

Automatic interpretation of cod otoliths using deep learning

Endre Moen^{1*}, Rune Vabø¹, Come Denechaud¹, Ketil Malde^{1,2},

1 Institute of Marine Research, Bergen, Norway

2 Department of Informatics, University of Bergen, Norway

* endre.moen@hi.no

Abstract

The age of individual cod (*Gadus morhua*) is determined by manually examining the layered structure of otoliths, a calcium carbonate structure of the inner ear.

Image-based methods have been tried to age otoliths with varying results, but recent developments in automatic image analysis techniques are promising. The objective of this paper is to investigate the accuracy in aging broken otolith images on state-of-the-art convolutional neural networks.

Introduction

Information on fish age constitutes one of the most important biological variables, which is used in the investigations of life history (e.g. growth, sexual maturation) and population dynamics Campana [2001]

– The need for more Automated analysis. On the future of data analysis. – On the importance of determining age distributions for fisheries and ecosystem management. – How Cod otoliths are measured manually. The importance of otoliths for age determination of various species. Specifics about cod. – Related work. Halibut otoliths. Salmon scales. The Greeks. Knowledge of fish age structure is central to the study of fish and stock dynamics. It informs on population growth and mortality and, with size distribution, is one of the main criteria used for determining the health of exploited

populations and monitoring the effects of selective fishing (Hidalgo et al., 2011; Brunel and Piet, 2013). Changes in the age distribution can track significant changes in population structure, such as a particularly strong year-class skewing the distribution (Reglero and Mosegaard, 2006), or the gradual truncation of older age classes as selective fishing mortality removes larger individuals (Siskey et al., 2016). Hard structures such as scales and otoliths are used worldwide as one of the primary sources of fish age estimates, due to their ability as natural physiological and environmental recorders to form regular, temporally resolved growth increments at the daily and annual levels (Campana, 2001; Francis and Campana, 2011; Albuquerque et al., 2019). While age is inferred from the “simple” counting of annual increments, the interpretation of this zonation pattern is species or even population-specific (Høie et al., 2009) and is based on precise knowledge of the timing of zone formation and of the correct identification of true and false zones (Panfili et al., 2002). This process therefore requires specific expertise and is subject to uncertainties in both between-reader precision and “true” age accuracy (Francis and Campana, 2011). Because those estimates are central to stock assessment, ageing errors or wrong interpretation of otolith zonation can have dramatic effects on the evaluation of fish biology and consequently stock size and structure (Tyler et al., 1989; Beamish and McFarlane, 1995; Ragonese, 2018). Otolith reading is also time and resource consuming. Training of expert readers can take up to several years depending on the species, and otoliths often undergo a long processing phase before the final age estimates can be produced (Carbonara and Follesa, 2019). This is particularly true for demersal fish species, like Atlantic cod (*Gadus morhua*), that have large opaque otoliths that can't be read whole and need to be prepared. These routines vary between populations and institutes and range from direct reading of broken otoliths under a magnifying glass, to embedding, thin sectioning and finally imaging of the sections under a microscope. There has been a variety of methods proposed to automatically interpret otoliths, which range from one-dimensional data analysis like intensity transects (Mahé, 2009) to the more recent effort toward developing machine learning (ML) frameworks (Moen et al., 2018; Politikos et al., 2021). Despite fast progress the results remain mixed and often yield lower precision and consistency than those obtained by trained human readers, which limits the application of automated methods in real conditions. However, one aspect that is often under considered by such studies

are the practical time and cost benefits that implementing a functional ML framework
51 would provide. As noted by Fisher and Hunter (2018) in their review of digital
52 techniques for otolith analysis, “costs for human and machine ageing systems are
53 broadly similar since a large part of the cost is associated with preparing the otolith
54 sections”. As such, the net benefit of automated ageing routines is directly dependent
55 on the ability to scale performance using a comparatively smaller number of samples
56 than human readers or, alternatively, to train them on “rougher” data that can be
57 produced faster and at a more efficient cost. MORE ON CNN ETC? In this study, we
58 develop a deep learning network for estimating Atlantic cod age using multi-exposure
59 images of broken otoliths set in place using simple plasticine. More on methods. Our
60 results are positive and show the potential for developing automated pipelines that
61 require minimum processing and could be able to produce near at-sea age estimates.
62

