# AIS Classifier
# (draft)

Thomas Nordli

# AIS Classifier (draft)

Thomas Nordli <tn@hive.no>

June 29, 2011

# Contents

# 1 Abstract

This report describes a system that makes it possible to explore the AIS data broadcasted by ships. The goal of such a system is to determine how/if it is possible to predict the ship type solely based on AIS messages sent by ships, ignoring the ships self proclaimed type.

The system is intended to be used to discover anomalous behavior by comparing the predicted type with the self proclaimed one. Eventually it is suppose to work online, processing live feeds of data. By now it is only working offline.

The report is to be considered a draft, as the prototype is still missing vital functions, and is not thoroughly tested.

# 2 Introduction

## 2.1 What is AIS?

Automatic Identification System, is a system that obligates some ships to transmit some data over VHF radio. The obligation is a part of the SOLAS convention and applies to all passenger ships, ships with tonnage of 300 or above in international waters and cargo ships trafficing national with tonnage of 500 or above [4].
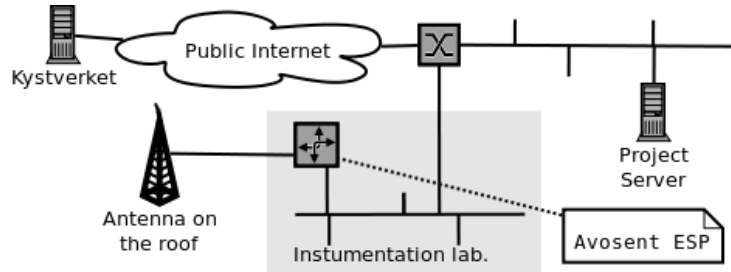
## 2.2 Problem definition

Is it possible to classify the ships based on the fields in the AIS messages (excluding the ship type) alone? The prototype system described here is a contribution to a framework for investigating this question.If the answer to the question turns out to be positive, a software module based on the prototype system developed, can be used for spotting anomalies and raising alarms. Anomaly detection systems have usually a high frequency of false positives, therefore the system must have parameters that may be tuned to lower the false positive rates. Such a software can be standalone or integrated with existing alarm systems.

## 2.3 Assumptions and limitations

This work is done under the assumption that there is detectable differences in the normal behavior of different classes of ships.AIS includes 27 message types. Only message of type 1, and 5 for self proclaimed type, is treated in this report, as these are the only one used in the prototype system.Since this project have not yet identified any effective discriminating features, only four of the fields from message type 1, are included. These fields are rotation (ROT), course over ground (COG), speed over ground (SOG) and heading (HEAD). The mean and variance of these are used as features, giving feature vectors with eight dimensions. These fields and the statistics were chosen because they were easy to implement, as examples to demonstrate the concept. They turned out to have poor ciscriminatroy capabilities.

# 3  Overview of the system.

Kystverket
Public Internet
Project
Server
Antenna on
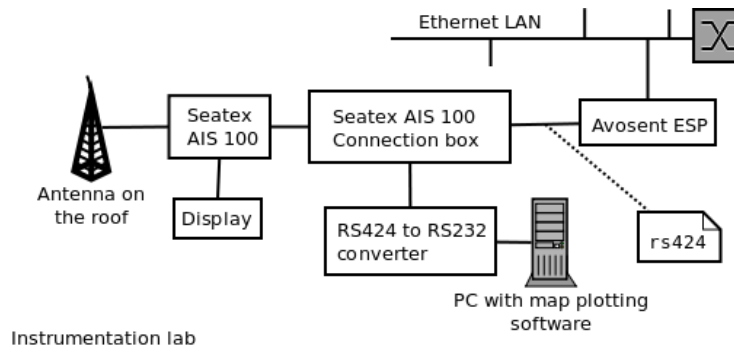the roof
Instumentation lab.
Avosent ESP

## 3.1  Kystverket

Kystverket, the Norwegian Coastal Administration, is the national agency for coastal management, maritime safety and -communication in Norway. They provide live AIS data for this project via a TCP server.

