

First assignment (2021/2022, this is the one you need to do, due 09-01-2022)

Merkle-Hellman cryptosystem (part 1)

In 1978, Merkle and Hellman proposed an interesting way to cipher information, based on an NP-complete problem (the knapsack problem). It turned out that this cipher is insecure. The techniques used to break it are outside the scope of this course. Below, we describe only one simplification of the cipher, where the knapsack problem is replaced by a subset-sum problem. (Hint: one of the techniques used to solve problem 2 of the homework of lesson P.02 may be useful here). In this assignment you will attempt to solve small instances of this problem.

The main idea of the Merkle-Hellman cryptosystem is to construct a subset-sum problem that is trivial to solve using a greedy algorithm, kept secret, and then to disguise it as a hard to solve problem, that is published. It is thus a public-key cryptosystem. The following information is kept secret:

- A super-increasing sequence $W = (w_1, w_2, \dots, w_n)$ of n positive integers. In a super-increasing sequence we have $w_k > \sum_{i=1}^{k-1} w_i$ for $2 \leq k \leq n$.
- A modulus m such that $m > \sum_{i=1}^n w_i$.
- A scrambling integer a such that $\gcd(a, m) = 1$.

Let $W' = (w'_1, w'_2, \dots, w'_n)$, where $w'_i = (aw_i) \bmod m$, for $1 \leq i \leq n$. Recall that $a \bmod b$ denotes the remainder of the division of a by b . The following information is published:

- The sequence $P = (p_1, \dots, p_n)$, which is the sequence W' sorted in increasing order.

To cipher the n bits $B = (b_1, \dots, b_n)$ compute and send $C = \sum_{k=1}^n b_k p_k$. Knowing W , m , and a , it is easy to decode what was sent. Without that information, it is necessary to solve the hard subset-sum problem (see example in the next page).

Merkle-Hellman cryptosystem (part 2, example)

The following example show the Merkle-Hellman cryptosystem in action.

- Secret data: $W = (1, 3, 5, 12, 22, 47)$, $m = 100$, and $a = 13$.
- Intermediate data: $W' = (13, 39, 65, 56, 86, 11)$.
- Published data: $P = (11, 13, 39, 56, 65, 86)$. Note that 11 is associated to w_6 , 13 to w_1 , 39 to w_2 , 56 to w_4 , 65 to w_3 , and 86 to w_5 .
- Unencrypted message: $B = (b_1, \dots, b_6) = (0, 0, 1, 1, 0, 1)$.
- Encrypted message: $C' = 0 \times 11 + 0 \times 13 + 1 \times 39 + 1 \times 56 + 0 \times 65 + 1 \times 86 = 181$.

Knowing the secret information it is easy to decipher the message:

- Compute $C = (C'a^{-1}) \bmod m$. Since $(w'_i a^{-1}) \bmod m = w_i$, this transforms the hard problem into the easy problem. Because $13^{-1} \bmod 100 = 77$ — this is easy to check: $(13 \times 77) \bmod 100 = 1$ — we get $C = (181 \times 77) \bmod 100 = 37$.
- Now the problem is easy: which elements of W sum to 37? 47 is clearly out, but 22 must be in. $37 - 22 = 15$, so 12 must also be in. $15 - 12 = 3$, so 3 is also in. So, $37 = 22 + 12 + 3$.
- To finish decoding, observe that $22 = w_5$ is associated to $86 = p_6$, $12 = w_4$ is associated to $56 = p_4$, and $3 = w_2$ is associated to $39 = p_3$, and so the non-zero bits are b_3 , b_4 , and b_6 .

Of course, in your written assignment you will not be able to do this last part, because you will not be given the secret information. You will have to solve the hard subset-sum problem.