Method and materials

Data Collection

We sampled a data set of 5150 cod otolith images, which has been collected on different
63 cruises and read by otolith experts. The images are taken during cruises in the period
64 2012-2018 conducted by Institute of Marine Research (IMR). There are six images with
65 three light exposures and one rotation. The expert readers has varied during this time
66 period as has the configuration for photographing the otoliths.
67

Figure 1. Otolith from 2016, read age: 6. With light exposure: medium, low, high,
68 then rotated 180 degrees and three new images



The images is of size 3744 x 5616 which are re-scaled for training to between 384x384
70 to 512x512. The image light exposure varies depending on light condition outside, and
71

are stored in the property 'ExposureTime' of the JPG file. Typically the exposure order
 72 is middle, dark, or light then a rotation of 180 degrees, and then middle, light, dark
 73 again. But the order might change, and the given order is recovered by reading the
 74 metadata property of the jpeg and sorting the exposure time.
 75

The otoliths are prepared for imaging by breaking them. The process also involves a
 76 camera setup, a folder structure referencing age, survey and station number, lighting
 77 setup, mounting of camera, and finally camera capture. More information about this
 78 process can be found in S.C. Myers [2019]
 79

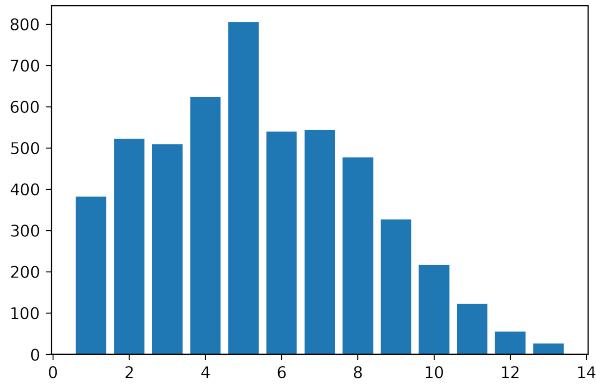


Figure 2. Age distribution of all 5150 images

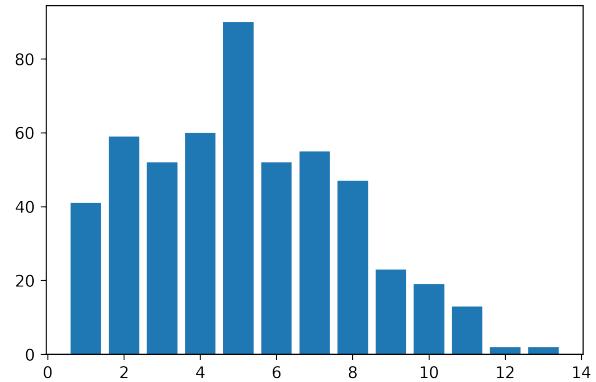


Figure 3. Age distribution of 515 images from the test set

Convolutional neural network architecture

There are two families of models used, EfficientNet B4-B6 (Tan and Le, 2019) and
 80 EfficientNet V2 medium, and Large (Tan and Le, 2019)
 81

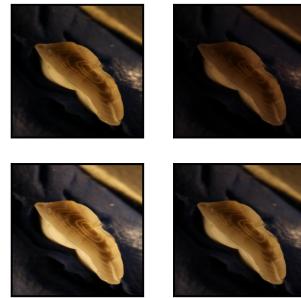
Each CNN was trained using transfer learning by loading ImageNet (Deng et al.,
 82

2009) weights. The image size varied between 380x380 and 528x528 for EfficientNet Bx,
84 and 384x384 for EfficientNetV2 size medium, large. While test-set size prediction has
85 been done both on 384x384 and larger resolutions 480x480 and 512x512 as described in
86 the paper. To investigate the image-taking protocol described in (S.C. Myers, 2019) we
87 have also training on 9-channel images. Three images are stacked to produce a
88 9-channel image. Using Timm(Wightman, 2019) the imageNet weights were duplicated
89 on the input layer to accommodate 9 channels. The 3 images used are of dark, medium
90 and light exposure of the first orientation.
91

