squarefree and divisible by $q$ or not divisible by $q$ are, respectively,

$$\frac{1}{q+1} \frac{6}{\pi^2} \quad \text{and} \quad \frac{q}{q+1} \frac{6}{\pi^2}. \tag{2}$$

A simple application of conditional probability shows that

$$P(\text{sqfree}|\text{q-even}) = \frac{P(\text{sqfree and q-even})}{P(\text{q-even})} = \frac{q}{q+1} \frac{6}{\pi^2}$$

$$P(\text{sqfree}|\text{q-odd}) = \frac{P(\text{sqfree and q-odd})}{P(\text{q-odd})} = \frac{q^2}{q^2 - 1} \frac{6}{\pi^2}.$$

I use the adhoc shorthand $q$-even to mean divisible by $q$, and $q$-odd to mean not divisible by $q$.

Let's quickly experimentally verify this. We make squarefree numbers with yet another Eratosthenes-type sieve.

```python
def squarefree_up_to(X):
    """
    Eratosthenes-like.
    """
    arr = [True] * (X + 1)
    arr[0] = False
    ps = primes_up_to(int(X**.5) + 1)
    for p in ps:
        for j in range(p*p, X + 1, p*p):
            arr[j] = False
    ret = []
    for i in range(X + 1):
        if arr[i]:
            ret.append(i)
    return ret


sfree = squarefree_up_to(10_000_000)
```

```
19  print(len(sfree)/10_000_000)
20  # 0.6079291
21
22  import math
23  print(6./math.pi**2)
24  # 0.6079271018540267
```

As an aside, I note that this converges very quickly. Look at how close that is! One useless application of the Riemann Hypothesis is that is would guarantee how quickly the density of the number of squarefree numbers up to $X$ would converge to $6/\pi^2$.

```
1   def ratio_sqfree_with(filterfunc):
2       return sum(1 for n in sfree if filterfunc(n))/len(sfree)
3
4   def is_even(x):
5       return 1 if x % 2 == 0 else 0
6   def is_odd(x):
7       return 1 if x % 2 == 1 else 0
8
9   # even and sqfree
10  ratio_sqfree_with(is_even)
11  # 0.3333309756022536
12
13  # odd and sqfree
14  ratio_sqfree_with(is_odd)
15  # 0.6666690243977463
16
17  def is_3even(x):
18      return 1 if x % 3 == 0 else 0
19  def is_3odd(x):
20      return not is_3even(x)
21
22  ratio_sqfree_with(is_3even)
23  # 0.24999839619455624
24
25  ratio_sqfree_with(is_3odd)
26  # 0.7500016038054438
```

This agrees with the claim above that $1/(q+1)$ of squarefree numbers are divisible by the prime $q$, and $q/(q+1)$ are not. The converse probabilities follow from basic probability, but to make sure:

```
1   sqfree_set = set(sfree)  # for quick inclusion checking
2
3   def prob_sqfree_given(filterfunc):
4       sqfree_count = 0
5       total_count = 0
```

```python
6        for n in (x for x in range(10_000_000) if filterfunc(x)):
7            total_count += 1
8            if n in sqfree_set:
9                sqfree_count += 1
10       if total_count == 0:
11           return 0.0
12       return sqfree_count / total_count
13
14   # P(sqfree | divis by 2)
15   prob_sqfree_given(is_even)
16   # 0.4052832
17
18   # P(sqfree | not divis by 2)
19   prob_sqfree_given(is_odd)
20   # 0.810575
21
22   # P(sqfree | divis by 3)
23   prob_sqfree_given(is_3even)
24   # 0.45594380881123825
25   3/4 * 6/math.pi**2
26   # 0.45594532639052
27
28   # P(sqfree | not divis by 3)
29   prob_sqfree_given(is_3odd)
30   # 0.6839217683921769
31   9/8 * 6/math.pi**2
32   # 0.68391798958578
```

These are very close to the theoretical computations above — again, it turns out that convergence is very quick.

