

Hashing

Databases and keys. A *database* is a collection of *records* with various attributes. It is commonly represented as a *table*, where the *rows* are the records, and the *columns* are the attributes:

Number	Name	Dept	Alias
20090612	오재훈	산디과	alpha0401
20100202	강상익	무학	scala
20100311	손호진	무학	python.is.great

Often one of the attributes is designated as the *key*. Keys must be unique—there can be only one row with a given key. A table with keys is called a *keyed table*.

A database index allows us to quickly find a record with a given key. In data structure terms, this is just a map from the key type to the record type.

Basic hashing with chaining. Let's implement a database of all the students in the class. We start with a `Student` class: Each record will be an object of this class.

```
case class Student(name: String, id: Int, dept: String, alias: String)
```

We will use the `id` field as the key, and want to build an index so that we can quickly find the `Student` object with a given key.

A simple idea is to make an array with 100 slots, called a *hash table*. We use the last two digits of the student id as the index into this array. This is the *hash function* of the student id. Unfortunately, this doesn't quite work, because there are students with identical last two digits, like these:

20100874	정민수	무학	ubuntu
20080174	방태수	산디과	apple

So what we do instead is to store in each slot of the array a *linked list* of (key, record) - pairs, as in Fig. 1.

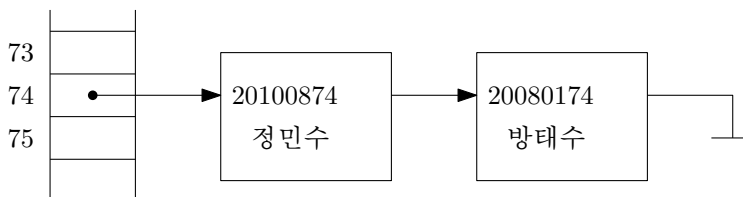


Figure 1: Chaining

To search for a key k in this data structure, we take the hash function of k (that is, the last two digits of k) to obtain the head of a linked list from the array. We then search the entire linked list by comparing k with each key in this linked list.

Implementation. The implementation consists of these private members:

```
private case class Node(val key: Int, var value: Value, val next: Node)
private val table = new Array[Node](100)
private def hash(key: Int): Int = key % 100
```

`contains` and `apply` make use of a private method `findNode`:

```
private def findNode(key: Int): Node = {
    val hashCode = hash(key)
    var n = table(hashCode)
```

```

while (n != null) {
  if (n.key == key)
    return n
  n = n.next
}
null
}

def contains(key: Int): Boolean = (findNode(key) != null)

def apply(key: Int): Value = {
  val n = findNode(key)
  if (n != null)
    n.value
  else
    throw new NoSuchElementException
}

```

To insert a (key, record) pair, we first check if the key already exists:

```

def update(key: Int, value: Value) {
  val hashCode = hash(key)
  var n = table(hashCode)
  while (n != null) {
    if (n.key == key) {
      n.value = value
      return
    }
    n = n.next
  }
  // add new node
  table(hashCode) = Node(key, value, table(hashCode))
}

```

Analysis of hash tables. How can we analyze this data structure? Clearly, in the worst case all students could have the same last two digits, and then the running time of insert and search operations would be $O(n)$, where n is the number of students.

So to analyze hashing, we need to make the assumption that the hash function (in our case, mapping the keys to their last two digits) is good: it should distribute the keys onto the slots of the array uniformly.

Formally, we will pretend as if the hash function was *random*: We will analyze the *expected running time* of insertions and search operations, where the expectancy is with respect to the random “choice” made by the hash function. In reality, of course, the hash function cannot be random: Given the same key twice, it must return the same hash value each time. Nevertheless, it turns out that with a good hash function, this analysis reflects the behavior of hash tables well.

