

Disaggregation of energy system optimization models using machine learning for identification of active constraints

David Cardona-Vasquez^{a,b,*}, Alexander M. Konrad^a, Yannick Werner^a, Sonja Wogrin^{a,b}

^a*Institute of Electricity Economics and Energy Innovation (IEE), Graz University of Technology, Inffeldgasse 18, Graz, Austria*

^b*Research Center ENERGETIC, Graz University of Technology, Rechbauerstraße 12, Graz, Austria*

1. Appendix

In this Appendix, we provide an additional description of the model used to represent the Austrian power system for the test case and a more in-depth analysis of the machine learning training process; finally, we expand the results concerning the comparison between the complete and disaggregated optimization model results.

1.1. Representation of the Austrian Power System

In our representation of the Austrian power system, we use a 10-node transport model (i.e., pure network flow) with fixed transport costs. This representation allows us to isolate the effect of the time-linking constraints and leverage the Basis-Oriented approach's theoretical foundation. Concerning the imports and exports, and given the geographical location of Austria, it is important to take into account its interconnections with other countries. To achieve this, we divide the neighboring countries into net exporters to Austria (imports) and importers from Austria (exports), where the exporters are represented as a power plant with an installed capacity up to the interconnection's limit, and the importers are modeled as an additional node with a demand corresponding to the historical exports data.

We also subtract long-duration storage from the whole power system demand using historical data, as is commonly done in some ESOM. As mentioned in the main part of the manuscript, long-duration storage leads to extended linked periods,

*Corresponding author:

URL: david.cardonavasquez@tugraz.at (David Cardona-Vasquez)

which hinder the methodology's computational advantage. Dealing with this problem is part of our future research.

With these considerations in mind, the complete model formulation can be found in Eq (1) and Table 1 describes the symbols used in the model.

$$\begin{aligned}
& \min_{p_{g,k}, nsp_{i,k}, f_{k,i,j}} \sum_{k,g} C^g p_{g,k} + \sum_{k,i} C^{nsp} nsp_{i,k} + \sum_{k,i,j} C^N f_{k,j,i} & (1a) \\
& \text{s.t.} \quad \sum_j f_{k,j,i} - \sum_j f_{k,i,j} + nsp_{i,k} + \sum_{g \in i} p_{g,k} = D_{k,i} \quad \forall k, i & (1b) \\
& 0 \leq p_{w,k} \leq CF_k \bar{P}_g \quad \forall k, w & (1c) \\
& 0 \leq p_{t,k} \leq \bar{P}_t \quad \forall k, t & (1d) \\
& f_{k,i,j} \leq \bar{F}_l \quad \forall k, i, j & (1e) \\
& f_{k,j,i} \leq \bar{F}_l \quad \forall k, i, j & (1f) \\
& p_{t,k} - p_{t,k-1} \leq RU_t \quad \forall k, t & (1g) \\
& p_{t,k-1} - p_{t,k} \leq RD_t \quad \forall k, t & (1h) \\
& & (1i)
\end{aligned}$$

Table 1: Description of the symbols used in the energy system optimization model

Symbol	Type	Description
g	Index set	Represents all the generation resources
t	Index set	Represents the ramp-constrained generation resources
w	Index set	Represents the generation resources that are constrained by a capacity factor
i, j	Index set	Represent the nodes
k	Index set	Represents the time index
C^g, C^{nsp}, C^N	Input data	Represent, respectively, the cost coefficients of the energy production the non-supplied power at each, and the transport
\bar{P}_g, \bar{P}_g	Input data	Represents the upper limit for generation resources
\bar{F}_l	Input data	Represents the maximum transport capacity between the nodes
RU_t, RD_t	Input data	Represent, respectively, the ramp-up and ramp-down limits of ramp-constrained generation resources
$f_{k,i,j}$	Variable	Represents the flow from i to j at time k
$p_{g,k}$	Variable	Represents the generation of resource g at time k
$nsp_{i,k}$	Variable	Represents the non-supplied power at node i and time k

1.2. Classifier Training and Validation Details

In the classifier selection and training process we used a two-step approach, using a low-variance, high-bias classifier to identify relevant features, and then a high-variance, low-bias for the actual prediction. We did it this way because we wanted to use features which can be explained and interpreted from a modeler perspective, and not only based on the improvements they achieve on the accuracy of the model; the complete list of candidate features is presented in Table 2, the *Importance* column

corresponds to the absolute relative importance the logistic regression assigned to the feature.