We apply state-of-the-art CNNs based on performance on the ImageNet benchmark.
92 The imageNet benchmark has a classification while we treat ageing the cod-otoliths as
93 a regression problem. So the last layer of the CNNs has been modified to output a
94 linear output. In the EfficientNetV2 family we have done this by applying multilayer
95 perceptron layers going from 1280 output of last hidden layer to dense 256-layer, then a
96 leakyRelu (Xu et al., 2015) layer, then dense 32-layer, then a leakyRelu layer, and
97 finally linear output layer. While in EfficientNetV1 we only change the softmax layer to
98 a linear out layer.
99

To each fold we apply Standard scalar, which Standardize the age by removing the
100 mean and scaling to unit variance, on the training set. The standard scalar is then
101 applied to validation and test-set. To get the test-set predictions we inverse transform
102 the standard scalar.
103

Figure 4. Otolith from 2013, read age: 6. With light exposure: medium, low, high,
and expectation per channel of the three exposures.



Implementation and training

EfficientNetV1 B4, B5, and B6 was implemented with tensorflow (Abadi et al., 2016) and keras (Chollet and others, 2018) software packages in python. Computation was done using CUDA 11.1 and CuDNN with Nvidia(Nvidia Corp., Santa Clara, California) A6000 accelerator card with 48 GB of GPU memory, EfficientNetV2 medium, and large was implemented with the pytorch (?) and TIMM citeprw2019timm software package. Computation was done on P100 cards with 12 GB of GPU memory. Pretrained weights for EfficientNetV1 was available through Keras, and pretrained weights for EfficientNetV2 was available through Timm.

Augmentation was applied to the training-set. The images were augmented using rotation between 0 and 360 degrees, and reflection by the vertical axis. The pixel values has a range between 0 and 255 which was normalized to between 0 and 1. No other augmentation techniques was used like cropping, shifting or shearing as it can result in loss of age structure information.

The augmented data set can produce $360*2*5150 = 3.708.000$ possible images. Depending on the augmentation factor and the number of images in a training cycle, the model will likely never see the same image twice.

The primary metric used for training the models is mean squared error (MSE) while the primary metric used for evaluating the models and comparing it to expert readers is accuracy. Accuracy is obtained by rounding the floating point number predictions to nearest integer and comparing the age classification against the true labels. To reach human level accuracy a score of 85% or higher is required (?).

To get the most out of a small data-set we applied 10-fold cross-validation on 90% of the data-set, 4635 images. Each fold of the 10 folds consists of 90% of the cross-validation set and 81% of the whole data-set, 4172 images for training. Each fold had then 463 images for validation which is 10% of the cross-validation set, and 9% of the whole data-set. Each model is training on the 4172 images and the model with the best MSE on the 463 images in the validation set is chosen. The best model on the validation set was then used to predict the age on the test-set, and the metric for accuracy and MSE was recorded. The test-set is chosen at random, while the 10-fold split is chosen using stratified-kfold split which preserves the distribution of the whole

cross-validation set in each validation set. So the 463 images while have similar age 135
distribution to that of the 4635 images in the cross-validation set. Both the test-set and 136
the whole data-set follows a normal distribution with larges age-class, 5 year old, but 137
with deviation in frequency in other age-classes. 138

The CNN hyper-parameters configurations varies a little between the two families of 139
networks, but are kept the same within the families. Some hyper-parameters that has 140
been tuned are batch size, learning rate, k-fold size, weight decay, step size, number of 141
epochs, early stopping, and patience. To keep track of all these parameters we wrote the 142
configuration to a JSON file, config.json. The configuration file has been written for 143
each model trained, and the exact configuration can be found in the results section of 144
the github page of this project 145