So this is our setting: Given a key x , the hash value $h(x)$ is a random element of $\{0, 1, \dots, N - 1\}$, where N is the size of the hash table. The random choices must be *independent*, so that given two keys x and y , the probability that $h(x) = h(y)$ is exactly $1/N$. (Further down, we will even need independence for the entire set of keys hashed into the table.)

Analysis of chaining. We assume that a set Y of n items has already been hashed into a hash table of size N . We now consider the insertion of a new key x . The running time of this insertion is proportional

to the length of the linked list in slot $h(x)$. So what is the expected length of this list, with respect to the random choices made by the hash function?

Every element $y \in Y$ contributes one to the length of this list if $h(y) = h(x)$. Since we assume that hash values are independent, this happens with probability $1/N$. By additivity of expectation, this means that the expected number of items $y \in Y$ with $h(y) = h(x)$ is n/N , and so the expected length of the linked list in slot $h(x)$ is n/N .

Let's define the *load factor* λ of a hash table to be the ratio $\lambda = n/N$. By the above, the expected running time of insertions is $O(\lambda)$, and the same analysis applies to search operations and deletions.

Closed and open addressing. Hashing with chaining works well, and is often implemented in practice. However, the linked lists are not very memory-efficient, as we need space for node objects. We could make the data structure much more compact if we could store *all* the data inside the hash table. *Open addressing* means that we allow an item to be stored at a slot that is different from the “correct” slot indicated by its hash function. On the contrary, *closed addressing* means that items have to be stored at the slot given by their hash function, that is by chaining.

Linear probing. The most commonly used form of open addressing is *linear probing*. It also seems to be the best method to be used in practice, as long as the load factor remains low enough.

In linear probing, the first choice for a key x is the slot $h(x)$. If this slot is already in use, we try the slot $h(x) + 1$. If this is also already in use, we continue with $h(x) + 2$, and so on. For an insertion, we continue until we find an empty slot, and insert the item there. For a search operation, we continue until we either find the item, or find an empty slot. In the latter case, we know that the key is not in the hash table.

Figure 2 shows the result of a sequence of insertions into a hash table of size $N = 10$. The hash function of key x is its last digit, that is $h(x) = x \bmod 10$. Keys 89 and 18 are inserted in their “correct” slot. Key 49

	insert 89	insert 18	insert 49	insert 58	insert 9	insert 30
0			49	49	49	49
1				58	58	58
2					9	9
3						30
4						
5						
6						
7						
8		18	18	18	18	18
9	89	89	89	89	89	89

Figure 2: Insertions with linear probing.

should go in slot 9, but that is already occupied, so it goes into slot 0. Key 58 needs to try 3 slots before finding the empty slot 1, and the same is true for keys 9 and 30.

Consider the search for key 61: We first try slot $h(61) = 1$, but $58 \neq 61$. We then try slot 2, but $9 \neq 61$. We continue to slot 3, but $30 \neq 61$. We then check slot 4, which is empty, and so we can conclude that key 61 is not in the hash table.

Implementation. The implementation relies on the private method `findSlot`, which searches the hash table for a slot containing the given key, or otherwise an empty slot:

```

private def findSlot(key: Int): Int = {
  var h = hash(key)
  while (true) {
    val n = table(h)
    if (n == null || n.key == key)
      return h
    h += 1
    if (h == SIZE)
      h = 0
  }
  0 // make compiler happy
}

```

Insertions and search operations are now easy to implement in terms of `findSlot`:

```

def update(key: Int, value: Value) {
  val h = findSlot(key)
  if (table(h) != null)
    table(h).value = value
  else
    table(h) = Entry(key, value)
}

```

```

def contains(key: Int): Boolean = (table(findSlot(key)) != null)

```

```

def apply(key: Int): Value = {
  val h = findSlot(key)
  if (table(h) != null)
    table(h).value
  else
    throw new NoSuchElementException
}