4.1. **Compound Probabilities.** We'll compute joint probabilities theoretically in a moment. But we'll also experimentally find them.

Let's look at the probability using the small-prime strategy for the primes $2, 3, 5$: if $n$ is divisible by one of these, guess that $n$ is not squarefree; otherwise guess that $n$ is squarefree.

```python
1    def not_divis_by_small_prime(n):
2        for p in (2, 3, 5):
3            if n % p == 0:
4                return False
5        return True
6
7    A = prob_sqfree_given(not_divis_by_small_prime)
8    print(A)
9    # 0.9498902374725594
10
```

```python
11  def prob_notsqfree_given(filterfunc):
12      notsqfree_count = 0
13      total_count = 0
14      for n in (x for x in range(10_000_000) if filterfunc(x)):
15          total_count += 1
16          if n not in sqfree_set:
17              notsqfree_count += 1
18      if total_count == 0:
19          return 0
20      return notsqfree_count / total_count
21
22
23  def divis_by_small_prime(n):
24      return not not_divis_by_small_prime(n)
25
26  B = prob_notsqfree_given(divis_by_small_prime)
27  print(B)
28  # 0.5164203621436034
```

The density of numbers not divisible by 2, 3, or 5 is $(1-1/2)(1-1/3)(1-1/5) \approx 0.2666$. Thus 0.2666 of the time, $n$ isn't divisible by 2 or 3 or 5 and we would guess that $n$ is squarefree; this is correct about 0.9498 of the time. And the 0.7333 of the time when $n$ is divisible by at least one of 2 or 3 or 5, we guess that $n$ is not squarefree; this is correct 0.5164 of the time.

In total, we expect that this strategy is correct with density

$$0.2666 \cdot 0.9498 + 0.7333 \cdot 0.5164 \approx 0.6318.$$

Let's check:

```python
1  not_divis_prob = (1 - 1/2)*(1 - 1/3)*(1 - 1/5)
2  corr = not_divis_prob * A + (1 - not_divis_prob) * B
3  print(corr)
4  # 0.6320123288979917
```

If you look, you'll see that this does better than the naive guess (always guess squarefree) but is worse than guessing based only on mod 2 data. This is because we're ignoring all of the various cross-correlations. Clearly incorporating cross-correlations can never do worse than only using the mod 2 data.

Suppose we look instead at all the probabilities for all $2^\ell$ possibilities of $n$ being divisible or not by the first $\ell$ primes. Here, I use the first 4 primes, and the strategy is simple: compute whether it is more likely for $n$ to be squarefree or not given each divisibility pattern, and guess that one.

```python
1  def binary_to_prime_sets(n, length=4):
2      assert length <= 25
3      b = bin(n)[2:]
4      b = "0" * (length - len(b)) + b
5      is_divis = []
```

```
6        not_divis = []
7        ps = primes_up_to(100)[:length]
8        for l, p in zip(b, ps):
9            if l == "1":
10               is_divis.append(p)
11           else:
12               not_divis.append(p)
13       return is_divis, not_divis
14
15   def divis_rules(is_divis, not_divis):
16       def filterfunc(n):
17           for p in is_divis:
18               if n % p != 0:
19                   return False
20           for p in not_divis:
21               if n % p == 0:
22                   return False
23           return True
24       return filterfunc
25
26   def density_given(is_divis, not_divis):
27       filterfunc = divis_rules(is_divis, not_divis)
28       count = 0
29       for n in (x for x in range(10_000_000) if filterfunc(x)):
30           count += 1
31       return count/10_000_000
32
33   correct = 0
34   exp = 4
35   for n in range(2**exp):
36       is_divis, not_divis = binary_to_prime_sets(n, length=exp)
37       ff = divis_rules(is_divis, not_divis)
38       psqfree = prob_sqfree_given(ff)
39       density = density_given(is_divis, not_divis)
40       prob = max(psqfree, 1 - psqfree)
41       correct += density * prob
42       print(
43           is_divis, not_divis, n,
44           psqfree, density, prob, density * prob, correct
45       )  # my own diagnostics
46   print(correct)
47   # 0.7031860000000001
```

Remarkably this almost no better than just 2 alone! Before performing this computation, I had assumed that it would be notably better. Instead, it's close

enough that it might actually be the same as using 2 alone, combined with numerical imprecision.