(<https://github.com/emoen/Deep-learning-for-regression-of-cod-otoliths>) 146

As previously mentioned, we trained 10 models using 10-fold cross-validation which 147
resulted in 10 predictions on the test-set. We then produce an ensemble prediction of 148
the 10 models. Typically the ensemble prediction is better than any single fold 149
prediction. Ensembles are better because they improve performance. An ensemble can 150
make better predictions and achieve better performance than any single contributing 151
model, just as more experts will produce higher accuracy in predicting a single otolith. 152
Robustness; An ensemble reduces the spread or dispersion of the predictions and model 153
performance. This result can be improved further by taking ensemble predictions of 154
ensembles. We look at all ensembles from tuple-ensembles, consisting of 2 models, which 155
produces an ensemble of 20 models, and triplet-ensembles consisting of 3 models, to 156
ensemble of all models which produces an ensemble consisting of 180 models. 157

By choosing the best model we are over fitting to the test-set, but selecting a subset 158
of the best of these ensembles should produce a candidate ensemble of ensemble which 159
will produce the best prediction on a hold-out test-set. 160

Results

We have conducted a series of experiments on the EfficientNet family of CNNs, with 162
different hyper-parameters and on images with light exposures from the set light, 163
medium, dark. Training has been done on the 10-fold cross validation set which 164

produced 10 models.

In table 1 and table 2 are the accuracy and MSE metrics for the ensembled predictions from the 10 fold training. It can be observed that in the efficientNetV1 family, larger networks has better MSE, while accuracy is more fluctuating. A similar pattern can be observed for the efficientNetV2 networks. However it seems like efficientNetV1 is better than V2 in both metrics unlike the results observed on ImageNet.

Table 1. Accuracy by light exposure and CNN architectures

MSE:light/CNN	B4	B5	B6	Medium	Large	Xtra Large
min	72.8	74.4	73.4	67.0*	-	-
medium	-	-	74.4	72.4	71.8	-
max	-	-	-	-	-	-
9 channels	-	-	-	-	71.7	-

Table 2. MSE by light exposure and CNN architectures

ACC:light/CNN	B4	B5	B6	Medium	Large	Xtra Large
min	.277	.277	.272	.331*	-	-
medium	-	-	.262	.292	.280	-
max	-	-	-	-	-	-
9 channels	-	-	-	-	.281	-

We compare the 10 fold prediction accuracy of all the models in a box plot in figure 6, and for MSE in 11. The red line is the ensemble accuracy or MSE. The ensemble metric is either better than all the folds or in the upper quantile.

By comparing the models on MSE we can see that larger models are better, e.g B6 has higher mean than B5 and B4, and large is better than medium. We also see that the EfficientNetV2 networks has higher mean than the first generation EfficientNet. However, this is not true for the ensemble predictions (red line) nor for the fold-mean or ensemble of the accuracy. We can also see that the effect of adding 3 images, creating 9 channels, on the model is that the variance is reduced, the fold mean metric increases, but the ensemble metric is reduced.

The box plots are produced from the folds given in table 3 and 4.

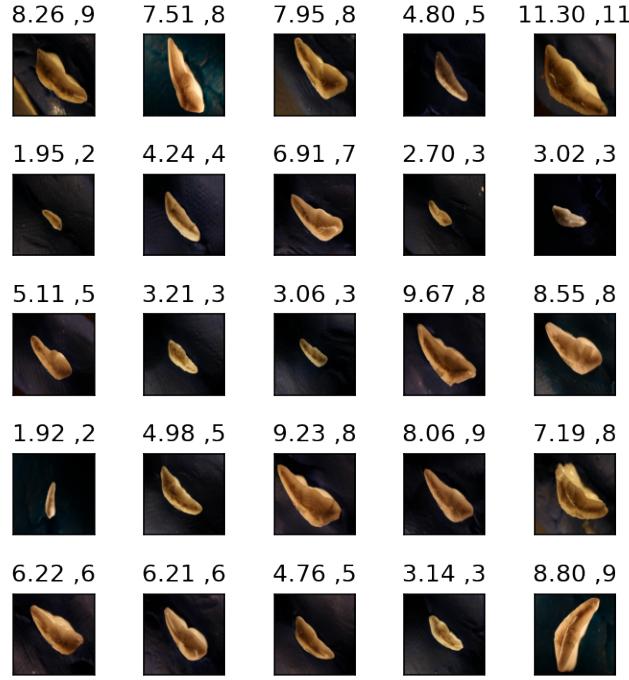


Figure 5. Sample of 25 predictions on a fold of training on EfficientNetV2 size medium with minimum light exposure, left number is prediction, and right number is age read

