

Tutorial on how to use tpmPy toolbox

Step 0: Key imports...

```
In [1]: import numpy as np
from numpy.linalg import pinv
from tensorlib.tucker import hosvd
from tensorlib import dtensor, unfolded_dtensor
from tensorlib.tpconvex import to_cno, to_irno
from tensorlib.grid import coordinate_grid, map_points_to_grid
from tensorlib.draw import draw_weighting_system
from tensorlib.inference import reconstruct, infer_ats, rulify, infer_from_ruleset, genetic_algo_rulify
from tensorlib.metrics import calc_R2, calc_nRMSE_maxmin, calc_nRMSE_iqr, calc_cindex
from funtools import reduce

#####
## Constants, file paths
#####
CSVFILE = "./app_data_demo/data_v2.csv"
DATAFILE = "./app_data_demo/data_v2.txt"
AGGREGATED_VALUES_FILE = "./app_data_demo/data_v2.npy"
```

Step 1: read the data from file or elsewhere

```
In [2]: import pandas as pd

df = pd.read_csv(CSVFILE, delimiter=',', encoding='utf-8')

## you can keep only a subset of columns if you want...
##df = df[['input1', 'outputcol']]

cols_to_export = ['input1', 'input2', 'input3', 'outputcol']
df[cols_to_export].to_csv(DATAFILE, sep=' ', header=False, index=False)

print(f"Number of rows in final dataset: {len(df)}")
df.head()
```

Number of rows in final dataset: 15

```
Out[2]:
```

	input1	input2	input3	outputcol
0	8	5	9	69
1	12	4	7	4
2	2	9	0	13
3	3	3	-3	12
4	1	9	8	10

```
In [3]: df.describe() ## stats on columns
```

```
Out[3]:
```

	input1	input2	input3	outputcol
count	15.000000	15.000000	15.000000	15.000000
mean	8.600000	7.200000	29.133333	59.266667
std	6.400893	6.270111	64.357114	112.262745
min	0.000000	0.000000	-29.000000	-18.000000
25%	2.500000	3.000000	0.500000	5.000000
50%	9.000000	5.000000	7.000000	13.000000
75%	12.500000	9.000000	14.500000	46.000000
max	19.000000	19.000000	229.000000	342.000000

Step 2: Create discretization grid

Explanation on coordinate grid

To do TP modeling, we need an N-dimensional grid of coordinates.

In our case, we have 2 input dimensions and 1 output dimension. So our grid will be 2-dimensional.

Each of the dimensions ranges from a minimum to a maximum value, at a pre-determined (or indeed variable) step size.

One way to create the grid would be by listing all coordinates in all dimensions. If we had regular step intervals, and a large number of coordinates, we could also use the `create_from_polyranges` constructor.

