

## YMP re-training on 2016-2021 data

Yield Map Prognosis

Exported on 12/15/2021

## Table of Contents

1 Cleaning of yield maps (YMCA) .....	4
2 Fetch satellite images (YMP) .....	5
3 Retraining of YMP .....	6
4 Overall conclusion .....	9
5 Future work .....	10

## Table of Content

- [Cleaning of yield maps \(YMCA\)\(see page 4\)](#)
- [Fetch satellite images \(YMP\)\(see page 5\)](#)
- [Retraining of YMP\(see page 6\)](#)
- [Overall conclusion\(see page 9\)](#)
- [Future work\(see page 10\)](#)

# 1 Cleaning of yield maps (YMCA)

- It was complicated to use the code from 2019, as it relied on a different data structure and code-environment.
- There is no matching polygon/yield map sets in the timespan 2018-2020 (both included) resulting in no useful data for this period, for 2021 only 7 unique yield maps was found.
- Use of field polygons 'Dansk Mark Database' (DMDB) instead of 'Internet Mark Kort' (IMK) resulted in a vastly different amount of matching polygon/yield map sets, even when adding the 7 new yield maps (100 unique sets in 2019, and 77 now)
- The location of the fields in the old and new set of yield maps also differs, as the new yield maps also covers the area of Randers and Ringkøbing, as seen the these plots:



Old yield maps:



- New yield maps:

- **Conclusion:** For this pipeline to be automated, and to ensure an easy and reliable way to retrain and expand the data, this solution needs a lot of work. We strongly recommend a full restructure on a more automated pipeline.
  - Ensure all needed field polygons exists in DMDB - else add the field polygons from IMK to DMDB, such that we can still rely on our SEGES database, instead of an external data source.
  - Heighten the consistency of the yield data, a way to achieve this is by using the ISOXML files from CropManager instead of the current shape files (manually received from Danish farmers).
    - The developer time spend on normalizing the inconsistencies in the data, makes the current setup unfeasible for reuse in coming years.
  - Our opinion is that YMCA should become a data product, owned and maintained by a dedicated data-team, thus avoiding the technical debt which prevented improvements in the current rerun.

## 2 Fetch satellite images (YMP)

- It was (again) complicated to use the code from 2019, as it relied on a different data structure and code-environment.
- Instead of downloading the satellite images for each rerun, we should rely on all ready downloaded images, stored in an easily accessible way - e.g. the 'sat-lake'-solution, which pre-cashes all sentinel 2 images in .zarr format on Azure Data Lake.

### 3 Retraining of YMP

- It was complicated to use the code from 2019, as it relied on a different data structure and code-environment.
- Model score results:

	OLD - trained on 2016-2017 yield maps				NEW - trained on 2016-2021 yield maps				Difference				
Field overall scores	April 10 <sup>1</sup>	May 10 <sup>2</sup>	June 5 <sup>3</sup>	August 1 <sup>4</sup>	April 10 <sup>5</sup>	May 10 <sup>6</sup>	June 5 <sup>7</sup>	August 1 <sup>8</sup>		April 10	May 10	June 5	August 1
train_field_mean_mae_kg/ha	327.9	333.9	290.8	283.1	354.9	362.7	331.2	223.7		27,00	28,80	40,40	-59,40
train_field_mean_std_kg/ha	451.9	454.3	419.7	381.4	472.4	472.6	458.4	301		20,50	18,30	38,70	-80,40
train_field_mean_samples	79	79	79	79	61	61	61	61		-18	-18	-18	-18
val_field_mean_mae_kg/ha	835.7	926.3	962.8	717.5	648.2	747	608.6	674.7		-187,50	-179,30	-354,20	-42,80
val_field_mean_std_kg/ha	1127	1220.7	1240.1	863.7	879	969.5	799.9	859.4		-248,00	-251,20	-440,20	-4,30
val_field_mean_samples	22	22	22	22	15	15	15	15		-7	-7	-7	-7
Position overall scores													

1 <http://localhost:5000/#/experiments/4/runs/36759b7dea8045a1ab8ab1ed2b6bc427>

2 <http://localhost:5000/#/experiments/4/runs/8e24b64ef9194a10aeaa776f1d5cc6ae>

3 <http://localhost:5000/#/experiments/4/runs/a16a72324a5849819ae37f4e5e7e04ee>

4 <http://localhost:5000/#/experiments/4/runs/8616813e09c64e44a6e8e3ea3a566f56>

5 <http://localhost:5000/#/experiments/4/runs/642826c95b66489db8199587bd920d94>

6 <http://localhost:5000/#/experiments/4/runs/b840288d70a34d8c8e8299dc40d0c130>

7 <http://localhost:5000/#/experiments/4/runs/cf7a340a973e4accb98a3f1951f6aa7b>

8 <http://localhost:5000/#/experiments/4/runs/9673b10a2eaf43a9b17a7b576a57039d>

<b>train_position_mae_kg/ha</b>	1410.6	1359.8	1295.5	1223.6	1410.5	1349.9	1283.6	1199.2		-0,10	-9,90	-11,90	-24,40
<b>train_position_samples</b>	383594	383594	383594	383594	295733	295733	295733	295733		-87861	-87861	-87861	-87861
<b>val_position_mae_kg/ha</b>	1729.6	1714.4	1707.5	1473.7	1462.2	1552.2	1327.5	1271.9		-267,40	-162,20	-380,00	-201,80
<b>val_position_samples</b>	81880	81880	81880	81880	77411	77411	77411	77411		-4469	-4469	-4469	-4469