Merkle-Hellman cryptosystem (part 3, what is given)

In the archive A01.tgz you will find the problem instances you need to solve. Each student will have her/his own problem instances, stored in a .h file with the student's number. For each value of n , $10 \leq n \leq 57$, the .h file will contain a p sequence and 20 sums obtained by summing some terms of that sequence. The goal is to find out, for each sum, which terms were used to get the sum.

Please place your code in the subset_sum.c file. Each group of students should solve as many problems as they can, using the .h files corresponding to their student numbers — three students, three .h files! You can, and should, **delete** the .h files that are of no interest to you (the files for the students not in your group).

Things have been set up to make this easy. For example, for a student with number 100000, the program can be compiled using

```
cc -Wall -O2 -DSTUDENT_H_FILE=\"100000.h\" subset_sum_problem.c
```

so there is no need to hardwire the .h file name in the C file. As an alternative, you can edit the makefile, replacing the .h file name in the line

```
STUDENT_H_FILE=000000.h
```

There exists a special .h file, named 000000.h. It contains a full set of problems. The solutions to these problems are stored in comments after each sum; in each comment, the bits appear in the order b_1, b_2, \dots, b_n . Use these solutions to check the correctness of your own program.

If the Schroepel and Shamir technique is successfully implemented by a group of students (see next slide), it becomes feasible to break the 64-bit barrier, and to solve problem instances with somewhat larger values of n (up to 64). The program will have to deal with integers of type unsigned __int128. For those students, larger sets of problems instances will be available upon request.

Merkle-Hellman cryptosystem (part 4, what has to be done)

Each group should try the following techniques to solve the subset-sum problems:

- Brute force, using a recursive function.
- Clever brute force (called branch and bound), using a recursive function that avoids further recursions when it can be proved that with the current data no solution is possible because the partial sum is either too small or too large.
- Horowitz and Sahni technique (meet-in-the-middle technique). In this technique the p data is split in two nearly equal parts, all possible subset sums of each of the two parts are stored in sorted order in two arrays, and, in a single pass over these sorted arrays, one looks for the solution. In particular,
 1. Let a be the first sorted array, with size n_a , and let b be the second sorted array, with size n_b .
 2. Let i be the a -array index and let j be the b -array index; i starts at 0 and j starts at $n_b - 1$.
 3. The goal is to find i and j such that $a_i + b_j = s$, where s is the desired sum. In each loop iteration of the algorithm either the solution is found or either i is increased if $a_i + b_j < s$ — that is the only way of increasing the sum — or j is decreased if $a_i + b_j > s$ — that is the only way of decreasing the sum.
 4. If one of the indices goes out of bounds, no solution exists.
- Schroepel and Shamir technique (meet-in-the-middle using less memory). This is similar to the Horowitz and Sahni technique but the two arrays a and b are not kept in memory. Instead, the elements of each one are generated on-the-fly using a min-heap (for the a array) and a max-heap (for the b array). The p data has to be subdivided in four nearly equal parts. This is the more difficult part of the assignment, and should be attempted only by good students. No further details are given here. You are on your own.

Dynamic programming techniques, which appear in numerous sites on the internet, are **not** feasible for almost all of the problems in the .h files.

Merkle-Hellman cryptosystem (part 5, the written report)

The written report must have:

- A title page (front page) with the name of the course, the name of the report, date, the numbers and names of the students, and an estimate of the percentage each student contributed to the project. The sum of percentages should be 100%.
- A short introduction describing the problem; the source of any material adapted from the internet must be properly cited.
- A **small** description of the methods used to find the solutions.
- A description of the solutions found; this should include graphs of the execution time of the program as a function of the size of the problem.
- Comments or attempts at explanations of the results found. This can be placed near where the results are presented in the report.
- An appendix with all the code and data (bit vectors) of all solutions. Use a small font.
- **Deliverable: one PDF file. No archives, please!**

Suggestion: to simplify placing the solutions in the report, make your program output them in a format that makes the work of pasting them in the report easy. The same applies to the execution times (both for making tables to to making graphs).