But in our case, we will choose a different path. **The idea is that we want to try many different gridding alternatives and find the best reconstruction over all of them, using a genetic algorithm.**

```
In [4]: arr = np.loadtxt(DATAFILE)
data = dtensor(arr)
data[:5, :] ## first 5 rows

Out[4]: dtensor([[ 8.,  5.,  9., 69.],
 [12.,  4.,  7.,  4.],
 [ 2.,  9.,  0., 13.],
 [ 3.,  3., -3., 12.],
 [ 1.,  9.,  8., 10.]])

In [5]: dim_names=['input1', 'input2', 'input3'] ## outputcol is skipped because that's an output

In [6]: ## this is the genetic algorithm:
## see df.describe above. Minimum is 0 or 1 and maximum is 19
gen_algo = coordinate_grid.genetic_algo_search_for_coordinates_r2(data, [(0, 19), (0, 19), (-29, 229)], population_size=1000, validation_error_callback=lambda x: print(x))

Generation 1 [#####>] 100% | Elapsed: 00h 00m 01s | Remaining: 00h 00m 00s
Validation error for top-top entity so far: 0.7484859538409422
... with params [[0.20564552090191154, 0.8431310455894512, 0.9445232980841441], [0.02362695643729425, 0.21508008006594298, 0.6580191764858615], [0.27672301602620347, 0.46056205108523296, 0.7769607667700054, 0.9346798994226699], (0, 19), (0, 19), (-29, 229), 'mean', 2]

Generation 3 [#####>] 100% | Elapsed: 00h 00m 01s | Remaining: 00h 00m 00s
Validation error for top-top entity so far: 0.7978809098239812, because old < new: True
... with params [[0.24656413623274354, 0.786294749312991, 0.796294749312991, 0.84831028517425, 0.8602628782948306], [0.501873535223912, 0.5335079887474733, 0.8386415109113635, 0.9058292562966329, 0.917520673286434], [0.1043712489168985, 0.5059788780515969, 0.6618362376805598], (0, 19), (0, 19), (-29, 229), 'median', 3]

Generation 4 [#####>] 100% | Elapsed: 00h 00m 00s | Remaining: 00h 00m 00s
Validation error for top-top entity so far: 0.895859232831446, because old < new: True
... with params [[0.14645757880138638, 0.570220872164866, 0.7464319382210569, 0.7979013474487685], [0.12015846720631969, 0.3666581355402387, 0.6685358223054563, 0.9037011645889055, 0.9305774907024656], [0.1773499487807818, 0.39376646974041585], (0, 19), (0, 19), (-29, 229), 'wsum', 3]
top top entity: 10 [#####>] 100% | Elapsed: 00h 00m 01s | Remaining: 00h 00m 00s
[[0.14645757880138638, 0.570220872164866, 0.7464319382210569, 0.7979013474487685], [0.12015846720631969, 0.3666581355402387, 0.6685358223054563, 0.9037011645889055, 0.9305774907024656], [0.1773499487807818, 0.39376646974041585], (0, 19), (0, 19), (-29, 229), 'wsum', 3]

In [7]: ## note that "Validation error" is in fact the R2 score between the reconstructed and original data
## so, a higher validation error is better
## we can see the results on the best entity found above, however, we
## can also retrieve all information from the returned object, gen_algo:
print(f"fitness: {gen_algo.top_top_entity.get_fitness()}")
print(f"phenotype: {gen_algo.top_top_entity.get_phenotype()}")

fitness: 0.895859232831446
phenotype: [[0.14645757880138638, 0.570220872164866, 0.7464319382210569, 0.7979013474487685], [0.12015846720631969, 0.3666581355402387, 0.6685358223054563, 0.9037011645889055, 0.9305774907024656], [0.1773499487807818, 0.39376646974041585], (0, 19), (0, 19), (-29, 229), 'wsum', 3]
```

Based on the above, we can use the `create_from_ranges_and_positions_as_percentages` constructor to create the grid that was found to be optimal

```
In [8]: mygrid = coordinate_grid.create_from_ranges_and_positions_as_percentages(
    [(0, 19), (0, 19), (-29, 229)],
    [
        [0.14645757880138638, 0.570220872164866, 0.7464319382210569, 0.7979013474487685],
        [0.12015846720631969, 0.3666581355402387, 0.6685358223054563, 0.9037011645889055, 0.9305774907024656],
        [0.1773499487807818, 0.39376646974041585]
    ],
    dim_names=dim_names)

print("Coordinates per dimension:")
```

```
for inx, dim in enumerate(mygrid.get_coords_per_dim()):
    print(f"{dim_names[inx]}:")
    print(dim)
    print(f"{len(dim)} coordinates", end="\n\n")
```

Coordinates per dimension:

input1:

(2.782693997226341, 10.834196571132455, 14.182206826200082, 15.160125601526602)

4 coordinates

input2:

(2.283010876920074, 6.966504575264535, 12.702180623803669, 17.170322127189205, 17.680972323346847)

5 coordinates

input3:

(16.75628678544171, 72.59174919302728)

2 coordinates

Next, we load the data, convert it to a dense tensor (dtensor) and map it onto the coordinate grid. To do this:

- First, we associate each data point with the closest grid point (in Euclidean terms). These are referred to as 'primary associations'
- In case there are grid points that do not have at least P_c data points associated with them, we select as many further data points that are closest to those grid points as necessary (these data points may then be associated with more than 1 grid point - through a primary association and one or more secondary associations)
- Next, we aggregate the output value of each data point associated with each grid point, to get the aggregates tensor (by taking the mean, median or weighted sum of the output values, or by taking the output value of the closest data point). We also get back a counts tensor, which reveals the number of primary associations at each grid point

```
In [9]: arr = np.loadtxt(DATAFILE)
data = dtensor(arr)
data[:5, :]
```

```
Out[9]: dtensor([[ 8.,  5.,  9., 69.],
 [12.,  4.,  7.,  4.],
 [ 2.,  9.,  0., 13.],
 [ 3.,  3., -3., 12.],
 [ 1.,  9.,  8., 10.]])
```

```
In [10]: ## Recall that from the genetic algorithm, the output said 'wsum' and 'Pc=3'
```

```
aggs, counts = map_points_to_grid(data, mygrid, agg='wsum', Pc=3) ## can be closest, mean, median or wsum
aggs_filled = np.nan_to_num(aggs, nan=0) ## now, because the frequency of some gridpoints is 0, we need to fill it
aggregates = dtensor(aggs_filled)
print(data.shape)
print(counts.shape)
```

```
Bucketing all datapoints by gridpoint based on primary associations [#####>] 100% | Elapsed: 00h 00m 00s | Remaining: 00h 00m 00s
Bucket further datapoints by gridpoint based on secondary associations (where necessary) [#####>] 100% | Elapsed: 00h 00m 00s | Remaining: 00h 00m 00s
Aggregate (wsum) [#####>] 100% | Elapsed: 00h 00m 00s | Remaining: 00h 00m 00s
(15, 4)
(4, 5, 2)
```

... and here are some ways in which you can interact with mygrid, counts and aggregates

```
In [11]: ## Random sampling from any numpy ndarray
def sample_from_ndarr(arr, k=1):
    """
    Samples k items from arr and also returns the coordinates of the samples
    """
    flat_indices = np.random.choice(arr.size, size=k, replace=False) ## flattened indices
    multi_indices = np.array(np.unravel_index(flat_indices, arr.shape)).T ## convert to n-dimensional indices
    values = arr[np.unravel_index(flat_indices, arr.shape)] ## get the values

    return (values, multi_indices)
```

```
In [12]: samples, indices = sample_from_ndarr(mygrid.get_grid_with_coords(), 1)

indices = indices[0][:~1]
```

```
print(f"Sample coordinate: {samples[0]} - at indices: {indices}")
print(f"Count for coordinate: {counts[*indices]}")
print(f"Aggregate for coordinate: {aggregates[*indices]}")
print(f"shape of each data structure: {mygrid.get_grid_with_coords().shape}, {counts.shape}, {aggregates.shape}")
print(f"Sum of all counts: {np.ravel(counts).std()}")
```

```
Sample coordinate: 17.680972323346847 - at indices: [1 4 1]
Count for coordinate: 0
Aggregate for coordinate: 9.054261207580566
shape of each data structure: (4, 5, 2, 3), (4, 5, 2), (4, 5, 2)
Sum of all counts: 0.6590713163232034
```

Recall that mygrid is a 4-by-5 grid. Each data point in the data tensor has now been associated with 1 (or more, if needed) grid point, and all data points at each grid point have been aggregated by choosing the closest data point and using its output as the 'value' of the grid point. This is what 'aggregates' represents: it is a 2-D tensor, such that each value is the output value of the data point that is closest to the given grid point...

Since the grid can be quite big, you might want to persist the aggregated values...

```
In [13]: np.save(AGGREGATED_VALUES_FILE, aggs_filled)
```

```
In [14]: ## you can re-load them like this:
         aggregates = dtensor(np.load(AGGREGATED_VALUES_FILE))
```