## 3.2  Project Server

A server running linux is set up to capture, prosess and store the data. The software described below is run on this server.

## 3.3  Instrumentation laboratory

Ethernet LAN
Antenna on
the roof
Seatex
AIS 100
Seatex AIS 100
Connection box
Avosent ESP
Display
RS424 to RS232
converter
PC with map plotting
software
rs424
Instrumentation lab

Vestfold University College has an maritime instrumentation lab with some AIS equippement connected to an antenna on the roof. While waiting for Kystverket to process our application for access to live AIS data, a Avosent ESP was purchased, and connected to Seatex AIS 100 via an Connection box. The Avosent ESP was configured to recieve data on an rs424 link and distributing it via TCP on the lab's ethernet LAN. Via the university college ethernet switch, the data was thereby availiable to the project server.

# 4  Expected Quality of the Data

Harati-Mokthari et al. discusses errors found in live data. Many of the ship types are very vage or misleading, such as "vessel" or "cargo". For several

types the regulations of use are ambiguous. One of the studies investigated by Harati-Mokthari et al., found vessels that were registered with the same vessel type in Lloyds register database, transmitted different vessel types in the AIS messages[2].

Kjerstad also describes some sources of errors that can be found in AIS data[4].

# 5 Software

## 5.1 netcat

Netcat [3] is called *the swiss army knife of TCP/IP*. It is in this project used to connect to the data soucres, recieve the data and printing it out to *standard output*.

Command: nc nmea.hive.no 4002

## 5.2 devtools_ais

*devtools_ais*, a AIVDM/AIVDO decoder, is part of *gpsd*[6]. It is accompanied with a detailed documentation of how to decode AIVDM/AIVDO sentences [5]. In this project it is used as a command line tool to convert AIVDM data, that comes from the AIS recievers, to a common data format with comma-separated values (CSV).

Command: `./devtools_ais.py -c`

## 5.3 TrackSplitter

The CSV file produced by devtools_ais, contains all recieved messages from all ships. The features that we need, should be per ship, therefore a little script is made, that reads the CSV file and splits it into multiple files named in a directory. Each of the files produced contains messages from one ship only. They are named with the ship's MMSI number. The script is called *tracksplitter*, and you can find it listed in appendix A.

## 5.4 TrackDescriber

Each of the files produced by tracksplitter, is to be processed by trackdescriber. Trackdescriber reads data from standard input (STDIN) and makes features that characterize a 'track'. By track is here meant all messages comming from one ship.As there is still not identified any discriminating features, only a few example features is implemented. Four fields are used: ROT, SOG, COG and HEAD. Each field is represented by an instance of a class called 'dataField'.

```
┌─────────────────────────┐
│        dataField        │
├─────────────────────────┤
│ sum : float             │
│ sumSqr : float          │
│ n : int                 │
├─────────────────────────┤
│ var()                   │
│ addData()               │
│ mean()                  │
└─────────────────────────┘
```

The class keeps track of three attributes: The sum, the sum of squares and the number of recieved messages (n). Each time a new message arrives, the value is passed to the method addData, and the three attributes are updated. If the script is left running reading real time data (from a live feed, and not from a file), it will eventuall produce an arithmetic overflow. The issue of arithmetic overflow is not yet dealt with.Four objects of this class is instantiated, one from each of the above mentioned AIS fields. STDIN is read line by line, adding data for each line. Each line represents one AIS message. When the features are to be generated (and printed to STDOUT, var() and mean() are called for, producing the mean and variance. This gives a feature vector of eigth dimensions. This printing is preceded with a aproriate header with meta data intended for reading with a classifier written using the python scripting interface from the data mining tool *orange*[1].

