

Modular Arithmetic

In several settings, such as error-correcting codes and cryptography, we sometimes wish to work over a smaller range of numbers. Modular arithmetic is useful in these settings, since it limits numbers to a predefined range $\{0, 1, \dots, N - 1\}$, and wraps around whenever you try to leave this range — like the hand of a clock (where $N = 12$) or the days of the week (where $N = 7$).

Example: Calculating the time When you calculate the time, you automatically use modular arithmetic. For example, if you are asked what time it will be 13 hours from 1 o'clock, you say 2 o'clock rather than 14. Let's assume our clock displays 12 as 0. This is limiting numbers to a predefined range, $\{0, 1, 2, \dots, 11\}$. Whenever you add two numbers in this setting, you divide by 12 and provide the remainder as the answer.

If we wanted to know what the time would be 24 hours from 2 o'clock, the answer is easy. It would be 2 o'clock. This is true not just for 24 hours, but for any multiple of 12 hours. What about 25 hours from 2? Since the time 24 hours from 2 is still 2, 25 hours later it would be 3. Another way to say this is that we add 1 hour, which is the remainder when we divide 25 by 12.

This example shows that under certain circumstances it makes sense to do arithmetic within the confines of a particular number (12 in this example). That is, we only keep track of the remainder when we divide by 12, and when we need to add two numbers, instead we just add the remainders. This method is quite efficient in the sense of keeping intermediate values as small as possible, and we shall see in later notes how useful it can be.

More generally we can define $x \bmod m$ (in words x modulo m) to be the remainder r when we divide x by m . i.e. if $x \bmod m = r$, then $x = mq + r$ where $0 \leq r \leq m - 1$ and q is an integer. Thus $29 \bmod 12 = 5$ and $13 \bmod 5 = 3$.

Computation

If we wish to calculate $(x + y) \bmod m$, we would first add $x + y$ and then calculate the remainder when we divide the result by m . For example, if $x = 14$ and $y = 25$ and $m = 12$, we would compute the remainder when we divide $x + y = 14 + 25 = 39$ by 12, to get the answer 3. Notice that we would get the same answer if we first computed $2 = x \bmod 12$ and $1 = y \bmod 12$ and added the results modulo 12 to get 3. The same holds for subtraction: $(x - y) \bmod 12$ is $-11 \bmod 12$, which is 1. Again, we could have directly obtained this as $2 - 1$ by first simplifying $x \bmod 12$ and $y \bmod 12$.

This is even more convenient if we are trying to multiply: to compute $xy \bmod 12$, we could first compute $xy = 14 \times 25 = 350$ and then compute the remainder when we divide by 12, which is 2. Notice that we get the same answer if we first compute $2 = x \bmod 12$ and $1 = y \bmod 12$ and simply multiply the results modulo 12.

More generally, while carrying out any sequence of additions, subtractions or multiplications $\bmod m$, we get the same answer even if we reduce any intermediate results $\bmod m$. This can considerably simplify the calculations.

Set Representation

There is an alternate view of modular arithmetic which helps understand all this better. For any integer m we say that x and y are *congruent modulo m* if they differ by a multiple of m , or in symbols,

$$x \equiv y \pmod{m} \Leftrightarrow m \text{ divides } (x - y).$$

For example, 29 and 5 are congruent modulo 12 because 12 divides $29 - 5$. We can also write $22 \equiv -2 \pmod{12}$. Notice that x and y are congruent modulo m iff they have the same remainder modulo m .

What is the set of numbers that are congruent to 0 mod 12? These are all the multiples of 12: $\{\dots, -36, -24, -12, 0, 12, 24, 36, \dots\}$. What about the set of numbers that are congruent to 1 mod 12? These are all the numbers that give a remainder 1 when divided by 12: $\{\dots, -35, -23, -11, 1, 13, 25, 37, \dots\}$. Similarly the set of numbers congruent to 2 mod 12 is $\{\dots, -34, -22, -10, 2, 14, 26, 38, \dots\}$. Notice in this way we get 12 such sets of integers, and every integer belongs to one and only one of these sets.