Next, we will perform HOSVD. When the program asks us how many eigenvalues we want to keep, let's type 2 (as long as there are at least 2 singular values in each dimension) and then hit Enter

```
In [16]: Us, S, eigvals = hosvd(aggregates, with_eigvals=True)

print("Eigenvalues in each dimension:")
for inx, eigval in enumerate(eigvals):
    print(f"Dimension {inx+1}: {eigval}")

rank_to_keep = int(input("How many singular values do you want to keep?"))

Us_tilde, S_tilde = hosvd(aggregates, rank=[rank_to_keep]*aggregates.ndim)

reconstructed = reconstruct(S_tilde, Us_tilde)
rmse_reconstruction = np.sqrt(np.mean((reconstructed - aggregates)**2))

r2_reconstruction = calc_R2(aggregates, reconstructed)
#c_reconstruction = calc_cindex(aggregates, reconstructed)
print(f"Reconstruction R2 is: {r2_reconstruction}")
#print(f"Reconstruction Concordance index is: {c_reconstruction}")

data_reconstructed = infer_ats(S_tilde, Us_tilde, mygrid, data[:, :4])

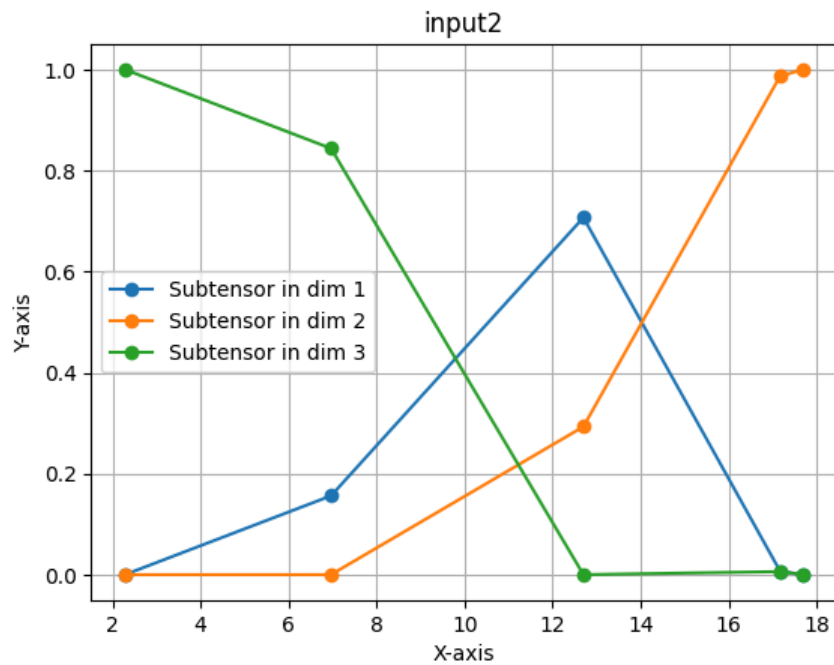
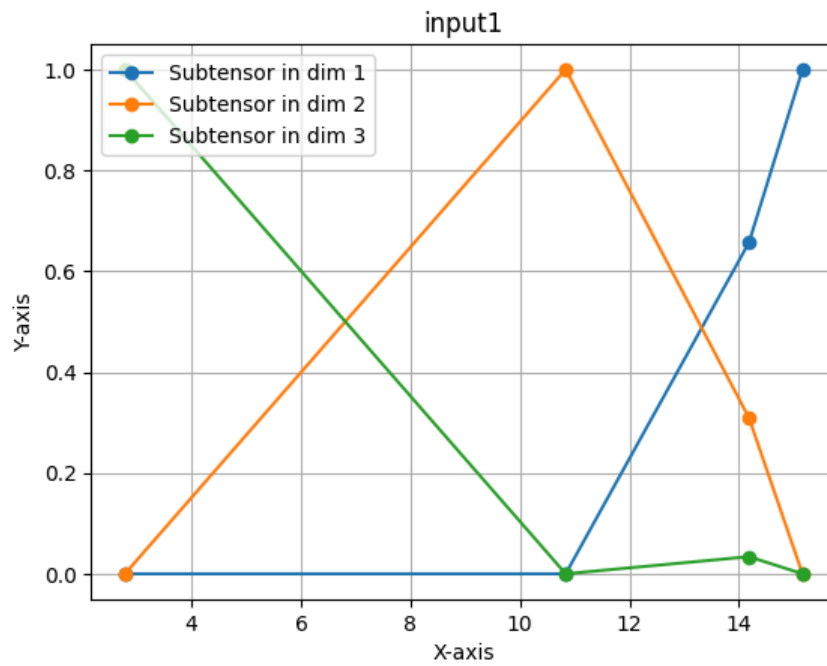
nrmse_mmin_dataset = calc_nRMSE_maxmin(data[:, -1].reshape(-1,1), data_reconstructed)
nrmse_iqr_dataset = calc_nRMSE_iqr(data[:, -1].reshape(-1,1), data_reconstructed)
r2_dataset = calc_R2(data[:, -1].reshape(-1,1), data_reconstructed)
c_dataset = calc_cindex(data[:, -1].reshape(-1,1), data_reconstructed)
print(f"Normalized RMSE on original data (dividing by max - min) is: {nrmse_mmin_dataset}")
print(f"Normalized RMSE on original data (dividing by interquartile range) is: {nrmse_iqr_dataset}")
print(f"R2 on original data is: {r2_dataset}")
print(f"Concordance index on original data is: {c_dataset}")
```