With this set up, we can compute the theoretical probabilities instead of using experimentally determined probabilities.

4.2. **Actual computation.** Let $\{p_1, \ldots, p_N\}$ and $\{q_1, \ldots, q_D\}$ denote two disjoint sets of primes. We want to compute the density of squarefree numbers that are divisible by each of the $p_i$ and not divisible by any of the $q_j$. Each of these local conditions are independent; the overall density is the product of the local densities as described in~(1) and~(2). That is, the density of integers divisible by the $p_i$ and not divisible by the $q_j$ is

$$\prod_{p_i} \Big( \frac{1}{p_i + 1} \Big) \prod_{q_j} \Big( \frac{q_j}{q_j + 1} \Big) \frac{6}{\pi^2}.$$

Recall the chain rule from probability, that says

$$P\Big( \bigcap_{i=1}^{k} E_i \Big) = P\Big( E_1 | \bigcap_{i=2}^{k} E_i \Big) = P\Big( E_1 | \bigcap_{i=2}^{k} E_i \Big) P\Big( \bigcap_{i=2}^{k} E_i \Big),$$

(and which could chain further). I write $P(\text{sqfree}, p_1, p_2, \widehat{q_1}, \widehat{q_2})$ to mean the probability that a number is squarefree, divisible by $p_1$ and $p_2$, and not divisible by $q_1$ or $q_2$ (with obvious notational generalization). Then

$$P(\text{sqfree}|p_1, \ldots, p_N, \widehat{q_1}, \ldots, \widehat{q_D}) = \frac{P(\text{sqfree}, p_1, \ldots, p_N, \widehat{q_1}, \ldots, \widehat{q_D})}{P(p_1, \ldots, p_N, \widehat{q_1}, \ldots, \widehat{q_D})}.$$

Divisibility by different primes are independent, so this simplifies to

$$P(\text{sqfree}|p_1, \ldots, p_N, \widehat{q_1}, \ldots, \widehat{q_D}) = \frac{P(\text{sqfree}, p_1, \ldots, p_N, \widehat{q_1}, \ldots, \widehat{q_D})}{P(p_1) \cdots P(p_N) P(\widehat{q_1}) \cdots P(\widehat{q_D})}.$$

We also have that $P(p) = 1/p$ and $P(\widehat{q}) = (q-1)/q$.

Altogether, we compute that

$$P(\text{sqfree}|p_1, \ldots, p_N, \widehat{q_1}, \ldots, \widehat{q_D}) = \prod_{p_i} \Big( \frac{p_i}{p_i + 1} \Big) \prod_{q_j} \Big( \frac{q_j^2}{q_j^2 - 1} \Big) \frac{6}{\pi^2}.$$

Note that this generalizes the previous probabilities and is generically straightforward.

Let's quickly check by computing $P(\text{sqfree}|2, 3)$ and $P(\text{sqfree}|2, \widehat{3})$:

$$P(\text{sqfree}|2, 3) = \frac{1}{2} \frac{6}{\pi^2} \approx 0.3039,$$

$$P(\text{sqfree}|2, \widehat{3}) = \frac{3}{4} \frac{6}{\pi^2} \approx 0.4559.$$

```
1  divis_by_2_and_3 = divis_rules([2, 3], [])
2  print(prob_sqfree_given(divis_by_2_and_3))
3  # 0.30395813920837217
4
5  divis_by_2_not_3 = divis_rules([2], [3])
```

```
6   print(prob_sqfree_given(divis_by_2_not_3))
7   # 0.45594574559457457
```

Let's now compute the density of the following strategy being correct:

1. Fix a set of primes $P$.
2. For each partition of $P$ into two disjoint sets of primes $\{p_i\}$ and $\{q_j\}$:
3. Compute $P(\text{sqfree}|p_1,\ldots,p_N,\widehat{q_1},\ldots,\widehat{q_D})$.
4. For integers satisfying this set of prime divisibility rules, guess "square-free" if this probability is larger than 0.5; otherwise guess "not square-free".

```
1   def prob_sqfree_theoretical(ps, qs):
2       ret = 6/math.pi**2
3       for p in ps:
4           ret *= (p / (p + 1))
5       for q in qs:
6           ret *= (q*q/(q*q - 1))
7       return ret
8
9   prob_sqfree_theoretical([2], [3])
10  # 0.45594532639052
```

I assume we use the first $\ell$ primes and reuse some of the same logic as above.

```
1   def density_theoretical(ps, qs):
2       ret = 1
3       for p in ps:
4           ret *= (1/p)
5       for q in qs:
6           ret *= (q-1)/q
7       return ret
8
9   def strategy(ell):
10      correct = 0
11      for n in range(2**ell):
12          ps, qs = binary_to_prime_sets(n, length=ell)
13          psqfree = prob_sqfree_theoretical(ps, qs)
14          density = density_theoretical(ps, qs)
15          prob = max(psqfree, 1 - psqfree)
16          correct += density * prob
17      return correct
18
19  strategy(1)
20  # 0.7026423672846756
21
22  strategy(10)
23  # 0.7034137933079656
24
```

```
25  strategy(20)
26  # 0.7034211847385363
27
28  strategy(25)
29  # 0.7034221516869834
```

Computing this for anything much larger would be prohibitively computationally expensive. Without more sophisticated thinking, it seems we've hit a wall. Presumably this continues to grow, but perhaps is strictly bounded.

I pose this as an open question. It's not clear to me how hard it is to answer it.

**Question 1.** *What is the limiting behavior of this strategy? Can it be shown to be less than* 71 *percent?*

## 5. Extended remark on the utility of ML for pure mathematics

Machine learning operates as a block box. There may be many problems where it can achieve impressive results, but it might give no clues as to how it actually does its predictions.

But one place where it is useful is acting as a **one-sided oracle** to determine whether the inputs are enough to correctly evaluate an output. For example, if I wanted to determine if a particular set of inputs are sufficient to determine some behavior, I might try to feed these inputs along with known outcomes into a machine learning blackbox. If the ML soup acts with high accuracy using only those inputs, it seems more likely that those variables are indeed significant.

It's "one-sided" because the model might simply fail to model the function well. Failing to obtain high accuracy could reflect merely that the model wasn't strong enough, or there wasn't enough data, or any of a variety of points of failure that are independent of the underlying mathematical question.

And it's an "oracle" because there are no explanations for the insight. We can not ask for the workings behind the curtain.

With regard to the question above, I think I've tried such a variety of models and structures and learning rates that I suspect that knowing $n \bmod p$ for the first 100 primes isn't enough to guess $\mu(n)^2$ on random input $n$ more than 75 percent of the time. And this belief is bolstered by the failure of the ML to do better. (I've tried neural networks of various forms too, even though I haven't described those here).

But unfortunately that's not the direction the oracle sees and I don't believe in the strength of ML enough to make a conjecture or to draw a line in the sand.