```

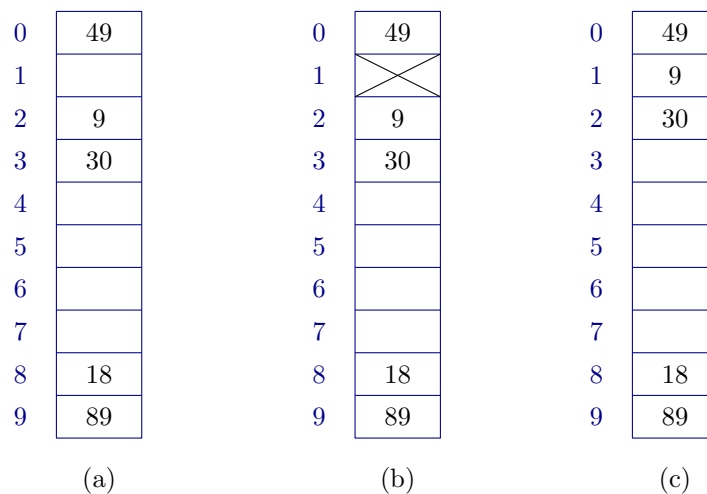


Figure 3: Deletions in linear probing: (a) this doesn't work, as we can no longer find 9 and 30. (b) mark deleted slot as "available". (c) moving items to fill the holes.

Deletions in linear probing. Deleting items from a hash table with linear probing is not as easy as for chaining. Consider the table of Figure 2, and assume we simply delete key 58. The result is shown in Figure 3(a). However, this doesn't work: Imagine if we now search for key 9. The search would terminate at slot 1 and report that 9 is not in the hash table!

There are two solutions to this problem: The easiest and common solution is shown in Figure 3(b): we mark slot 1 as “available”. This means that we consider the slot as empty during insertions (and so a new item may be placed here), but not during search operations (and so we can still correctly find keys 9 and 30).

With a bit more thinking, we can find a solution that doesn't need a special marker. Whenever we delete an item from a slot h , we create a *hole*. We then continue to scan the hash table from slot $h + 1$. If we find an item y in slot k , we move it to the hole at h if h lies in between $h(y)$ and k . This move will create a new hole at slot k , and so we continue scanning at slot $k + 1$ to fill the hole at slot k . The process continues until we find an empty slot. Figure 3(c) shows the result of deleting 58 from the hash table of Figure 2. Figure 4 shows a more interesting example: Key 89 is deleted and creates a hole in slot 9. This hole cannot be filled by key 30, but is filled with key 68. This creates a new hole in slot 1, which is filled by key 91. This creates a new hole in slot 2. Note that this hole is never filled, as the scan ends in slot 5 without finding an item that can be moved.

delete 89

0	30	30	30	30
1	68	68		91
2	91	91	91	
3	83	83	83	83
4	73	73	73	73
5				
6				
7				
8	18	18	18	18
9	89		68	68

Figure 4: Deletions in linear probing by filling the hole.

Simplified analysis of linear probing. We assume that a set Y of n items has already been hashed into a hash table of size N using linear probing. We now consider the insertion of a new key x . The running time of this insertion is proportional to the number of *probes* we need to make to find an empty slot. Here, a probe means that we check a slot to see if it is empty.

Since there are n items in a hash table of size N , “on average” a slot contains an item with probability $n/N = \lambda$. This means that a slot is free with probability $1 - \lambda$. So we are repeating the following experiment: We flip a coin. With probability λ , it comes up “occupied”, with probability $1 - \lambda$, it comes up “free”. The expected number of times we have to repeat this until we see the outcome “free” is $1/(1 - \lambda)$.

In other words, when $\lambda = 1/2$, we expect to check only two slots. When $\lambda = 3/4$, we already need to check four slots, and for $\lambda = 0.9$, it would be ten slots. This already shows that linear probing can only work well when the load factor remains significantly smaller than one.