In general if we work modulo m , then we get m such disjoint sets whose union is the set of all integers. We can think of each set as represented by the unique element it contains in the range $(0, \dots, m - 1)$. The set represented by element i would be all numbers z such that $z = mx + i$ for some integer x . Observe that all of these numbers have remainder i when divided by m ; they are therefore congruent modulo m .

We can understand the operations of addition, subtraction and multiplication in terms of these sets. When we add two numbers, say $x \equiv 2 \pmod{12}$ and $y \equiv 1 \pmod{12}$, it does not matter which x and y we pick from the two sets, since the result is always an element of the set that contains 3. The same is true about subtraction and multiplication. It should now be clear that the elements of each set are interchangeable when computing modulo m , and this is why we can reduce any intermediate results modulo m .

Here is a more formal way of stating this observation:

Theorem 5.1: If $a \equiv c \pmod{m}$ and $b \equiv d \pmod{m}$, then $a + b \equiv c + d \pmod{m}$ and $a \cdot b \equiv c \cdot d \pmod{m}$.

Proof: We know that there are integers k and ℓ such that $c = a + k \cdot m$ and $d = b + \ell \cdot m$. So $c + d = a + k \cdot m + b + \ell \cdot m = a + b + (k + \ell) \cdot m$, which means that $a + b \equiv c + d \pmod{m}$. The proof for multiplication is similar and left as an exercise. ■

What this theorem tells us is that we can always reduce any arithmetic expression modulo m into a natural number smaller than m . As an example, consider the expression $(13 + 11) \cdot 18 \pmod{7}$. Using the above Theorem several times we can write:

$$\begin{aligned} (13 + 11) \cdot 18 &\equiv (6 + 4) \cdot 4 \pmod{7} \\ &\equiv 10 \cdot 4 \pmod{7} \\ &\equiv 3 \cdot 4 \pmod{7} \\ &\equiv 12 \pmod{7} \\ &\equiv 5 \pmod{7}. \end{aligned}$$

In summary, we can always do basic arithmetic (multiplication, addition, subtraction, and division) calculations modulo m by reducing intermediate results modulo m .

Exponentiation

Another standard operation in arithmetic algorithms (this is used heavily in primality testing and RSA) is raising one number to a power modulo another number. How do we compute $x^y \pmod{m}$, where x, y, m

are natural numbers and $m > 0$? A naïve approach would be to compute the sequence $x \bmod m, x^2 \bmod m, x^3 \bmod m, \dots$ up to y terms, but this requires time exponential in the number of bits in y . We can do much better using the trick of *repeated squaring*:

```

algorithm mod-exp(x, y, m)
  if y = 0 then return(1)
  else
    z = mod-exp(x, y div 2, m)
    if y mod 2 = 0 then return(z * z mod m)
    else return(x * z * z mod m)

```

(In the algorithm, $y \text{ div } 2$ means the largest integer which isn't bigger than $\frac{y}{2}$, in other words $\frac{y}{2}$ “rounded down”. This is often written as $\lfloor \frac{y}{2} \rfloor$.) The algorithm uses the fact that any $y > 0$ can be written as $y = 2a$ or $y = 2a + 1$, where a is $\lfloor \frac{y}{2} \rfloor$, plus the facts

$$x^{2a} = (x^a)^2; \quad \text{and}$$

$$x^{2a+1} = x \cdot (x^a)^2.$$

As a useful exercise, you should use these facts to construct a formal inductive argument that the algorithm always returns the correct value.

What is its running time? The main task here, as is usual for recursive algorithms, is to figure out how many recursive calls are made. But we can see that the second argument, y , is being (integer) divided by 2 in each call, so the number of recursive calls is exactly equal to the number of bits, n , in y . (The same is true, up to a small constant factor, if we let n be the number of decimal digits in y .) Thus, if we charge only constant time for each arithmetic operation (`div`, `mod` etc.) then the running time of `mod-exp` is $O(n)$ ¹.