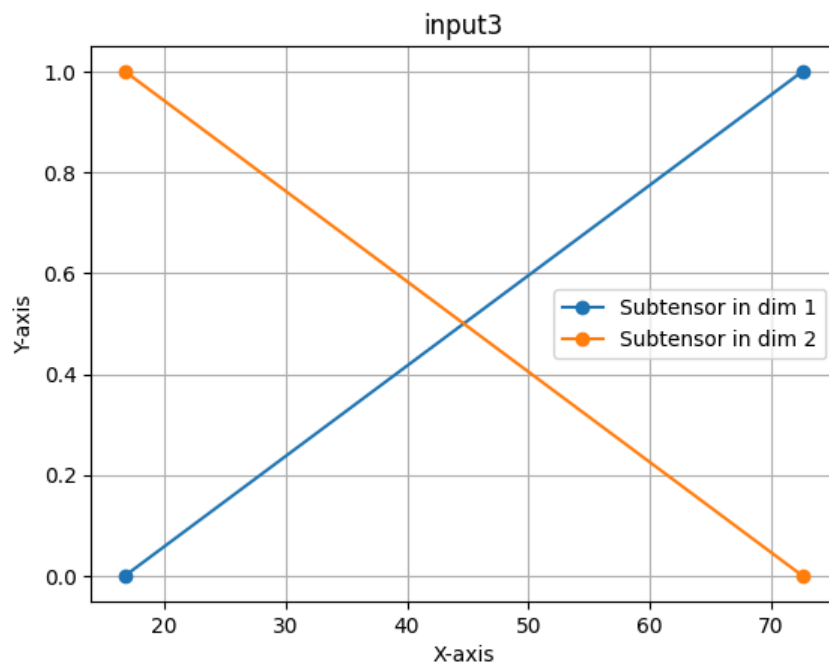
```
Eigenvalues in each dimension:
Dimension 1: [4.5967384e+05 3.8402939e+03 5.1213196e+01 7.3535290e+00]
Dimension 2: [4.5921919e+05 3.9294260e+03 4.1841306e+02 5.0179181e+00 6.1136317e-01]
Dimension 3: [462963.56      609.2066]
Reconstruction R2 is: 0.9986146688461304
Normalized RMSE on original data (dividing by max - min) is: 0.09148001113868504
Normalized RMSE on original data (dividing by interquartile range) is: 0.803239122193332
R2 on original data is: 0.9077960913491149
Concordance index on original data is: 0.7333333333333333
```

```
In [17]: S_cno, Us_cno = to_cno(S_tilde, Us_tilde)
```

```
Distance from NO: 5.688616369072868e-05
Distance from NO: 0.34676610294390253
```

```
In [18]: draw_weighting_system(Us_cno, mygrid)
```