If we use a special marker to implement deletions as in Figure 3(b), then the marked slots do not count as free in the computation of the load factor. It is therefore possible that a hash table is nearly empty, but still has a high load factor! Such a hash table should be cleaned up by *rehashing* all the items.

Real behavior. Unfortunately, the argument above is not correct, because the probabilities for each slot to be filled are not independent. Let's call a *run* a maximal sequence of consecutive filled slots (that is, a sequence of slots $h, h + 1, \dots, h + k - 1$ (modulo N) that are all filled, while $h - 1$ and $h + k$ are empty). Consider a run of length k . For every item that is hashed into the run, the run will grow by one element. This happens with probability $(k + 2)/N$, and so the probability grows with the length of the run. This effect causes *clusters* to appear in the hash table.

The file `clustering.scala` implements an experiment to verify this. For a given table size N and load factor λ , it performs an experiment: It first fills a table A of size N by independently filling each slot with probability λ . It then fills a second table B of size N by inserting $\lambda * N$ randomly chosen items using linear probing. Finally, it evaluates the expected search cost in both tables as follows: For each slot h , we count the number of probes needed to find an empty slot.

The following shows the result for $N = 10000$, taking the average over 1000 experiments to get more accurate results:

λ	Table A	Table B
0.5	2.0	2.5
0.7	3.3	6.0
0.9	10.0	49.5
0.95	20.0	182.1
0.99	100.0	1750.5

It turns out that the simplified argument above is far too optimistic. When the load factor increases, the expected search time grows much faster than $1/(1 - \lambda)$.

Complexity of linear probing. It can be shown that the expected number of probes needed for the insertion of a new item (or equivalently, for an unsuccessful search) is:

$$\frac{1}{2} \left(1 + \frac{1}{(1 - \lambda)^2} \right)$$

So the running time grows with the *square* of $1/(1 - \lambda)$, which explains why it is important to keep the load factor small. It should not exceed 0.5 or maybe 0.75.

Analysis of linear probing. We cannot prove the precise formula above here, but we want to at least show that the expected time for insertions and search operations is bounded by a constant independent of n , for the case $\lambda = 1/2$.

So we assume that we have a set Y of n items that have already been inserted into an array of length $N = 2n$ using linear probing. We consider the insertion of a new item with key x .

We start at slot $h = h(x)$, and visit slots $h, h + 1, h + 2, \dots, h + t$ (modulo N), where slot $h + t$ is the first free slot visited. The item will be inserted in slot $h + t$.

The number of probes needed for this insertion is $t + 1$. We observe that the number t depends only on the hash function $h = h(x)$ of the key x , so let's denote it as $t = t(h)$.

We assume that the hash function $h(x)$ is a random element of $\{0, 1, \dots, N - 1\}$, so to find the expected number E of probes for the insertion, we need to compute the average

$$E = \frac{1}{N} \sum_{h=0}^{N-1} t(h).$$

We note next that $t(h) = 0$ if slot h is an empty slot. The sum is therefore determined only by the occupied slots. A run of length k in the hash table contributes at most k^2 to the sum. (The precise contribution is $k(k + 1)/2$.)

We can therefore bound the sum above as

$$E = \frac{1}{N} \sum_{h=0}^{N-1} t(h) \leq \frac{1}{N} \sum_{h=0}^{N-1} \sum_{k=1}^n P(h, k) \cdot k^2,$$

where $P(h, k)$ is the probability that the slots $h, h + 1, h + 2, \dots, h + k - 1$ form a run of length k . (This probability is with respect to the insertions of the n items in Y into the hash table.)

When do these k slots form a run? A necessary condition is that exactly k of the n items have a hash index in the range $\{h, h + 1, \dots, h + k - 1\}$, while the remaining $n - k$ items have a hash index *not* in this range.