In a more realistic model (where we count the cost of operations at the bit level), we would need to look more carefully at the cost of each recursive call. Note first that the test on y in the `if`-statement just involves looking at the least significant bit of y , and the computation of $\lfloor \frac{y}{2} \rfloor$ is just a shift in the bit representation. Hence each of these operations takes only constant time. The cost of each recursive call is therefore dominated by the `mod` operation² in the final result. A fuller analysis of such algorithms is performed in CS170.

Inverses

We have so far discussed addition, multiplication and exponentiation. Subtraction is the inverse of addition and just requires us to notice that subtracting b modulo m is the same as adding $-b \equiv m - b \pmod{m}$.

What about division? This is a bit harder³. Over the reals dividing by a number x is the same as multiplying by $y = 1/x$. Here y is that number such that $x \cdot y = 1$. Of course we have to be careful when $x = 0$, since such a y does not exist. Similarly, when we wish to divide by $x \pmod{m}$, we need to find $y \pmod{m}$ such that $x \cdot y \equiv 1 \pmod{m}$; then dividing by x modulo m will be the same as multiplying by y modulo m . Such a y is called the *multiplicative inverse* of x modulo m . In our present setting of modular arithmetic, can we be sure that x has an inverse mod m , and if so, is it unique (modulo m) and can we compute it?

¹The notation $O(n)$ will be explained in discussion section. You can search the web for “big O notation” to find out more.

²You can analyze grade-school long-division for binary numbers to understand how long a `mod` operation would take.

³With the real numbers, inverting exponentiation uses logarithms. In modular arithmetic, the equivalent function is called the “discrete logarithm”. It is impossible to compute efficiently far as anyone knows, so we will not be talking about it.

As a first example, take $x = 8$ and $m = 15$. Then $2x = 16 \equiv 1 \pmod{15}$, so 2 is a multiplicative inverse of 8 (mod 15). As a second example, take $x = 12$ and $m = 15$. Then the sequence $\{ax \pmod{m} : a = 1, 2, 3, \dots\}$ is periodic, and takes on the values $\{12, 9, 6, 3, 0, \dots\}$ (check this!). Thus 12 has no multiplicative inverse mod 15 since the number 1 never appears in that sequence.

This is the first warning sign that working in modulo arithmetic might actually be a bit different than grade-school arithmetic. Two weird things are happening. First, no multiplicative inverse seems to exist for a number that isn't zero. In normal arithmetic, the only thing you have to worry about is dividing by zero. Second, the "times table" for a number that isn't zero has zero showing up in it. So 12 times 5 is equal to zero when we are considering numbers modulo 15. For grade-school arithmetic, zero never shows up in the multiplication table for any number other than zero.

So when *does* x have a multiplicative inverse modulo m ? The answer is: iff the greatest common divisor of m and x is 1. Moreover, when the inverse exists it is unique. Recall that the *greatest common divisor* of two natural numbers x and y , denoted $\gcd(x, y)$, is the largest natural number that divides them both. For example, $\gcd(30, 24) = 6$. If $\gcd(x, y)$ is 1, it means that x and y share no common factors (except 1). This is often expressed by saying that x and m are *relatively prime*.

Theorem 5.2: Let m, x be positive integers such that $\gcd(m, x) = 1$. Then x has a multiplicative inverse modulo m , and it is unique (modulo m).

Proof: Consider the sequence of m numbers $0, x, 2x, \dots, (m-1)x$. We claim that these are all distinct modulo m . Since there are only m distinct values modulo m , it must then be the case that $ax \equiv 1 \pmod{m}$ for exactly one a (modulo m). This a is the unique multiplicative inverse.

To verify the above claim, suppose that $ax \equiv bx \pmod{m}$ for two distinct values a, b in the range $0 \leq a, b \leq m-1$. Then we would have $(a-b)x \equiv 0 \pmod{m}$, or equivalently, $(a-b)x = km$ for some integer k (possibly zero or negative).