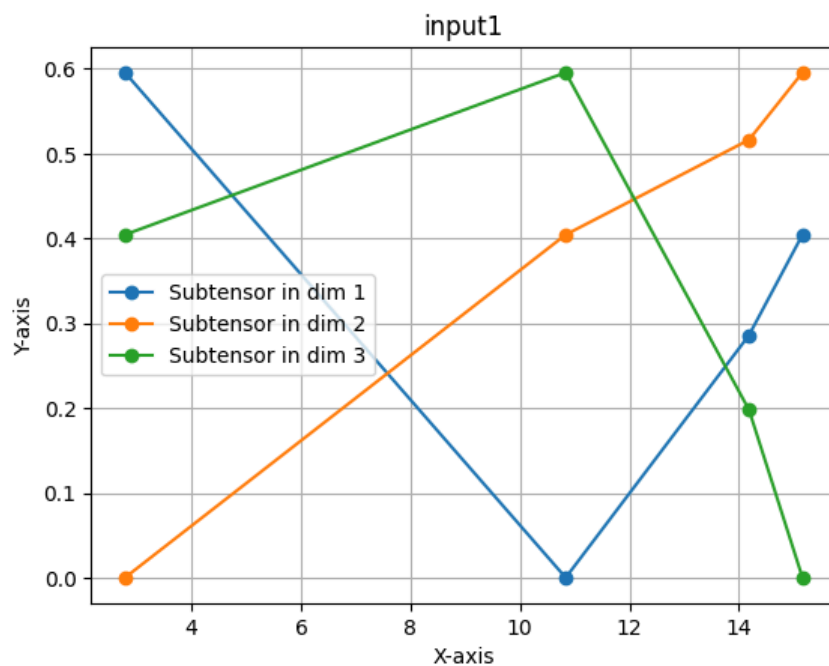


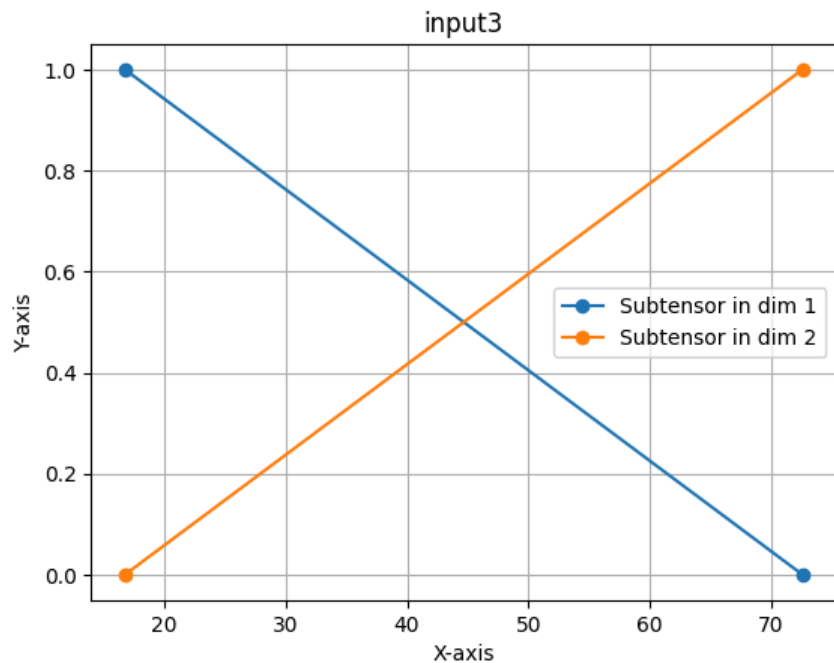
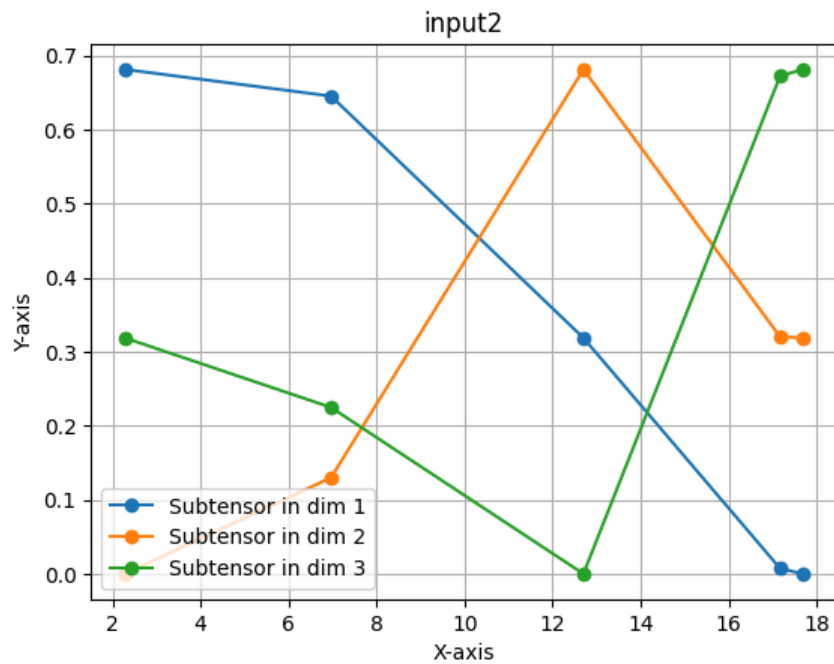
In [19]: `S_irno, Us_irno = to_irno(S_tilde, Us_tilde)`