For a fixed item y , the probability that it is in the range is k/N , and the probability that it is not in the range is $1 - k/N$. It follows that

$$P(h, k) \leq P_k = \binom{n}{k} \left(\frac{k}{2n}\right)^k \left(1 - \frac{k}{2n}\right)^{n-k}$$

Since P_k is independent of h , we can bound

$$E \leq \frac{1}{2n} \sum_{h=0}^{2n-1} \sum_{k=1}^n P_k k^2 = \sum_{k=1}^n P_k k^2.$$

We now need two inequalities:¹

$$\binom{n}{k} \leq \frac{n^n}{k^k (n - k)^{n-k}}, \tag{1}$$

$$\left(1 + \frac{x}{n}\right)^n \leq e^x \quad \text{for } x \geq 0. \tag{2}$$

It follows that

$$\begin{aligned} P_k &\leq \frac{n^n}{k^k (n - k)^{n-k}} \left(\frac{k}{2n}\right)^k \left(1 - \frac{k}{2n}\right)^{n-k} \\ &= \left(\frac{n}{k}\right)^k \left(\frac{k}{2n}\right)^k \left(\frac{n}{n - k}\right)^{n-k} \left(1 - \frac{k}{2n}\right)^{n-k} \\ &= \frac{1}{2^k} \left(\frac{n}{n - k} \frac{2n - k}{2n}\right)^{n-k} = \frac{1}{2^k} \left(\frac{2n - k}{2(n - k)}\right)^{n-k} = \frac{1}{2^k} \left(\frac{2(n - k) + k}{2(n - k)}\right)^{n-k} \\ &= \frac{1}{2^k} \left(1 + \frac{k}{2(n - k)}\right)^{n-k} = \frac{1}{2^k} \left(1 + \frac{k/2}{n - k}\right)^{n-k} \leq \frac{1}{2^k} e^{k/2} = \left(\frac{\sqrt{e}}{2}\right)^k. \end{aligned}$$

Since $\sqrt{e}/2 < 1$, the infinite series $\sum_{k=1}^{\infty} (\sqrt{e}/2)^k k^2$ converges to some constant C , and we have $E \leq C$.

Comparison of linear probing and chaining. Linear probing works well when the load factor remains low, say $\lambda \leq 0.5$. It is, however, more sensitive to the quality of the hash function—chaining is more robust when the hash function is not so good.

Linear probing works particularly well when all data is allocated inside the array in C or C++. In that case, it requires no dynamic memory allocation at all, and consecutive probes are extremely fast as they exploit caching in the processor (a processor can access consecutive memory locations much faster than data spread out over the memory). Open addressing is therefore the method of choice for hashing in devices like routers etc.

¹The first one follows by writing $n^n = (k + (n - k))^n = \sum_{i=0}^n \binom{n}{i} k^i (n - k)^{n-i}$, which implies $n^n \geq \binom{n}{k} k^k (n - k)^{n-k}$. The second is true because $(1 + x/n)^n$ is an increasing series whose limit is e^x , as you should have learnt in Calculus I.

Hash functions. When keys are integers, we only need to somehow map them to the range $\{0, \dots, N - 1\}$. In many applications, however, keys are other objects, such as strings. In that case, we first have to map the string to an integer.

In this case, we speak about the *hash code* as a function h_1 from keys to integers, and the *compression function* as a function h_2 from integers to the index range of the hash table, that is $\{0, 1, \dots, N - 1\}$. The *hash function* is then the function $x \mapsto h_2(h_1(x))$ from keys to indices.

Hash codes and compression functions are a bit of a black art. The ideal hash code function should map keys to a uniformly distributed random slot in $\{0, \dots, N - 1\}$. This ideal is tricky to obtain. In practice, it's easy to mess up and create far more collisions than necessary.

Rehashing. When the load factor has become too large, or when a hash table contains too many slots marked as “available” due to deletions, we need to *rehash* the table.

This means that we create a new table (typically of about twice the size of the previous one), and change the compression function to map to the index range of the new table. We consider all items in the current table one by one, and insert them into the new table.