TrackDescriber is listed in appendix B

## 5.5   TrackClassifier

The script trackclassifier reads the feature vectors generated by TrackDescriber, and predict a ship type. As there are still no classifiers developed, all the script does at this point in time, is to do some evalution of the discriminating powers of the perliminary features using some of the data mining techniques included in Orange [1]. It currently uses the *AIS ship type* as the class. The reason for this not being a good idea is explained in section 4.

TrackClassifier is listed in appendix C

Below you can the result of a run.

```
$ ./trackClassifier.py
Classes: 26
Attributes: 8

After feature subset selection with margin 0.010 (8 attributes):
Outliers are now removed
Before feature subset selection (8 attributes):
0.069 varCog
0.058 meanHead
0.049 varRot
0.048 meanSog
0.047 varSog
0.040 meanRot
0.039 meanCog
```

```
0.035 varHead
Classification accuracies:
bayes 0.0142857142857
tree 0.201388888889
Classification accuracies (leave one out):
bayes 0.0208333333333
tree 0.201388888889
$
```

# 6   Generation of the dataset

Data is comming in as a stream of text over TCP connections, one line per AIS packet. It is recieceved and logged to a file with the unix command 'nc' (netcat).

```
 $ head -n5 data/hive/nohup.out
!AIVDO,1,1,,,13o02wnP000gk;LQv:uh0?vv0000,0*07
!AIVDM,1,1,,B,D02R3f1HpNfq6D06D0,4*4F
!AIVDO,1,1,,,13o02wnP000gk;HQv:uP0?w20000,0*7E
!AIVDM,1,1,,A,13mDIj0P000h05FR0@:@0?w0000,0*36
!AIVDM,1,1,,B,13oE6P002Q0h;0VQwfTbS'Q200f,0*71
```

To start caputuring data the following command was used to log data comming in through the antenna on thr roof.nohup nc nmea.hive.no 4002 &

   The command nohup makes sure that the process created by the rest of the commandline is disconnected from the current session, so that it will continue after the session is terminated and the user that issued the command has logged out. The logging will now continue until an authorized user stops it, the systems shuts down or the storage space is filled up.The tool 'devtools_ais' extracts the packed bitfields from the aquired AIS messages and converted to a 'csv' (comma separated values) format.

```
$ head -n5 data/hive/nohup.out | ./devtools_ais.py -c
1,0,258999039,6,-128,0,0,6265198,35621623,0,511,31,0,0,0
20,0,2655160,1422,1,7,750,1125,1,7,1125,0
1,0,258999039,6,-128,0,0,6265196,35621622,0,511,33,0,0,0
1,0,257235400,0,-128,0,0,6291627,35655721,0,511,32,0,0,49183
1,0,259344000,0,0,161,0,6314003,35647122,2702,272,33,0,0,49198
```

Trackplitter reads the data produced by 'devtools_ais.py -c'. The data is split by mmsi, and written to temporary files in './tmpfiles'. Each file represents data from one ship, and is named with its mmsi number.

```
$  ls tmpfiles/ | head -n5
156198689
205203000
205585000
209324000
209350000
```

Each file in the directory contains information on one ship.

```
$ head -n5 tmpfiles/205203000
1,0,205203000,0,0,192,1,6336077,35873338,1810,180,54,0,0,81956
3,0,205203000,0,-15,193,1,6329874,35849725,1820,181,26,0,0,0
1,0,205203000,0,11,184,1,6369048,35786658,1740,174,54,0,0,81945
1,0,205203000,0,0,190,1,6370747,35779543,1730,173,12,0,0,21076
1,0,205203000,0,0,189,1,6371729,35775168,1730,173,36,0,0,81948
```

The track_describer Reads AIS messages from STDIN. It assumes that the messages is formatted as produced by './devtools.py -c', and that all the messages originates in the same ship. Features are calculated and output is printed on STDOUT.To be scalable, and to be able to continously processing data from live streams, It is desiged in a way that it does not have to store all the data. Only aggregations of data is stored (and updated for each incomming packet).The following command is used to write features of all the tracks to the file tabulator separated file called 'data.tab'.

```
for F in tmpfiles/* ;do ./track_describer.py  $F;done > data.tab
```