Table 2: Feature importance according to the coefficients from the logistic regression

Feature Name	Importance	Feature Name	Importance
<i>ren_demand_ratio</i>	0.134	<i>net_demand</i>	0.020
<i>ren_demand_ratio_lag_1</i>	0.133	<i>net_demand_lag_3</i>	0.019
<i>ren_demand_ratio_lead_1</i>	0.092	<i>max_gen_lag_3</i>	0.018
<i>ren_demand_ratio_lag_2</i>	0.071	<i>max_gen</i>	0.015
<i>ren_demand_ratio_lag_3</i>	0.067	<i>max_gen_lag_2</i>	0.015
<i>demand_lag_2</i>	0.049	<i>net_demand_lead_2</i>	0.014
<i>max_gen_lead_3</i>	0.046	<i>ren_demand_ratio_lead_3</i>	0.013
<i>max_gen_lead_2</i>	0.036	<i>max_gen_exp_cap_ratio</i>	0.008
<i>net_demand_lag_2</i>	0.034	<i>max_gen_exp_cap_ratio_lag_1</i>	0.007
<i>max_gen_exp_cap_ratio_lead_1</i>	0.030	<i>demand_lag_1</i>	0.007
<i>max_gen_lead_1</i>	0.029	<i>net_demand_lead_1</i>	0.006
<i>demand</i>	0.027	<i>net_demand_lag_1</i>	0.005
<i>demand_lead_1</i>	0.027	<i>demand_lead_2</i>	0.005
<i>net_demand_lead_3</i>	0.026	<i>ren_demand_ratio_lead_2</i>	0.003
<i>demand_lag_3</i>	0.022	<i>max_gen_lag_1</i>	0.002
<i>demand_lead_3</i>	0.021		

For the low-variance classifier we used a logistic regression, a widely used and tested tool in data analysis, while for the actual predictions we tried several models: support vector with linear and radial basis kernels, nearest neighbors, Ada boost and random forests. The results show that the performance of Ada boost and random forests was basically indistinguishable from one another, nevertheless we chose random forests because they allow for easier parallelization, thus are faster to train in a multi-processor setup like one with GPUs.

Table 3: Performance of other classifiers in the testing portion of the 2022 dataset

Algorithm	Accuracy	Balanced Accuracy
<i>SVC linear</i>	0.81	0.68
<i>SVC deg.=2</i>	0.81	0.69
<i>SVC deg.=3</i>	0.81	0.67
<i>Random Forest</i>	0.93	0.89
<i>Ada Boost</i>	0.93	0.89

Finally, in Table 4 we present the confusion matrices in relative terms for the years 2015 through 2021 for the classifier without any re-training. As we can see, the

performance of the classifier deteriorates the further we go back in time as shown by the increase of false positives.

Table 4: Performance of the classifier for additional years without re-training

	Predicted							
	2021		2020		2019		2018	
	<i>Linked</i>	<i>Unlinked</i>	<i>Linked</i>	<i>Unlinked</i>	<i>Linked</i>	<i>Unlinked</i>	<i>Linked</i>	<i>Unlinked</i>
<i>Linked</i>	0.44	0.08	0.48	0.08	0.40	0.10	0.35	0.10
<i>Unlinked</i>	0.01	0.46	0.01	0.43	0.01	0.49	0.02	0.53
	2017		2016		2015			
	<i>Linked</i>	<i>Unlinked</i>	<i>Linked</i>	<i>Unlinked</i>	<i>Linked</i>	<i>Unlinked</i>		
<i>Linked</i>	0.40	0.09	0.38	0.09	0.20	0.16		
<i>Unlinked</i>	0.02	0.49	0.01	0.52	0.03	0.61		

1.3. Additional Results

Let us further discuss the solving times of the individual submodels. First, the submodels with a length longer than 1 hour have an average solving time of 0.07 s, with a maximum of 0.4 s for a submodel of length 70. As for the submodels of 1-hour length, the average solution time was 0.01 seconds. In Fig 1 we see that, as expected, submodels with longer time horizons take longer to solve, with an apparent linear relationship.

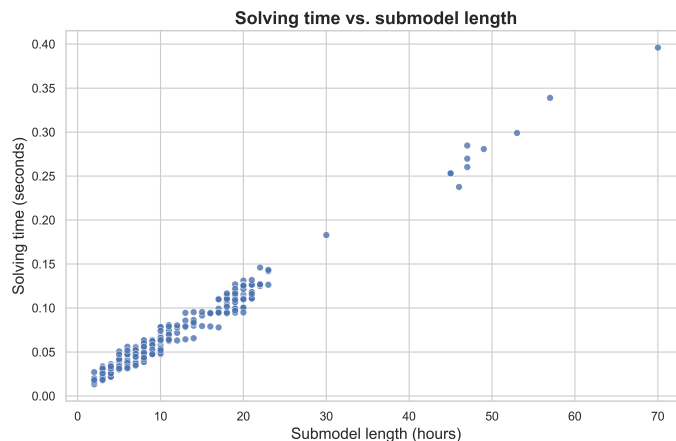


Figure 1: Solving time vs. submodel length, longer submodels took longer to solve

In Fig 2 we can see the distribution of solving time taking into account the submodel length; this result is consistent with the previous figure; however, in this

case we can also see how much more frequent smaller submodels are with respect to larger ones, this allows us to highlight the potential computational gains if one can correctly identify this submodels in the time index.