You *cannot* just copy the linked lists of a hash table with chaining to the same slots in the new table, because the compression functions of the two tables will certainly be incompatible. You have to rehash each entry individually.

You can also shrink hash tables (for instance when $\lambda < 0.25$) to free memory, if you think the memory will benefit something else. (In practice, it's only sometimes worth the effort.)

Obviously, an operation that causes a hash table to resize itself takes more than constant time; nevertheless, the *average* over the long run is still constant time per operation.

Compression functions. Let's consider compression functions first. Suppose the keys are integers, and each integer's hash code is itself, so $h_1(x) = x$. The obvious compression function is

$$h_2(x) = x \bmod N.$$

Hash codes are often negative, so remember that `mod` is not the same as Scala's `%` operator. If you compute `x % N`, check if the result is negative, and add `N` if it is.

But suppose we use this compression function, and $N = 10,000$. Suppose for some reason that our application only ever generates keys that are divisible by 4. A number divisible by $4 \bmod 10,000$ is still a number divisible by 4, so three quarters of the slots are never used! We have at least four times as many collisions as necessary. (This is an important example, because in many programming languages, such as C and C++, memory addresses of objects, when converted to integers, are always divisible by 4 or even by 8.)

The same compression function is much better if N is a prime number. With N prime, even if the hash codes are always divisible by 4, numbers larger than N often hash to slots not divisible by 4, so all slots can be used.

For reasons we won't explain

$$h_2(x) = ((a \cdot x + b) \bmod p) \bmod N$$

is a better compression function. Here, a , b , and p are positive integers, p is a large prime, and $p \gg N$. Now, the number N doesn't need to be prime.

Use a known good compression function like the two above instead of inventing your own. Unfortunately, it's still possible to mess up by inventing a hash code that creates lots of conflicts even before the compression function is used.

Hash codes. In Scala, every object automatically has a method `##` that returns a hash code. For basic types like `Int`, `Double`, `String`, `List`, `Set`, or `Map` this method returns a useful hash code.

When you create your own objects, you may need to design a hash code specially for this object. Here is an example of a good hash code for strings:


```
def hashCode(key: String): Int = {
  var hash = 0
  for (ch <- key)
    hash = (127 * hash + ch) % 16908799
  hash
}
```

By multiplying the hash code by 127 before adding in each new character, we make sure that each character has a different effect on the final result. The % operator with a prime number tends to “mix up the bits” of the hash code. The prime number is chosen to be large, but not so large that $127 * \text{hash} + \text{ch}$ will ever exceed the maximum possible value of an `Int`.

The `##`-method of Scala strings seems to be simpler than this, and appears to do the following:

```
def hashCode(key: String): Int = {
  var hash = 0
  for (ch <- key)
    hash = 31 * hash + ch
  hash
}
```

The best way to understand good hash codes is to understand why bad hash codes are bad. Here are some examples of bad hash codes on words.

- Sum up the values of the characters. Unfortunately, for English words the sum will rarely exceed 500 or so, and most of the entries will be bunched up in a few hundred buckets. Moreover, anagrams like “pat,” “tap,” and “apt” will collide.
- Use the first three letters of a word, in a table with 26^3 buckets. Unfortunately, words beginning with “pre” are much more common than words beginning with “xzq”, and the former will be bunched up in one long list. This does not approach our uniformly distributed ideal.
- Consider the “good” `hashCode()` function written out above. Suppose the prime modulus is 127 instead of 16908799. Then the return value is just the last character of the word, because $(127 * \text{hash}) \% 127 = 0$. That’s why 127 and 16908799 were chosen to have no common factors.

Why is the `hashCode()` function presented above good? Because we can find no obvious flaws, and it seems to work well in practice. (A black art indeed.)