### Appendix A. Parsing Logs

```python
1   # path and env name : THIS IS YOUR DUMP PATH
2   path = "~/scratch/"
3
4   # THE EXPERIMENTS YOU WANT TO PROBE AND THE ACCURACY INDICATOR
5   indicator = "valid_arithmetic"
6   xp_env=["dld_mu_modp_and_p_sqfree"]
7
8   # SET TO TRUE IF YOU USE BEAM SEARCH
9   has_beam=False
10
11  import os
12  import pickle
13  import matplotlib.pyplot as plt
14  import glob
15  import ast
16  from datetime import datetime
17  from tabulate import tabulate
18  import numpy as np
19  from operator import itemgetter
20
21  xp_id_filter=[]
22  xp_id_selector=[]
23  unwanted_args = ['dump_path']
24  var_args = set()
25  all_args = {}
26
27  # list experiments
28  xps = [(env, xp) for env in xp_env
29          for xp in os.listdir(path+'/'+env)
30          if (len(xp_id_selector)==0 or xp in xp_id_selector)
31          and (len(xp_id_filter)==0 or not xp in xp_id_filter)]
32  names = [path + env + '/' + xp for (env, xp) in xps]
33  print(len(names),"experiments found")
34
35  # read all args
36  pickled_xp = 0
37  for name in names:
38      pa = name+'/params.pkl'
39      if not os.path.exists(pa):
40          print("Unpickled experiment: ", name)
41          continue
42      pk = pickle.load(open(pa,'rb'))
43      all_args.update(pk.__dict__)
```

```python
44       pickled_xp += 1
45   print(pickled_xp, "pickled experiments found")
46   print()
47
48   # find variable args
49   for name in names:
50       pa = name+'/params.pkl'
51       if not os.path.exists(pa):
52           continue
53       pk = pickle.load(open(pa,'rb'))
54       for key,value in all_args.items():
55           if key in pk.__dict__ and value == pk.__dict__[key]:
56               continue
57           if key not in unwanted_args:
58               var_args.add(key)
59
60   print("common args")
61   for key in all_args:
62       if key not in unwanted_args and key not in var_args:
63           print(key,"=", all_args[key])
64   print()
65   print(len(var_args)," variables params out of", len(all_args))
66   print(var_args)
67
68   def vars_from_env_xp(env, xp):
69       res = {}
70       pa = path+env+'/'+xp+'/params.pkl'
71       if not os.path.exists(pa):
72           print("pickle", pa, "not found")
73           return res
74       pk = pickle.load(open(pa,'rb'))
75       for key in var_args:
76           if key in pk.__dict__:
77               res[key] = pk.__dict__[key]
78           else:
79               res[key] = None
80       return res
81
82   def get_start_time(line):
83       parsed_line = line.split(" ")
84       dt = datetime.strptime(parsed_line[2]+'
         ↪  '+parsed_line[3],"%m/%d/%y %H:%M:%S")
85       try:
86           idx = parsed_line.index("epoch")
87           curr_epoch = int(parsed_line[idx+1])
```

```python
88          except ValueError:
89              curr_epoch = ""
90          return dt, curr_epoch
91
92  def read_xp(env, xp, indics, max_epoch=None):
93      res = {"env":env, "xp": xp, "stderr":False, "log":False,
        ↪  "error":False}
94      stderr_file = os.path.join(os.path.expanduser("~"),
        ↪  'workdir/'+env+'/*/'+xp+'.stderr')
95      nb_stderr =len(glob.glob(stderr_file))
96      if nb_stderr > 1:
97          print("duplicate stderr", env, xp)
98          return res
99      for name in glob.glob(stderr_file):
100         with open(name, 'rt') as f:
101             res["stderr"]=True
102             errlines = []
103             cuda = False
104             terminated = False
105             forced = False
106             for line in f:
107                 if line.find("RuntimeError:") >= 0:
108                     errlines.append(line)
109                 if line.find("CUDA out of memory") >= 0:
110                     cuda = True
111                 if line.find("Exited with exit code 1") >=0:
112                     terminated = True
113