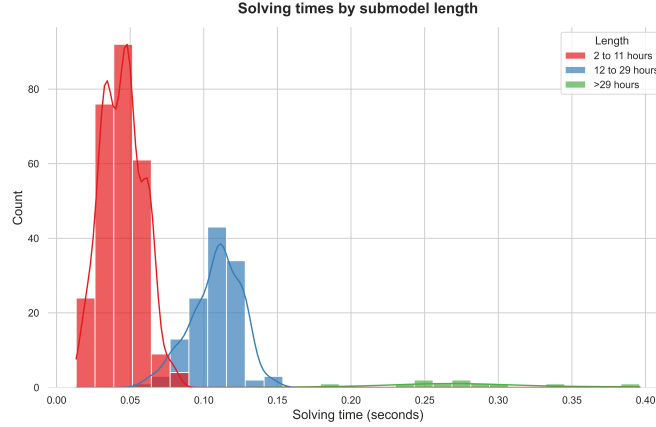


Figure 2: Distribution of solving time by model length, longer models are highly infrequent

Finally, in Table 5 we can see the aggregated production values for each of the generation resources in our representation of the Austrian power system. The values correspond to the validation year (2023) and the *Error* is calculated as the absolute value of the relative error.

Generation Resource	Aggregated Model (MWh)	Complete Model (MWh)	Error (%)
<i>BGLONSHORE</i>	2,378,723.68	2,378,723.68	-
<i>BGLPV</i>	732,351.16	732,351.16	-
<i>BGLROR</i>	5,649.07	5,659.77	0.19
<i>KTNONSHORE</i>	16,210.54	16,210.54	-
<i>KTNPV</i>	645,040.59	645,040.59	-
<i>KTNROR</i>	3,014,895.81	3,014,942.30	0.00
<i>KTNTKWCNG</i>	196,877.13	195,601.01	0.65
<i>NOECZWCNG</i>	-	-	-
<i>NOEONSHORE</i>	2,502,945.92	2,502,945.92	-
<i>NOEPV</i>	1,636,761.09	1,636,564.74	0.01
<i>NOEROR</i>	6,231,501.71	6,231,478.16	0.00
<i>NOETKWCNG</i>	869,256.29	874,173.70	0.56
<i>OOEDEWCNG</i>	-	-	-
<i>OOEONSHORE</i>	55,246.57	55,246.57	-
<i>OOEPV</i>	2,327,632.49	2,327,632.49	-
<i>OOEROR</i>	10,553,935.10	10,553,935.10	-
<i>OOETKWCNG</i>	57,496.88	57,861.45	0.63
<i>OTRONSHORE</i>	-	-	-
<i>OTRPV</i>	43,746.04	43,746.04	-
<i>OTRROR</i>	815,809.82	815,809.82	-
<i>STMKONSHORE</i>	366,935.34	366,935.34	-
<i>STMKPV</i>	1,446,430.86	1,446,785.26	0.02
<i>STMKROR</i>	3,229,830.64	3,230,062.87	0.01
<i>STMKTKWCNG</i>	5,432,464.97	5,424,945.25	0.14
<i>SZBGONSHORE</i>	-	-	-
<i>SZBGPV</i>	474,389.49	474,386.27	0.00
<i>SZBGROR</i>	1,706,265.09	1,706,288.63	0.00
<i>SZBGTKWCNG</i>	3,260.06	3,260.06	-
<i>TIRCHWCNG</i>	2,054.35	2,498.60	17.78
<i>TIRDEWCNG</i>	-	-	-
<i>TIRONSHORE</i>	-	-	-
<i>TIRPV</i>	408,857.22	408,857.22	-
<i>TIRROR</i>	2,831,743.16	2,831,743.16	-
<i>VBGCHWCNG</i>	-	-	-
<i>VBGDEWCNG</i>	-	-	-
<i>VBGONSHORE</i>	-	-	-
<i>VBGPV</i>	388,673.96	388,673.96	-
<i>VBGROR</i>	525,675.02	525,675.02	-
<i>WIEONSHORE</i>	13,420.09	13,420.09	-
<i>WIEPV</i>	429,295.97	429,295.97	-
<i>WIEROR</i>	920,567.19	920,567.19	-
<i>WIETKWCNG</i>	1,073,668.63	1,076,293.98	0.24

Table 5: Comparison of the decision variables from the disaggregated and complete model runs for each generation resource.