Equality and hash codes. Hash tables only work correctly if equal objects have the same hash code. The following example shows what can go wrong:

```
scala> class Point(val x: Int, val y: Int) {
  |   override def equals(rhs: Any): Boolean = {
  |     rhs match {
  |       case q: Point => x == q.x && y == q.y
  |       case _ => false
  |     }
  |   }
  | }
scala> val p = new Point(3, 5)
p: Point = Point@1535ac
scala> val q = new Point(3, 5)
q: Point = Point@15ddf5
scala> p == q
```

```

res0: Boolean = true
scala> var s = new scala.collection.mutable.HashSet[Point]
scala> s += p
scala> s contains p
res1: Boolean = true
scala> s contains q
res2: Boolean = false

```

We define a new type `Point` with an equality operator. As you can see, `p == q` is true, because both objects contain the same coordinates. However, if we add `p` to a set and then test if it contains `q`, the answer is *no*.

What went wrong? Here is a clue:

```

scala> p.##
res3: Int = 1389996
scala> q.##
res4: Int = 1433077

```

`p` and `q` have different hash codes! And so the `HashSet` (which uses a hash table) has no chance to find `p` when using `q` as a key. It simply searches in the wrong slot!

To fix this problem, we need to redefine our class:

```

scala> class Point(val x: Int, val y: Int) {
  |   override def hashCode: Int = 41 * x + y
  |   override def equals(rhs: Any): Boolean = {
  |     rhs match {
  |       case q: Point => x == q.x && y == q.y
  |       case _ => false
  |     }
  |   }
  | }
scala> val p = new Point(3, 5)
scala> val q = new Point(3, 5)
scala> p == q
res0: Boolean = true
scala> p.##
res1: Int = 128
scala> q.##
res2: Int = 128
scala> var s = new scala.collection.mutable.HashSet[Point]
scala> s += p
scala> s contains p
res3: Boolean = true
scala> s contains q
res4: Boolean = true

```

By overriding the method `hashCode`, both objects now have the same hash code when they are equal, and the `HashSet` works correctly.

We could have achieved this much more easily by defining `Point` as a case class. Case classes automatically define equality and hash code using the fields of the object:

```

scala> case class Point(x: Int, y: Int)
scala> val p = Point(3, 5)
p: Point = Point(3,5)
scala> val q = Point(3, 5)

```

```

q: Point = Point(3,5)
scala> p == q
res0: Boolean = true
scala> p.##
res1: Int = -1839273401
scala> q.##
res2: Int = -1839273401

```

The lesson is: hash tables require that keys satisfy the following “contract”:

If `obj1 == obj2` then `obj1.## == obj2.##`.

Mutable keys. A similar problem happens when keys are mutable. Consider the following example:

```

scala> case class Point(var x: Int, var y: Int)
scala> val p = Point(3, 5)
p: Point = Point(3,5)
scala> val s = new scala.collection.mutable.HashSet[Point]
s: scala.collection.mutable.HashSet[Point] = Set()
scala> s += p
res3: s.type = Set(Point(3,5))
scala> s contains p
res4: Boolean = true
scala> p.x = 4
p.x: Int = 4
scala> p
res5: Point = Point(4,5)
scala> s
res6: scala.collection.mutable.HashSet[Point] = Set(Point(4,5))
scala> s contains p
res7: Boolean = false

```

Even though `s` clearly contains `Point(4,5)`, calling `s contains p` returns false. The reason is that `p`'s hash code has changed after it was added to the hash table, so `p` is simply in the wrong slot of the hash table!

The lesson here: *Never modify keys* after they were added to a hash table.

In fact, I would go further and recommend: Never use mutable objects as keys in a hash table. This is yet another example why immutable objects make programming safer and easier.

Acknowledgments. Once again I made use of Jonathan Shewchuk's lecture notes. I learnt the analysis of linear probing presented here from David Eppstein, and proofs of Inequality (1) from Antoine Vigneron, Jingil Choi, and Sarel Har-Peled.