114                 if line.find("Force Terminated") >=0:
115                     forced = True
116             res["forced"] = forced
117
118             res["terminated"] = terminated
119             if len(errlines) > 0:
120                 res["error"] = True
121                 res["runtime_errors"] = errlines
122                 res["oom"] = cuda
123                 if not cuda:
124                     print(stderr_file,"runtime error no oom")
125
126     pa = path+env+'/'+xp+'/train.log'
127     if not os.path.exists(pa):
128         return res
129     res["log"] = True
130     with open(pa, 'rt') as f:
```

```
131            series = []
132            train_loss=[]
133            for ind in indics:
134                series.append([])
135            best_val = -1.0
136            best_xel = 999999999.0
137            best_epoch = -1
138            epoch = -1
139            val = -1
140            ended = False
141            nanfound = False
142            res["curr_epoch"]=-1
143            res["train_time"]=0
144            res["eval_time"]=0
145            res["pred_nr"]=[]
146            nb_sig10 = 0
147            nb_sig15 = 0
148            counter = 0
149            counting = False
150            for line in f:
151                try:
152                    if counting:
153                        counter += 1
154                    if line.find("Signal handler called with signal
      ↪    10") >= 0:
155                        nb_sig10 += 1
156                    if line.find("Signal handler called with signal
      ↪    15") >= 0:
157                        nb_sig15 += 1
158                    if line.find("Stopping criterion has been below
      ↪    its best value for more than") >=0:
159                        ended = True
160                    elif line.find("============ Starting epoch")
      ↪    >=0:
161                        dt, curr_epoch = get_start_time(line)
162                        if curr_epoch == max_epoch: break
163                        res["start_time"] = dt
164                        if curr_epoch >0 and curr_epoch ==
      ↪    res["curr_epoch"]+1:
165                            res["eval_time"] += (dt -
      ↪    res["end_time"]).total_seconds()
166                        res["curr_epoch"] = curr_epoch
167                    elif line.find("============ End of epoch")
      ↪    >=0:
```

```
168                     dt, curr_epoch = get_start_time(line)
169                     if curr_epoch != res["curr_epoch"]:
170                         print("epoch mismatch",
                          ↪  curr_epoch,"in", env,",", xp)
171                     else:
172                         res["end_time"] = dt
173                         res["train_time"] +=
                          ↪  (dt-res["start_time"]).total_seconds()
174                 elif line.find("- model LR:") >=0:
175                     loss = line.split(" ")[-5].strip()
176                     train_loss.append(None if loss == 'nan'
                      ↪  else float(loss))
177                 elif line.find("- LR:") >=0:
178                     loss = line.split(" ")[-4].strip()
179                     if loss == "predictions.":
180                         print(line)
181                     else:
182                         train_loss.append(None if loss == 'nan'
                          ↪  else float(loss))
183                 elif line.find('- test predicted pairs') >=0:
184                     counter = 0
185                     counting = True
186                 else:
187                     pos = line.find('__log__:')
188                     if pos >=0:
189                         counting = False
190                         res['pred_nr'].append(counter/100.0)
191                         if line[pos+8:].find(': NaN,') >= 0:
192                             nanfound = True
193                             line = line.replace(': NaN,',':