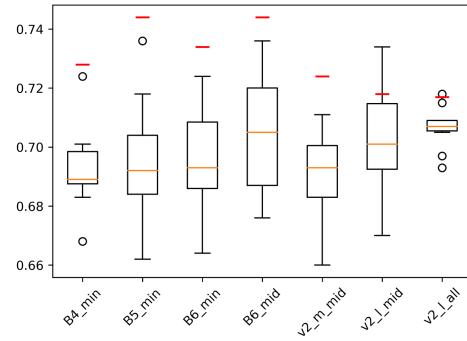


Figure 6. Accuracy score of 5 models and red line is ensemble prediction accuracy

Prediction by age class and residuals

The figures below shows the predictions per age group on the test-set. We can see that the prediction follows a linear trend $y = x$ except for the 2-3 last years, when the mean drops below $y = x$. This is even more obvious in the residual plots where the prediction drops below $y = 0$ for the last 2-3 age groups.

Figure 9 shows scatter plots of all predictions that results in a misclassification.

That is predictions that error greater than 0.5 in magnitude. Predictions that miss by

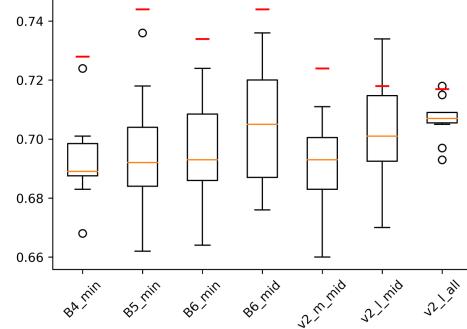


Figure 7. MSE score of 5 models and red line is ensemble prediction MSE

Table 3. MSE per CNN and per fold

CNN/fold	1	2	3	4	5	6	7	8	9	10	mse
B4, min	.320	.318	.306	.313	.322	.314	.315	.316	.306	.302	.277
B4, middle											
B4, max											
B5, min	.324	.322	.325	.336	.291	.314	.320	.331	.33	.317	.277
B5, middle											
B5, max											
B6, min	.325	.329	.334	.293	.312	.290	.320	.3.	.276	.306	.272
B6, middle	.323	.301	.312	.268	.294	.266	.309	.311	.278	.289	.262
B6, max											
medium, min											
med., mid.	.321	.377	.332	.285	.285	.325	.311	.348	.295	.373	.292
medium, max											
medium, all											
large, min											
large, middle	.301	.281	.299	.318	.282	.305	.280	.334	.3	.310	.280
large, max											
large, all	.292	.289	.289	.326	.307	.327	.283	.30	.335	.295	.281
xl, min											
xl, middle											
xl, max											
xl, all											

more than 1.5 in magnitude are shown with red dots.

192

Ensemble of ensembles

193

We search the space of ensemble of ensemble predictions which are given by $\sum_{k=1}^N \binom{N}{k}$ 194 where $N = 22$ and $k \in 1..N$ and find three ensemble of ensembles which produce the 195 best results overall with accuracy of 75.9%, 76.1%, and 76.9% and MSE 0.247, 0.248, 196 and 0.248 from ensemble of all networks, ensemble of B4, B5 and B6 with min exposure, 197