The file 'data.tab' has to be provided with a header, giving names and types to the columns. The column indicating the class (ship type) is also marked. This is done giving the file å header.

```
$ head -n5 data.tab
meanSog stdSog  meanHead stdHead meanCog stdCog  meanRot stdRot type
c       c       c        c       c       c       c       c      d
                                                                class
0.188781  29.984842 188.666667 27.537797   1802.664509 17832.525553  \
180.325782 173.444507  35
-0.452115 48.898434 41.386802   1782.776197 2222.091371 860973.688331 \
112.157699 6503.311616 70
```

The first line contains the field names, The second names declears type (c = continous, d = discrete). The third line marks the column that indicates the column that the classifier should aim to predict.The file 'data.tab' can now be used by the tool trackClassifier.

# 7   Result

The project has developed a system for producing data set, preprocessing the data, extracting features and analyzing them. The software works only offline, but is developed with the goal of beeing able to work both offline and online.Two datasets are produced. One based on local messages recived from the antenna on the roof. And one based on data comming from Kystverket.

# 8   Future work

This project has produced a framework for further investigation. In this section you will find some a brief description of some of the work that should be done, if the project is continued.

## 8.1 Features

As there is not yet identified any features, this is something to do further work on. When features are identified more preprossessing on some of the fields have to be done. The details of the domains of fields can be found in *AIVDM/AIVDO protocol decoding*[5]

## 8.2 Classes

Search for classes − the AIS field ship type, is not usable as a class field in a classifing algorithm. This is due to the unambigous use and regultations of the field, as explained in section 4.

## 8.3 Data

Add a message timestamp to each recieved message. Consider other types of data sources, in addition to AIS. This may be data from radar or from external databases.

## 8.4 Implementation

Put the data in relational database. This will make it easier to query the data. The issue of arithmetic overflow in the trackdesciber.dataField.addData() must be dealt with if the system is to be used in online processing.

# 9 Conclusion

A lot of work have been done to prepeare the ground for investigating the possibilities of data mining the AIS messages. To answer the question whether it is possible to classify the ships based on AIS data, while distrusting the ship type field, the project have to be continued.

# A Tracksplitter

The script is written as a part of this project to demultiplex the messages, based on their origin.

Listing 1: tracksplitter.py

```
1  #!/usr/bin/python
2  """ Reads data from standard input.
3  Data should be produced by 'devtools_ais.py −c '.
4  The data is split by mmsi, writing files
5  in the './tmp' directory
6  or the directory given with option '−d '.
7
8  Each file represents data from one ship ,
9  and is named with the mmsi number """
10
11
12 import os,sys,getopt
13
14 if __name__ == "__main__" :
```

```
15
16        datafiledir='tmp'
17        o, a = getopt.getopt(sys.argv[1:], 'd:')
18        opts = {}
19        for k,v in o:
20            if k == '-d':
21                datafiledir=v
22
23
24        ships = {}
25
26        # Makes sure the temporary directory is empty
27        if os.path.isdir(datafiledir):
28            for f in os.listdir(datafiledir):
29                os.unlink("%s/%s"%(datafiledir,f))
30        else:
31            os.mkdir(datafiledir)
32
33        while 1:
34
35            line = sys.stdin.readline() # Read from stdin
36            if not line:                        #EOF
37                break
38
39            mmsi=line.split(',')[2] # finds mmsi
40
41            if not ships.has_key(mmsi):
42
43                # stores file objects in dictionary
44                ships[mmsi] = open("%s/%s" % (datafiledir, mmsi), 'a')
45
46            ships[mmsi].write(line)
```

# B    Track Describer

The script is written as a part of this project to produce the (preliminary)
features describing a ship based on its AIS brodcasts.