However, x and m are relatively prime, so x cannot share any factors with m . This implies that $a-b$ must be an integer multiple of m . This is not possible, since $a-b$ ranges between 1 and $m-1$. ■

Actually it turns out that $\gcd(m, x) = 1$ is also a *necessary* condition for the existence of an inverse: i.e., if $\gcd(m, x) > 1$ then x has no multiplicative inverse modulo m . You might like to try to prove this using a similar idea to that in the above proof. (*Hint: Think about when zeros show up in multiplication tables.*)

Since we know that multiplicative inverses are unique when $\gcd(m, x) = 1$, we shall write the inverse of x as $x^{-1} \pmod{m}$. Being able to compute the multiplicative inverse of a number is crucial to many applications, so ideally the algorithm used should be efficient. It turns out that we can use an extended version of Euclid's algorithm, which computes the gcd of two numbers, to compute the multiplicative inverse.

Computing the Multiplicative Inverse

Let us first discuss how computing the multiplicative inverse of x modulo m is related to finding $\gcd(x, m)$. For any pair of numbers x, y , suppose we could not only compute $\gcd(x, y)$, but also find integers a, b such that

$$d = \gcd(x, y) = ax + by. \tag{1}$$

(Note that this is not a modular equation; and the integers a, b could be zero or negative.) For example, we can write $1 = \gcd(35, 12) = -1 \cdot 35 + 3 \cdot 12$, so here $a = -1$ and $b = 3$ are possible values for a, b .

If we could do this then we'd be able to compute inverses, as follows. We first find the integers a and b such that

$$1 = \gcd(m, x) = am + bx.$$

But this means that $bx \equiv 1 \pmod{m}$, so b is the multiplicative inverse of x modulo m . Reducing b modulo m gives us the unique inverse we are looking for. In the above example, we see that 3 is the multiplicative inverse of 12 (mod 35). So, we have reduced the problem of computing inverses to that of finding integers a, b that satisfy equation (1). Remarkably, Euclid's algorithm for computing gcd's also allows us to find the integers a and b described above. So computing the multiplicative inverse of x modulo m is as simple as running Euclid's gcd algorithm on input x and m !

Euclid's Algorithm for computing the GCD

If we wish to compute the gcd of two numbers x and y , how would we proceed? If x or y is 0, then computing the gcd is easy; it is simply the other number, since 0 is divisible by everything (although of course it divides nothing). The algorithm for other cases is ancient, and although associated with the name of Euclid, is almost certainly a folk algorithm invented by craftsmen (the engineers of their day) because of its intensely practical nature⁴. This algorithm exists in cultures throughout the globe.

The algorithm for computing $\text{gcd}(x, y)$ uses the following theorem to eventually reduce to the case where one of the numbers is 0:

Theorem 5.3: Let $x \geq y$ and let q, r be natural numbers such $x = yq + r$ and $r < y$. Then $\text{gcd}(x, y) = \text{gcd}(r, y)$.

Proof: This is because any common divisor of x and y is also a common divisor of y and r and vice versa. To see this, if d divides both x and y , there exist integers z and z' such that $zd = x$ and $z'd = y$. Therefore $r = x - yq = zd - z'dq = (z - z'q)d$, and so d divides r . The other direction follows in exactly the same way. ■

Given this theorem, let's see how to compute $\text{gcd}(16, 10)$:

$$\begin{aligned} 16 &= 10 \times 1 + 6 \\ 10 &= 6 \times 1 + 4 \\ 6 &= 4 \times 1 + 2 \\ 4 &= 2 \times 2 + 0 \\ 2 &= 0 \times 0 + 2 \end{aligned}$$

In each line, we write the larger number x as $yq + r$, where y is the smaller number. The next line then replaces the larger number with y , and the smaller number with r . This preserves the gcd, as shown in the theorem above. Therefore, $\text{gcd}(16, 10) = \text{gcd}(2, 0) = 2$. Or if you wish you can stop a step earlier and say that the gcd is the last non-zero remainder: i.e. you can stop at the step $6 = 4 \times 1 + 2$, since at the next step the remainder is 0.

