# REPORT LAB 5 and 6: CS 5542

File  Edit  View  Navigate  Code  Analyze  Refactor  Build  Run  Tools  VCS  Window  Help

scala-library-2.11.7-sources.jar  ›  scala  ›  collection  ›  MapLike.scala

```
     *  The method implemented here throws an exception,
     *  but it might be overridden in subclasses.
     *
     *  @param key  the given key value for which a binding is missing.
     *  @throws NoSuchElementException
     */
    def default(key: A): B =
      throw new NoSuchElementException("key not found: " + key)

    protected class FilteredKeys(p: A => Boolean) extends AbstractMap[A, B] with DefaultMap[A, B] {
      override def foreach[C](f: ((A, B)) => C): Unit = for (kv <- self) if (p(kv._1)) f(kv)
      def iterator = self.iterator.filter(kv => p(kv._1))
      override def contains(key: A) = self.contains(key) && p(key)
      def get(key: A) = if (!p(key)) None else self.get(key)
    }

    /** Filters this map by retaining only keys satisfying a predicate.
     *  @param p    the predicate used to test keys
     *  @return an immutable map consisting only of those key value pairs of this map where the key satisfies
     *          the predicate `p`. The resulting map wraps the original map without copying any elements.
     */
    def filterKeys(p: A => Boolean): Map[A, B] = new FilteredKeys(p)
```

Run  SparkDecisionTree

```
16/03/02 23:21:14 INFO MemoryStore: MemoryStore cleared
16/03/02 23:21:14 INFO BlockManager: BlockManager stopped
16/03/02 23:21:14 INFO BlockManagerMaster: BlockManagerMaster stopped
16/03/02 23:21:14 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
16/03/02 23:21:14 INFO SparkContext: Successfully stopped SparkContext
16/03/02 23:21:14 INFO ShutdownHookManager: Shutdown hook called
16/03/02 23:21:14 INFO ShutdownHookManager: Deleting directory C:\Users\ry6d3\AppData\Local\Temp\spark-c9eb28ab-e9f8-4ac0-86e5-11e24102d28f

Process finished with exit code 1
```

Terminal    Java Enterprise    4: Run    6: TODO

All files are up-to-date (yesterday 11:20 PM)                264:1   LF   UTF-8