Listing 2: trackdescriber.py

```
1   #!/usr/bin/python
2
3   """Reads AIS messages from STDIN.
4   Input format: CSV (as from './devtools.py -c | ./tracksplitter.py')
5
6   Output: Features printed on STDOUT
7
8   """
9
10  import sys, os
11
12  class dataField:
13      """ Class representing one of the fields of an AIS message. """
14
15      def __init__(self):
16          self.sumSqr    = 0.0
17          self.sum       = 0.0
18          self.n         = 0
19
20      def addData(self,newData):
21          self.sum += newData
```

```
22              self.sumSqr += newData**2
23              self.n +=1
24
25         def mean(self):
26              return self.sum / self.n
27
28         def var(self):
29              return ( self.sumSqr - self.sum*(self.sum/n) ) / (self.n -1
                    )
30
31         def __str__(self):
32              return "%f\t%f" % ( self.mean(), self.var() )
33
34  def readData():
35
36         mmsi = ''  # Maritime Mobile Service Identity
37         # navgitation status
38         rot  = 0   # Rate of turna
39         sog  = 0   # Speed over ground
40         # position accuracy
41         # longitude
42         # latitude
43         cog  = 0   # Course over ground
44         head = 0   # True heading
45         # time stamp
46         # manouver indicator
47         # RAIM flag
48         # Radio status
49
50         sType = int('100') # Ship type (100 is out of range)
51         n     = 0   # Number of samples
52
53         fieldnames = ['rot', 'sog', 'cog', 'head']
54         flds        = {}
55         for f in fieldnames:
56             flds[f] = dataField();
57
58
59         while 1:
60
61             line = sys.stdin.readline()   # Read from stdin.
62             if not line:                   # EOF
63                 break
64
65             fields = line.split(',')
66
67             messageType = int(fields[0])
68
69             if messageType == 5:            #
                    !!!!!!!!!!!!!!!!!!!!!!!!!!!!!!
70                 sType = int( fields[7] )
71                 continue
72
73             if messageType != 1:
74                 continue
75
76             n += 1
77             flds['rot'].addData(  int( fields[4]  ) )
78             flds['sog'].addData(  int( fields[5]  ) )
79             flds['cog'].addData(  int( fields[9]  ) )
80             flds['head'].addData( int( fields[10] ) )
81
```

```
82        mmsi  =  fields[2]
83        return  (flds, sType, n, mmsi)
84
85  if __name__ == "__main__":
86
87        (flds, sType, n, mmsi) = readData()
88
89        if (n>1):
90            print "%s\t%s\t%s\t%s\t%d" % (flds['rot'], flds['sog'],
                  flds['cog'], flds['head'], sType)
91            #print "%s (%d):|t%s|t%s|t%s|t%s|t%d" % (mmsi,n,flds['rot
                  '], flds['sog'], flds['cog'], flds['head'], sType)
```

## C    TrackClassifier

The script is written as a part of this project to start the development of a script
that constitutes the module responisble for classifying the ships based on the
feature vectors comming from TrackDescriber. As of now, it only do some test
on the preliminary features' discriminatory capabilities. It also uses the *AIS
ship type* as class. As explained in section 4.

Listing 3: trackClassifier.py

```
1   #!/usr/bin/python
2   import orngOutlier
3   import orngTree
4   import orange, orngFSS
5   import orngDisc, orngTest, orngStat
6   import orngClustering
7   import os
8
9   def accuracy(test_data, classifiers):
10       correct = [0.0]*len(classifiers)
11       for ex in test_data:
12           for i in range(len(classifiers)):
13               if classifiers[i](ex) == ex.getclass():
14                   correct[i] += 1
15       for i in range(len(correct)):
16           correct[i] = correct[i] / len(test_data)
17       return correct
18
19
20  def cross_validation(data, learners, k=10):
21       acc = [0.0]*len(learners)
22       selection = orange.MakeRandomIndicesCV(data, folds=k)
23       for test_fold in range(k):
24           train_data = data.select(selection, test_fold, negate=1)
25           test_data = data.select(selection, test_fold)
26           classifiers = []
27           for l in learners:
28               classifiers.append(l(train_data))
29           acc1 = accuracy(test_data, classifiers)
30           for j in range(len(learners)):
31               acc[j] += acc1[j]
32       for j in range(len(learners)):
33           acc[j] = acc[j]/k
34       return acc
35
36
37  def leave_one_out(data, learners):
```