This algorithm can be written recursively as follows:

```
algorithm gcd(x, y)
  if y = 0 then return(x)
  else return(gcd(y, x mod y))
```

Note: This algorithm assumes that $x \geq y \geq 0$ and $x > 0$.

⁴This algorithm is used for figuring out a common unit of measurement for two lengths. You can imagine how this is extremely important for building something up from a scale model. Different lengths in a design can be expressed as integer multiples of a common length, and then a new measuring stick can be found for the scaled-up design. We will see how the algorithm itself can be executed without literacy or symbolic notation. It is fundamentally *physical* in its intuition and you should figure out how this can be executed using threads. In the homework, you will see how this algorithm reveals the secret hidden in plain sight within the Pentagram.

Let's go through a quick example of this recursive implementation of Euclid's algorithm. We wish to compute $\text{gcd}(32, 10)$:

$$\begin{aligned}\text{gcd}(32, 10) &= \text{gcd}(10, 2) \\ &= \text{gcd}(2, 0) \\ &= 2\end{aligned}$$

Theorem 5.4: The algorithm above correctly computes the gcd of x and y .

Proof: Correctness is proved by (strong) induction on y , the smaller of the two input numbers. For each $y \geq 0$, let $P(y)$ denote the proposition that the algorithm correctly computes $\text{gcd}(x, y)$ for all values of x such that $x \geq y$ (and $x > 0$). Certainly $P(0)$ holds, since $\text{gcd}(x, 0) = x$ and the algorithm correctly computes this in the `if`-clause. For the inductive step, we may assume that $P(z)$ holds for all $z < y$ (the inductive hypothesis); our task is to prove $P(y)$. The key observation here is that $\text{gcd}(x, y) = \text{gcd}(y, x \bmod y)$ — that is, replacing x by $x \bmod y$ does not change the gcd. This is because a divisor d of y also divides x if and only if it divides $x \bmod y$ (divisibility by d is not affected by adding or subtracting multiples of d , and y is a multiple of d). Hence the `else`-clause of the algorithm will return the correct value provided the recursive call `gcd(y, x mod y)` correctly computes the value $\text{gcd}(y, x \bmod y)$. But since $x \bmod y < y$, we know this is true by the inductive hypothesis. This completes our verification of $P(y)$, and hence the induction proof. ■

How long does this algorithm take? In terms of arithmetic operations on integers, it takes time $O(n)$, where n is the total number of bits in the input (x, y) .

You should be able to see the intuitive connection to exponentiation-by-repeated-squaring. It is obvious that the arguments of the recursive calls become smaller and smaller (because $y \leq x$ and $x \bmod y < y$). The question is, how fast?

We shall show that, in the computation of $\text{gcd}(x, y)$, after two recursive calls the first (larger) argument is smaller than x by at least a factor of two (assuming $x > 0$). There are two cases:

1. $y \leq \frac{x}{2}$. Then the first argument in the next recursive call, y , is already smaller than x by a factor of 2, and thus in the next recursive call it will be even smaller.
2. $x \geq y > \frac{x}{2}$. Then in two recursive calls the first argument will be $x \bmod y$, which is smaller than $\frac{x}{2}$.

So, in both cases the first argument decreases by a factor of at least two every two recursive calls. Thus after at most $2n$ recursive calls, where n is the number of bits in x , the recursion will stop (note that the first argument is always a natural number).

Note that the above argument only shows that the *number of recursive calls* in the computation is $O(n)$. We can make the same claim for the running time if we assume that each call only requires constant time. Since each call involves one integer comparison and one mod operation, it is reasonable to claim that its running time is constant. In a more realistic model of computation, however, we should really make the time for these operations depend on the size of the numbers involved. This will be discussed in CS170.

Extended Euclid's Algorithm

In order to compute the multiplicative inverse, we need an algorithm which also returns integers a and b such that:

$$\text{gcd}(x, y) = ax + by.$$

Now since this problem is a generalization of the basic gcd, it is perhaps not too surprising that we can solve it with a fairly straightforward extension of Euclid's algorithm.