                              ↪  -1.0,')
194                         dic = ast.literal_eval(line[pos+8:])
195                         epoch = dic["epoch"]
196                         if not indicator+"_"+indics[0] in dic:
197                             continue
198                         if not indicator+"_"+indics[1] in dic:
199                             continue
200                         val = dic[indicator+"_"+indics[0]]
201                         xel = dic[indicator+"_"+indics[1]]
202                         if xel < best_xel:
203                             best_xel= xel
204                         if val > best_val:
205                             best_val = val
206                             best_epoch = epoch
```

```python
207                                res["best_dic"] = dic
208                            for i, indic in enumerate(indics):
209                                if indicator+"_"+indic in dic:
210
                                  ↪  series[i].append(dic[indicator+"_"+indic])
211
212              except Exception as e:
213                  print(e, "exception in", env, xp)
214                  continue
215              except:
216                  print(line)
217                  continue
218          res["nans"] = nanfound
219          res["ended"] = (ended or (nb_sig15 > nb_sig10))
220          res["last_epoch"] = epoch
221          res["last_acc"] = "{:.2f}".format(val)
222          res["best_epoch"] = best_epoch
223          res["best_acc"] = float("{:.2f}".format(best_val))
224          res["best_xeloss"] = "{:.2f}".format(best_xel)
225          res["train_loss"]=train_loss
226          res["avg_d"] = np.median(res['pred_nr'])
227          res["last_d"] = res['pred_nr'][-1] if
            ↪  len(res['pred_nr']) > 0 else -1
228          if epoch >=0:
229              res["train_time"] /= (epoch+1)
230              res["eval_time"] /= (epoch+1)
231          res["train_time"] = int(res["train_time"]+0.5)
232          res["eval_time"] = int(res["eval_time"]+0.5)
233
234          for i,indic in enumerate(indics):
235              res["last_"+indic] = "{:.2f}".format(series[i][-1])
                 ↪  if len(series[i])>0 else '0'
236              res["best_"+indic] =
                 ↪  "{:.2f}".format(max(series[i])) if
                 ↪  len(series[i])>0 else '0'
237              res[indic] = series[i]
238              if len(series[i])!= epoch + 1:
239                  print("mismatch in nr of epochs",env, xp,
                     ↪  epoch+1, len(series[i]), indic)
240      return res
241
242  data = []
243  indics = ["beam_acc" if has_beam is True else "acc","xe_loss"]
```

```
244  indics.extend(["correct", "perfect", "beam_acc_d1",
     ↪  "beam_acc_d2",
245  "beam_acc_nb", "additional_1","additional_2","additional_3"])
246
247  for (env, xp) in xps:
248      res = read_xp(env, xp, indics, None)  # USE THE LAST
         ↪  PARAMETER IF YOU WANT TO LIMIT READ TO N EPOCHS
249      res.update(vars_from_env_xp(env, xp))
250      data.append(res)
251
252  print(len(data), "experiments read")
253  print(len([d for d in data if d["stderr"] is False]),"stderr
     ↪  not found")
254  print(len([d for d in data if d["error"] is True]),"runtime
     ↪  errors")
255  print(len([d for d in data if "oom" in d and d["oom"] is
     ↪  True]),"oom errors")
256  print(len([d for d in data if "terminated" in d and
     ↪  d["terminated"] is True]),"exit code 1")
257  print(len([d for d in data if "forced" in d and d["forced"] is
     ↪  True]),"Force Terminated")
258  print(len([d for d in data if "last_epoch" in d and
     ↪  d["last_epoch"] >= 0]),"started XP")
259  print(len([d for d in data if "ended" in d and d["ended"] is
     ↪  True]),"ended XP")
260  print(len([d for d in data if "best_acc" in d and
     ↪  float(d["best_acc"]) > 0.0]),"began predicting")
```

