

The Killer Shrimp Invasion Challenge on Kaggle:

An online competition tackling the spread of invasive marine species through machine learning

IMDIS Submission 43

Background

- invasive species a major challenge on environmental, social and economic levels
- Crowdsourcing: utilize the public's workforce and creativity
- Machine learning: increasingly advanced capabilities and availability

-> How to utilize crowdsourcing to engage the public in using machine learning (ML) in a marine context?

- This poster illustrates how Ocean Data Factory Sweden, a publicly funded innovation consortia, organized a challenge on Kaggle, a ML-focused crowdsourcing platform
- Goals:
 - o To build awareness
 - o To improve existing ML model (-> see submission 46)
 - o To attract interested people for ODF Sweden

Prologue

- ODF had already collected a large data set on the Killer Shrimp in the Baltic Sea
- Overarching question was
 - o How does the Shrimp spread in the Baltic Sea under different conditions?

Preparing and launching the Kaggle Challenge

- Register challenge on Kaggle
 - o "official" challenges are extremely expensive – good alternative are free "in-class" challenges
- define problem
 - o "to predict the presence of the Killer Shrimp in areas of the Baltic Sea"
 - o a catchy name "Killer Shrimp Invasion"
- upload
 - o training dataset with locations, physical parameters and shrimp presence to enable competitors to develop their ML models
 - o test dataset to evaluate the competitors' models on unseen data, i.e. predicting the presence of the Killer Shrimp in various areas, with 30% of the test set used to calculate the public leaderboard during the competition and the remaining 70% used to calculate the private leaderboard, i.e., the winners, at the end of the competition
 - o example of a submission file to enable competitors to understand what they should submit
- choose
 - o evaluation metric of "Area Under Receiver Operating Curve" (AUROC) between the predicted probability and the observed target
 - o reward (150€)
 - o timeline (~3 months)
- moderate

Results

- 30 participants, mainly data scientists
 - o Interest in "real-life" problems
 - o A catchy name helps
 - o Difficult with zero data science knowledge
 - o needs good data
- Kaggle is a good platform...
 - o For very streamlined, ML-focused competitions
 - o Not to attract non-data scientists
 - o But there are other interesting crowdsourcing alternatives (e.g. Zooniverse)
- Good moderation is key!
- Participants can be niftier than you think
 - o Figured out how to "cheat" and received a perfect score with a poor ML model
 - o This was done more out of curiosity than out of greed, and participants reacted positively to a subsequent rule change

Reflections

- This was an interesting learning experience for ODF and participants
 - o Yes, this type of competition works
 - o Practical value in this case? Not so much, since it was designed more as proof of concept
- Advice for others
 - o Think about whether and how much value you can get out of both ML and this type of open challenge
 - o Prepare good and solid data sets
 - o Think ahead how you want to spread the word about your challenge

Authors:

Adrian Bumann, Chalmers University of Technology (Sweden), adrian.bumann@chalmers.se
Robin Teigland, Chalmers University of Technology (Sweden), robin.teigland@chalmers.se
Jannes Germishuys, Combine AB (Sweden), jurie.germishuys@combine.se
Benedikt Ziegler, Combine AB (Sweden), benedikt.ziegler@combine.se
Martin Mattson, Medins Havs och Vattenkonsulter AB (Sweden), martin.mattsson@medinsab.se
Eddie Olsson, RISE - Research Institutes of Sweden (Sweden), eddie.olsson@ri.se
Robert Rylander, RISE - Research Institutes of Sweden (Sweden), robert.rylander@ri.se
Yixin Zhang, University of Gothenburg (Sweden), yixin.zhang@ait.gu.se
Torsten Linders, University of Gothenburg (Sweden), torsten.linders@gu.se



Figure 1 - the "Killer Shrimp"

• Open datasets

- 1) port locations in Europe (EMODNET Human Activities)
- 2) ocean surface temperatures and salinity
 - for Baltic Sea (SMHI - Swedish Meteorological and Hydrological Institute)
 - and North Sea regions (SeaDataNet)
- 3) presence data of *D. Villosus* (GBIF, approx. 3000 data points)
 - Pseudo-absence data from Baltic Sea (approx. 2.8m data points)
- 4) marine data layers (Bio-Oracle)
- 5) ocean temperature and salinity (Marine Copernicus)

Figure 2 - Data used

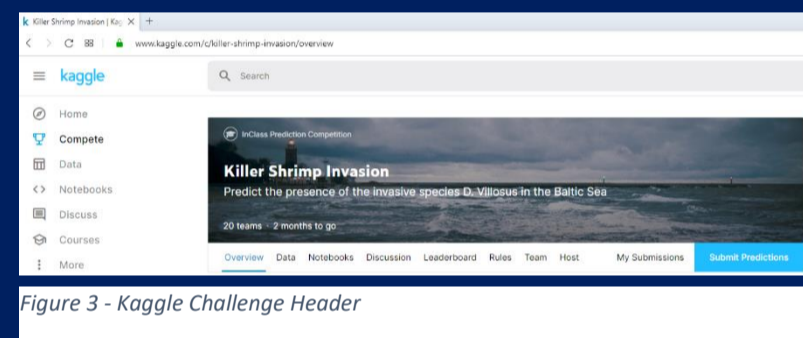


Figure 3 - Kaggle Challenge Header

Rank	Team Name	Notebooks	Team Members	Score	Entries	Last
1	Killer Shrimp Withdrawal			1.00000	1	1d
2	Rookiet27			1.00000	3	5mo
3	姜群尧 (Jiang Qun Yao)			1.00000	4	5mo
4	Dmitry			0.99954	11	5mo
5	Curtis Thompson			0.99938	33	5mo
6	Chris X			0.99770	15	5mo
7	Fabio Chiusano			0.99758	10	1d
8	Fred Meyler			0.98811	2	1d
9	Google ArtificialMoron			0.99543	2	5mo
10	Max Egorov			0.98527	26	10mo
11	kristaut			0.99384	13	10mo
12	katou1110			0.98996	7	5mo
13	LucasLau			0.98938	2	5mo
14	Two Pebbles			0.98515	33	1d
15	Alex Mina			0.98021	23	1d
16	naturina			0.97923	1	10mo
17	Ivan Z			0.97281	4	1d
18	Shrimp-Killa			0.96021	1	10mo
19	Ramón Echeverría			0.95175	1	1d
20	AndreaBianchi			0.90805	4	1d
21	Oskar Nykvist			0.90213	2	5mo
22	Andrew Tran			0.83277	9	5mo
23	Bastien1234			0.74889	1	1d

Figure 4- Overview of Participant Leaderboard



Check out the challenge yourself!