Examples

Let's first see how we would compute such numbers for $x = 6$ and $y = 4$. We'll need the equations from our example above, copied here for reference:

$$\begin{aligned}16 &= 10 \times 1 + 6 \\10 &= 6 \times 1 + 4 \\6 &= 4 \times 1 + 2 \\4 &= 2 \times 2 + 0\end{aligned}$$

From the last two equations it follows that $\gcd(6,4) = 2$. But now the second last equation gives us the numbers a, b , since we just rearrange that equation to say $2 = 6 \times 1 - 4 \times 1$. So $a = 1$ and $b = -1$.

What if we started with $x = 10$ and $y = 6$? Now we would write the last three equations to determine that $\gcd(10,6) = 2$. But how do we find a, b ? Start as above and write $2 = 6 \times 1 - 4 \times 1$. But we want 10 and 6 on the right hand side, not 6 and 4. But notice that the third from the last equation allows us to write 4 as a linear combination of 6 and 10 and so we can just back substitute: we rewrite that equation as $4 = 10 \times 1 - 6 \times 1$ and substitute to get:

$$2 = 6 \times 1 - 4 \times 1 = 6 \times 1 - (10 \times 1 - 6 \times 1) = 6 \times 2 - 10 \times 1.$$

If we started with $x = 16$ and $y = 10$ we would back substitute again using the first equation rewritten as $6 = 16 - 10$ to get:

$$2 = 6 \times 2 - 10 \times 1 = (16 - 10) \times 2 - 10 = 16 \times 2 - 10 \times 3. \text{ So } a = 2 \text{ and } b = -3.$$

Algorithm

The following recursive algorithm implements the idea used in the examples above. It takes as input a pair of natural numbers $x \geq y$ as in Euclid's algorithm, and returns a triple of integers (d, a, b) such that $d = \gcd(x, y)$ and $d = ax + by$:

```
algorithm extended-gcd(x, y)
  if y = 0 then return(x, 1, 0)
  else
    (d, a, b) := extended-gcd(y, x mod y)
    return((d, b, a - (x div y) * b))
```

Note that this algorithm has the same form as the basic gcd algorithm we saw earlier; the only difference is that we now carry around in addition the required values a, b . You should hand-turn the algorithm on the input $(x, y) = (16, 10)$ from our earlier example, and check that it delivers correct values for a, b .

Let's now look at why the algorithm works. We just need to generalize the back substitution method we used in the example above.

In the base case ($y = 0$), we return the gcd value $d = x$ as before, together with values $a = 1$ and $b = 0$ which satisfy $ax + by = d$. If $y > 0$, we first recursively compute values (d, a, b) such that $d = \gcd(y, x \bmod y)$ and

$$d = ay + b(x \bmod y). \tag{2}$$

Just as in our analysis of the vanilla GCD algorithm, we know that this d will be equal to $\gcd(x, y)$. So the first component of the triple returned by the algorithm is correct.

What about the other two components? We need to update these values of a and b , say to A and B .

What should their values be? Well, from the specification of the algorithm, they must be integers that satisfy

$$d = Ax + By. \tag{3}$$

To figure out what A and B should be, we need to rearrange equation (2), as follows:

$$\begin{aligned} d &= ay + b(x \bmod y) \\ &= ay + b(x - \lfloor x/y \rfloor y) \\ &= bx + (a - \lfloor x/y \rfloor b)y. \end{aligned}$$

(In the second line here, we have used the fact that $x \bmod y = x - \lfloor x/y \rfloor y$ — check this!) Comparing this last equation with equation (3), we see that we need to take $A = b$ and $B = a - \lfloor x/y \rfloor b$. This is exactly what the algorithm does, and this is why the algorithm works. The ideas here can be made more formal to get a full proof of correctness.

Since the extended gcd algorithm has exactly the same recursive structure as the vanilla version, its running time will be the same up to constant factors (reflecting the increased time per recursive call). So once again the running time on n -bit numbers will be $O(n)$ arithmetic operations. This means that we can find multiplicative inverses efficiently.