- We do not see significance changes in any (i.e. field-level or position-level) of the training performance scores. The small changes are likely to be a coincidence, given the changes in the dataset.
- However, both field-level and position-level validation Mean Absolute Error (MAE) scores are improved significantly (~200kg/ha) for all prediction dates.
- But due to the smaller sample size we decided to perform an ad hoc cross-validation to conclude if other train/val splits results in the same model improvements.
- 5-fold cross-validations results of val\_position\_mae\_kg/ha (see notebook [YMP-678\\_crossvalidation\\_of\\_2016-2021\\_yield\\_data.ipynb](https://bitbucket.seges.dk/projects/DDS/repos/yield_map_prognosis/browse/data_analyses/YMP-678_crossvalidation_of_2016-2021_yield_data.ipynb)<sup>9</sup>):

	<b>Support (# pixels)</b>	<b>04-10_MAE</b>	<b>04-10_STD</b>	<b>05-10_MAE</b>	<b>05-10_STD</b>	<b>06-05_MAE</b>	<b>06-05_STD</b>	<b>08-01_MAE</b>	<b>08-01_STD</b>
<b>fold_num_1</b>	108092	1715.81	2169.07	1745.27	2298.67	1617.63	2143.24	1369.40	1907.88
<b>fold_num_2</b>	178274	1558.85	2063.96	1493.90	2028.16	1578.41	2035.51	1505.94	1991.16
<b>fold_num_3</b>	155436	1726.96	2337.27	1722.46	2322.18	1727.66	2307.62	1461.70	1989.68
<b>fold_num_4</b>	190402	2901.53	3368.22	2791.30	3317.64	2762.46	3348.74	2449.48	3166.34
<b>fold_num_5</b>	114084	1742.17	2303.13	1708.38	2178.57	1594.68	2100.75	1458.95	1902.38

<sup>9</sup>[https://bitbucket.seges.dk/projects/DDS/repos/yield\\_map\\_prognosis/browse/data\\_analyses/YMP-678\\_crossvalidation\\_of\\_2016-2021\\_yield\\_data/YMP-678\\_crossvalidation\\_of\\_2016-2021\\_yield\\_data.ipynb?at=refs%2Fheads%2FYMP-678-cross-validation-notebook-on-new-ymp-models-with-2021-data](https://bitbucket.seges.dk/projects/DDS/repos/yield_map_prognosis/browse/data_analyses/YMP-678_crossvalidation_of_2016-2021_yield_data/YMP-678_crossvalidation_of_2016-2021_yield_data.ipynb?at=refs%2Fheads%2FYMP-678-cross-validation-notebook-on-new-ymp-models-with-2021-data)

<b>Mean over folds</b>	---	<b>1929.06</b>	<b>2448.33</b>	<b>1892.26</b>	<b>2429.04</b>	<b>1856.17</b>	<b>2387.17</b>	<b>1649.09</b>	<b>2191.49</b>
<b>Scores over total pixels</b>	---	<b>1987.18</b>	<b>2678.49</b>	<b>1941.71</b>	<b>2647.04</b>	<b>1919.75</b>	<b>2623.27</b>	<b>1710.49</b>	<b>2397.65</b>
<b>Scores over yield maps</b>	---	<b>1023.69</b>	<b>1555.24</b>	<b>1138.50</b>	<b>1607.36</b>	<b>1177.33</b>	<b>1705.19</b>	<b>902.73</b>	<b>1359.97</b>

- We see the MAE for each fold is higher than for the previous trained model, thus the mean over the folds are also significant higher (300-500 kg/ha higher). Thus the specific train/val split used previously benefitted the model, but its performance is not generalizable.
- The standard deviation (STD) of each prediction error, shows that 68% of all validation samples has less than +/-2300kg/ha in prediction error.
- However, our cross validation still shows that more data (closer to harvest) improves the model performance.
- **Conclusion:**
  - As a larger, and more consistent, pool of data greatly benefits the model, we think this should be the main focus going forward.
    - A way to achieve this by ensuring all data fetching to the YMP pipeline comes from data products used across SEGES, i.e. yield data from CropManager.



## 4 Overall conclusion

- Our rerun of the YMCA pipeline resulted in fewer yield maps, i.e. 100 unique fields in 2019, and 77 now. However, the new fields also covers the area of Randers and Ringkøbing, i.e. more divergent soil/climate types.
- Our retrained YMP models on the new 77 yield maps, show models improvements of ~200kg/ha. However, following cross validations, the results were worse than our train/val model. This shows the dataset is not generalizable, thus resulting in better model-performance than expected in real life application.
- The smaller sample size is troubling, but due to the slightly improved train/val model results, we would recommend deploying the new models to production.

## 5 Future work

1. Implement automated pipeline for YMCA based on ISOXML yield maps.
2. Implement YMP training pipeline to fetch all geo-data from SEGES data cache store ([Spatio-Temporal Data Store Definition](https://confluence.seges.dk/display/YMCA/Spatio-Temporal+Data+Store+Definition)<sup>10</sup> like Sat-Lake)
3. Move YMP training/deployment pipeline to Azure Machine learning.
4. Implement cross-validation in to the YMP training pipeline.

---

<sup>10</sup> <https://confluence.seges.dk/display/YMCA/Spatio-Temporal+Data+Store+Definition>