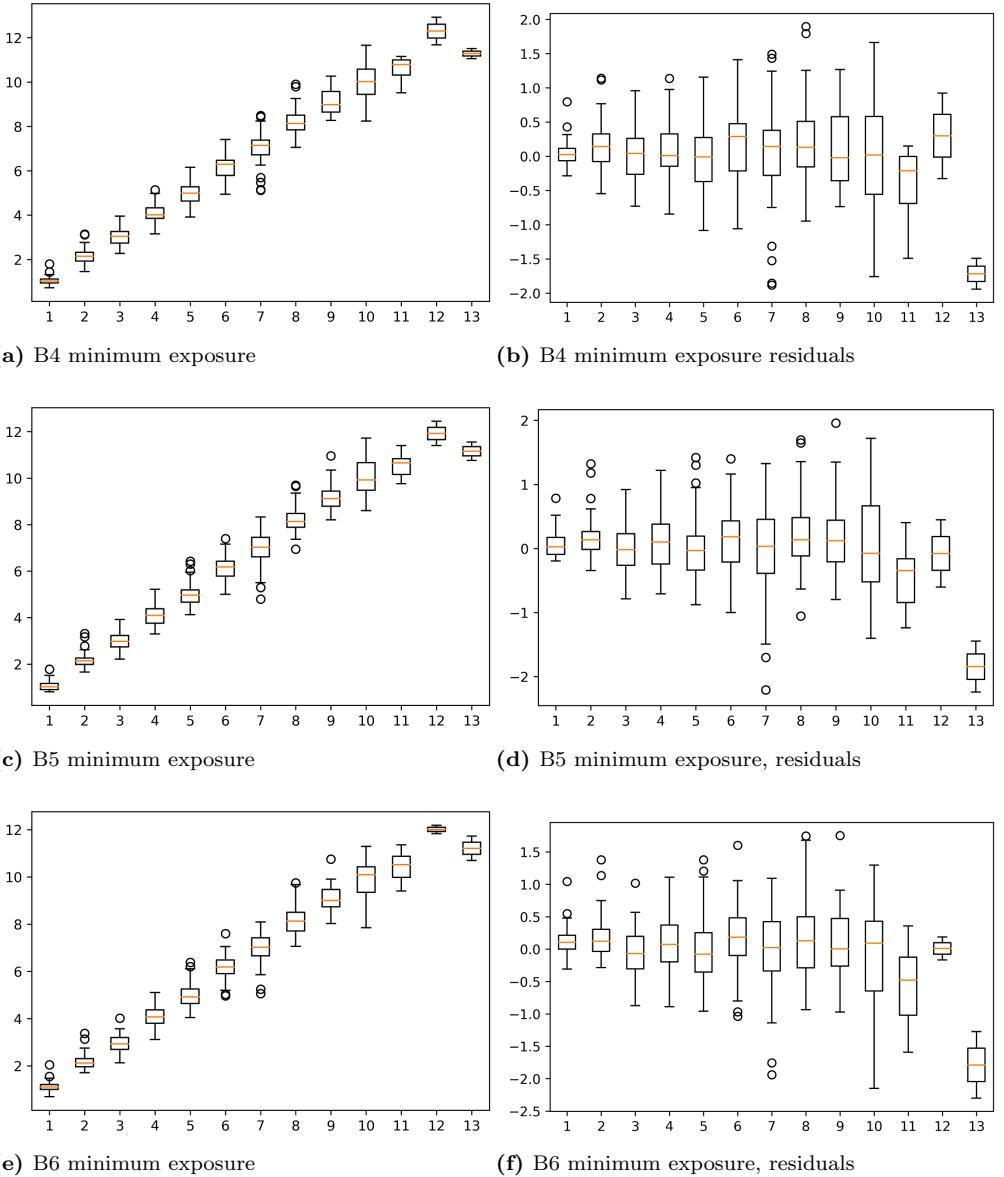
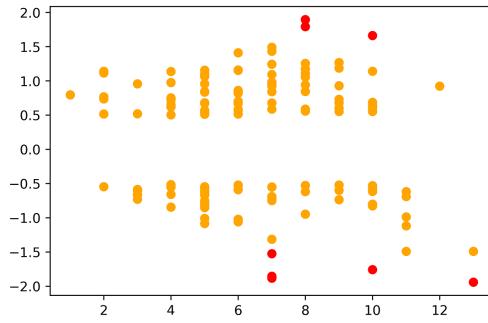
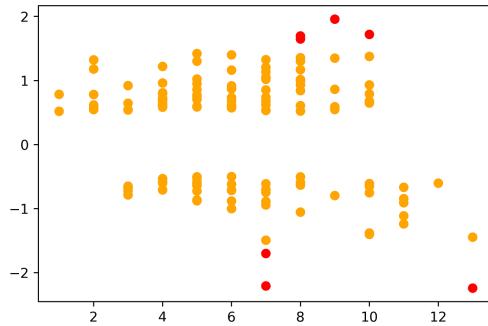