11

```python
38
39      acc = [0.0]*len(learners)
40      selection = [1] * len(data)
41      last = 0
42      for i in range(len(data)):
43          selection[last] = 1
44          selection[i] = 0
45          train_data = data.select(selection, 1)
46          for j in range(len(learners)):
47              classifier = learners[j](train_data)
48              if classifier(data[i]) == data[i].getclass():
49                  acc[j] += 1
50          last = i
51
52      for j in range(len(learners)):
53          acc[j] = acc[j]/len(data)
54      return acc


57  def report_relevance(data):
58      m = orngFSS.attMeasure(data)
59      for i in m:
60          print "%5.3f_%s" % (i[1], i[0])


63  def setMargins():
64      global data
65      marg = 0.01
66      filter = orngFSS.FilterRelief(margin=marg)
67      ndata = filter(data)
68      data = ndata
69      print "\nAfter_feature_subset_selection_with_margin_%5.3f_(%d_
              attributes):" % \
70          (marg, len(data.domain.attributes))


73  def removeOutliers():
74
75      global data
76      outlierDet = orngOutlier.OutlierDetection()
77      ndata=orange.ExampleTable(data.domain)
78
79      outlierDet.setExamples(data)
80      z=outlierDet.zValues()
81      for i in range (len(data)):
82          if abs(z[1])<1.5:
83              ndata.append(data[i])
84
85      print "Outliers_are_now_removed"
86      data=ndata


89  def kmeans():
90      global data
91      for k in range(2,18):
92
93          km = orngClustering.KMeans(data, k, initialization=
                  orngClustering.kmeans_init_diversity)
94          score = orngClustering.score_silhouette(km)
95          print km.clusters, k, score
96
97
```

```
98   if __name__ == "__main__":
99
100      data = orange.ExampleTable("data.tab")
101
102      # report on number of classes and attributes
103
104      print "Classes:", len(data.domain.classVar.values)
105      print "Attributes:", len(data.domain.attributes)
106
107      setMargins()
108      removeOutliers()
109
110      print "Before_feature_subset_selection_(%d_attributes):" % len(
             data.domain.attributes)
111
112      report_relevance(data)
113
114      # commented out -- reporting result is not implemented, giving
             verbose output
115      # kmeans()
116
117
118      bayes = orange.BayesLearner()
119      tree = orngTree.TreeLearner(mForPruning=2)
120
121      bayes.name = 'bayes'
122      tree.name = 'tree'
123      learners = [bayes,tree]
124
125
126      acc = cross_validation(data, learners, k=10)
127
128      print "Classification_accuracies:"
129      for i in range(len(learners)):
130          print learners[i].name, acc[i]
131
132
133          acc = leave_one_out(data, learners)
134
135      print "Classification_accuracies_(leave_one_out):"
136      for i in range(len(learners)):
137          print learners[i].name, acc[i]
```

# References

[1] Orange. http://orange.biolab.si/.

[2] A. Harati-Mokhtari, A. Wall, P. Brooks, and J. Wang. Automatic Identification System (AIS): A Human Factors Approach. In *The Nautical Institute AIS Forum*, pages 1–11. Citeseer, 2007.

[3] hobbit@avian.org hobbit@avian.org. netcat.

[4] Norvald Kjerstad. *Elektroniske og akustiske navigasjonssystemer: for maritime studier*. Number 3. utg. Tapir akademisk forl., Trondheim, 2008.

[5] Eric S. Raymond. Aivdm/aivdo protocol decoding v1.25, June 2010. http://gpsd.berlios.de/AIVDM.html.

[6] Eric S. Raymond et al. gpsd — a gps service daemon.