And to make some graphs displaying various things, I would run the following. Or rather, I would run the above and below in a notebook, so the graphs display inline. (Otherwise I guess I would save them).

In practice, it was sufficient to look at the `tail` of the running log and to extract learning rate failures and accuracies on test sets.

```
1  import numpy as np
2
3  def compose(f,g):
4      return lambda x : f(g(x))
5
6  def print_table(data, args, sort=False):
7      res = []
8      for d in data:
9          line = [d[v] if v in d else None for v in args]
10         res.append(line)
11     if sort:
```

```python
12          res = sorted(res, key=compose(float,itemgetter(0)),
            ↪   reverse=True)
13      print(tabulate(res,headers=args,tablefmt="pretty"))
14
15  def speed_table(data, args, indic, sort=False, percent=95):
16      res = []
17      for d in data:
18          if indic in d:
19              line = [d[v] if v in d else None for v in args]
20              val= 10000
21              for i,v in enumerate(d[indic]):
22                  if v >= percent and i < val:
23                      val = i
24              line.insert(1,val)
25              res.append(line)
26      e= args.copy()
27      e.insert(1,'first epoch')
28      if sort:
29          res = sorted(res, key=compose(float,itemgetter(1)),
            ↪   reverse=False)
30      print(tabulate(res,headers=e,tablefmt="pretty"))
31
32  def training_curve(data, indic, beg=0, end=-1, maxval=None,
    ↪   minval=None, export_to=""):
33      print(indic)
34      for d in data:
35          if indic in d:
36              if end == -1:
37                  plt.plot(d[indic][beg:],linewidth=1)
38              else:
39                  plt.plot(d[indic][beg:end],linewidth=1)
40      plt.ylim(minval,maxval)
41      plt.rcParams['figure.figsize'] = [10,10]
42      if export_to != '':
43          # print(export_to)
44          plt.savefig(export_to,bbox_inches="tight")
45      plt.show()
46
47  def filter_xp(xp, filt):
48      for f in filt:
49          if not f in xp:
50              return False
51          if not xp[f] in filt[f]:
52              return False
```

```python
53        return True
54
55  def xp_stats(data, splits, best_arg, best_value):
56      res_dic = {}
57      nb = 0
58      for d in data:
59          if d[best_arg] < best_value: continue
60          nb += 1
61          for s in splits:
62              if not s in d: continue
63              lib=s+':'+str(d[s])
64              if lib in res_dic:
65                  res_dic[lib] += 1
66              else:
67                  res_dic[lib]=1
68      print()
69      print(f"{nb} experiments with accuracy over {best_value}")
70      for elem in sorted(res_dic):
71          print(elem,' : ',res_dic[elem])
72      print()
73
74  xp_filter ={}
75
76  # CHANGE THESE TO FILTER THE EXPERIMENTS
77  #xp_filter.update({"n_enc_layers":[4]})
78  #xp_filter.update({"enc_emb_dim":[512]})
79
80  fdata = [d for d in data if filter_xp(d, xp_filter) is True]
81
82  oomtab = [d for d in fdata if d["error"] is True]
83  print(f"CUDA out of memory ({len(oomtab)})")
84  print_table(oomtab, var_args)
85
86  forcetab = [d for d in fdata if 'forced' in d and d["forced"]
    ↪  is True]
87  print(f"Forced terminations ({len(forcetab)})")
88  print_table(forcetab, var_args)
89
90  unstartedtab = [d for d in fdata if "last_epoch" in d and
    ↪  d["last_epoch"] < 0]
91  print(f"Not started ({len(unstartedtab)})")
92  print_table(unstartedtab, var_args)
93
94  crypto = False
```

```python
95  runargs = ["best_acc", "best_epoch","best_xeloss",  "ended",
    ↪  "last_epoch",
96  "last_acc", "last_xe_loss","nans", "error", "train_time",
    ↪  "eval_time"]
97
98  #runargs.extend(["best_acc_d1" , "best_acc_d2"])
99  for v in var_args:
100     runargs.append(v)
101 runningtab = [d for d in fdata if "last_epoch" in d and
    ↪  d["last_epoch"] >= 0]
102 print(f"Running experiments ({len(runningtab)})")
103
104 #splits = ['n_enc_layers','dec_emb_dim','reload_size']
105 #xp_stats(fdata, splits, 'best_acc',90.0)
106 print()
107 print_table(runningtab, runargs, sort=True)
108
109 training_curve(fdata, "beam_acc" if has_beam is True else
    ↪  "acc",0, -1, None, export_to = "")
110 training_curve(fdata, "perfect")
111 training_curve(fdata, "correct")
112
113 training_curve(fdata, "xe_loss", 0) #, None, 0.9* np.min([x for
    ↪  d in fdata for x in d["xe_loss"] if x >0.0]))
114 training_curve(fdata, "train_loss",0, -1, 2)
115 speed_table(runningtab, runargs, "beam_acc" if has_beam else
    ↪  "acc", sort=True,percent=99)
116 speed_table(runningtab, runargs, "beam_acc" if has_beam else
    ↪  "acc", sort=True,percent=50)
117 speed_table(runningtab, runargs, "beam_acc" if has_beam else
    ↪  "acc", sort=True,percent=55)
118 speed_table(runningtab, runargs, "beam_acc" if has_beam else
    ↪  "acc", sort=True,percent=60)
```

## REFERENCES

[DLD-General]  David Lowry-Duda, *General Report on Machine Learning Experiments for the Möbius Function*. 2024 October 21. (Cited on page 1)

[Int2Int]  *Int2Int*  Github  Repository,  https://github.com/f-charton/Int2Int. Accesssed 2024 October 20.

(Cited on page 1)