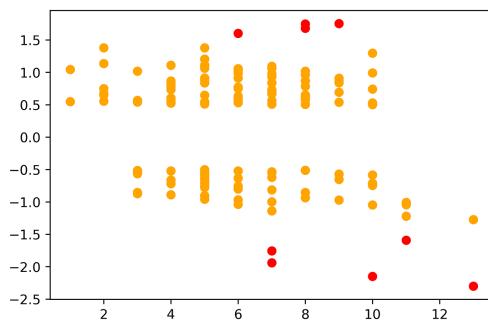
Figure 8. Comparing the models, looking at age per age class, and the residuals per prediction



(a) B4 minimum exposure



(b) B5 minimum exposure



(c) B6 minimum exposure

Figure 9. Comparing the models, looking at age per age class, and the residuals per prediction

Table 4. Accuracy per CNN and per fold

CNN/fold	1	2	3	4	5	6	7	8	9	10	acc
B4, min	69.9	68.9	68.7	68.3	68.9	70.1	69.7	66.8	68.9	72.4	72.8
B4, middle											
B4, max											
B5, min	71.8	69.1	69.3	66.8	73.6	70.7	66.2	68.3	69.5	68.7	74.4
B5, middle											
B5, max											
B6, min											
B6, middle	68.3	68.5	66.4	72.4	70.7	70.9	69.3	69.3	72	68.9	73.4
B6, max	68.5	69.9	67.6	73.6	72.8	72	68	69.3	72	71.1	74.4
medium, min											
med., mid.	68.7	67.6	68.3	71.1	70.1	70.5	69.9	68.3	69.9	66	72.4
medium, max											
medium, all											
large, min											
large, middle	69.7	73.4	69.1	67	71.8	69.9	72.6	68.2	70.5	70.3	71.8
large, max											
large, all	70.9	70.7	70.5	70.7	71.5	69.3	70.7	71.8	69.7	70.9	71.7
xl, min											
xl, middle											
xl, max											
xl, all											

and ensemble of B4, B5, B6 and middle with min exposure.

198

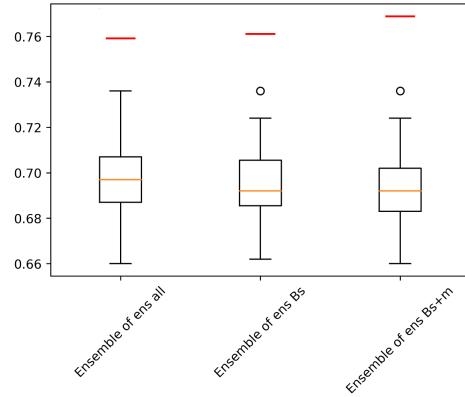

Figure 10. Ensemble of ensemble:
accuracy of the 3 best models

Table 5. Accuracy/MSE pr ensemble of ensemble. Eoe1 is ensemble of ensemble of all models, Eoe2 is for B4, B5 and B6, and Eoe3 is Eoe2 plus efficientNetV2 medium.

score/ensemble	eoe1	eoe2	eoe3
Acc	75.9	76.1	76.9
MSE	.247	.248	.248

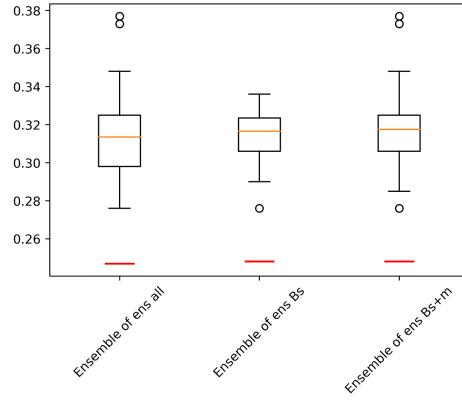


Figure 11. Ensemble of ensemble: mse of the 3 best models

199

Outliers

200

Looking at figure ?? we can see that the model under predicts the age of older otoliths. 201
 This pattern is especially observable for individuals read as 15 years and older. The 202
 oldest predication is 18 years while the test set contains individuals as old as 22 years. 203
 To better understand the bias, figure 11 shows the 4 largest outliers from the test set 204
 which come from two pairs 205

205

Table 6. Outliers with more than 1.5 year error. Index of image in test-set per model

V2-m,mid.	V2-m,mid.	V2-l,all	V2-l,mid.	B4,min	B5,min	B6,min	B6,mid.
				13	13	13	13
						48	
71	71	71	71	71	71	71	71
92	92						
				270	270		270
279	279	279	279	279	279	279	279
		312	312				
			320	320			
362	362	362	362	362	362	362	362
342	342	342	342	342	342	342	342
369	369	369	369	369		369	369
			393			393	393
423	423	423	423				
					444		
						502	502
7	7	7	9	8	7	9	9

Figure 9 are the most commonly misclassified images with greatest magnitude of

206

Table 7. Outliers with more than 1.5 year error. Prediction and true age, per model

Idx	V2-m,mid.	V2-l,all	V2-l,mid.	B4,min	B5,min	B6,min	B6,mid.	Age
13				9.79	9.64	9.74	9.58	8
48						7.6		6
71	4.96	4.98	4.94	5.14	4.79	5.06	5.12	7
92	10.95							13
270				11.66	11.71		11.53	10
279	9.93	9.79	9.75	9.89	9.69	9.67	9.7	8
312		9.42	9.38					11
320			5.44	5.47				7
362	5.11	5.14	5.23	5.11	5.29	5.24	5.15	7
342	10.35	10.6	10.61	11.05	10.75	10.69	10.84	13
369	8.17	8.13	8.23	8.24		7.85	8.29	10
393			10.53			10.75	10.83	9
423	5.39	5.69	5.43					8
444					10.95			9
502						9.4	9.43	11

Figure 12. The most common images miss predicted with more than 1.5 years



error.

207

Discussion

208

During initial training we trained a B4 network on ca 2000 images and obtained an accuracy of ca 60%, later another 3000 images was added and the same network was trained on ca 5000 images which resulted in accuracy of ca 70%. It could be interesting to investigating if adding another 3-5000 images would increase the accuracy to 80%.

209

210

211

References

213

References

214

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S.,
Davis, A., Dean, J., Devin, M., et al. (2016). Tensorflow: Large-scale machine
learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467.* 215
216
217
- Campana, S. (2001). Accuracy, precision and quality control in age determination,
including a review of the use and abuse of age validation methods. *Journal of fish
biology*, 59(2):197–242. 218
219
220
- Chollet, F. and others (2018). Keras 2.1.3. [https://github.com/fchollet/keras.](https://github.com/fchollet/keras) 221
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: A
large-scale hierarchical image database. In *Proceedings of IEEE Conference on
Computer Vision and Pattern Recognition*, pages 248–255. IEEE. 222
223
224
- S.C. Myers, A. Thorsen, J. G. K. M. N. H. (2019). An efficient protocol and data set for
automated otolith image analysis. *GeoScience Data Journal.* 225
226
- Tan, M. and Le, Q. V. (2019). Efficientnet: Rethinking model scaling for convolutional
neural networks. *CoRR*, abs/1905.11946. 227
228
- Wightman, R. (2019). Pytorch image models.
[https://github.com/rwightman/pytorch-image-models.](https://github.com/rwightman/pytorch-image-models) 229
230
- Xu, B., Wang, N., Chen, T., and Li, M. (2015). Empirical evaluation of rectified
activations in convolutional network. *CoRR*, abs/1505.00853. 231
232