C:\Dolgaim\development\tpmpy\tensorlib\tpconvex.py:40: RuntimeWarning: invalid value encountered in divide

$$U3 = U2 / (np.ones((n, 1)) * np.max(U2, axis=0))$$

In [20]: `draw_weighting_system(Us_irno, mygrid)`





```
In [21]: print(f"inference at random point using CNO system:\n {infer_at(S_cno, Us_cno, mygrid, np.array([[2, 5, 1], [11, 11, 1]])}")
print(f"inference at random point using IRNO system:\n {infer_at(S_irno, Us_irno, mygrid, np.array([[2, 5, 1], [11, 11, 1]])}")

inference at random point using CNO system:
[[50.19533763]
 [ 1.96805144]]
inference at random point using IRNO system:
[[50.19533763]
 [ 1.96805144]]

In [22]: print(f"Whereas the original system yields these results:\n {infer_at(S, Us, mygrid, np.array([[2, 5, 1], [11, 11, 1]])}")

Whereas the original system yields these results:
[[47.56450182]
 [ 7.22699615]]
```

Based on the above, we can already see what ranges of the different input dimensions are unique and separate from each other. Based on this, we can now generate a set of weighted rules

Step 3: Rule generation

Since there are quite a few hyperparameters, like support length to use, alpha-cut to use and maximum antecedent variables per input dimension, we can use a genetic algorithm again to find the rule set that can be best used to reconstruct the original data

```
In [23]: ga = genetic_algo_rulify(S_cno, Us_cno, mygrid, counts, data, population_sz=25, num_generations=10, dict_of_params={
    'min_wfun_val_ab': (0.5, 0.95), ## minimum alpha cut sampled from this range
    'min_rule_weight_ab': (0.001, 0.005), ## minimum rule representation sampled from this range
    'num_samples_ab': (1,4), ## number of samples used to generate rules sampled from this range
    'max_antecedent_nums_ab_per_dim': [(2,5), (2,5), (2,5)], ## maximum number of antecedent vars per dimension s
    'num_trials_to_aggregate': 5,
    'pct_data_points_ab': (0.1, 0.5) ## percentage of data points considered
})
```

Generation 1 [#####] 100% | Elapsed: 00h 00m 03s | Remaining: 00h 00m 00s

Validation error for top-top entity so far: 0.9128951817979347
 ... with params min wf val: 0.6135374006996508; min rule w: 0.0015296222391833407; max num of antecedents: [2, 5, 4]; samples for creating rules: 2; pct of datapts: 0.181175503693347; datapts considered: [14 5]; antecedent ranges: [(np.float64(-0.1), np.float64(0.1)), (np.float64(9.9), np.float64(10.1))], [(np.float64(5.9), np.float64(6.1)), (np.float64(17.9), np.float64(18.1))], [(np.float64(1.9), np.float64(2.1)), (np.float64(80.9), np.float64(81.1))]]

Generation 2 [#####] 100% | Elapsed: 00h 00m 01s | Remaining: 00h 00m 00s

Validation error for top-top entity so far: 0.9210924276426304, because old < new: True
 ... with params min wf val: 0.6135374006996508; min rule w: 0.0015296222391833407; max num of antecedents: [2, 5, 5]; samples for creating rules: 2; pct of datapts: 0.27093416960359906; datapts considered: [1 11 4 9]; antecedent ranges: [(np.float64(11.9), np.float64(12.1)), (np.float64(18.9), np.float64(19.1))], [(np.float64(3.9), np.float64(4.1)), (np.float64(17.9), np.float64(19.1))], [(np.float64(6.9), np.float64(11.1)), (np.float64(100.9), np.float64(101.1))]]
 top top entity: 10 [#####] 100% | Elapsed: 00h 00m 01s | Remaining: 00h 00m 00s
 min wf val: 0.6135374006996508; min rule w: 0.0015296222391833407; max num of antecedents: [2, 5, 5]; samples for creating rules: 2; pct of datapts: 0.27093416960359906; datapts considered: [1 11 4 9]; antecedent ranges: [(np.float64(11.9), np.float64(12.1)), (np.float64(18.9), np.float64(19.1))], [(np.float64(3.9), np.float64(4.1)), (np.float64(17.9), np.float64(19.1))], [(np.float64(6.9), np.float64(11.1)), (np.float64(100.9), np.float64(101.1))]]

```
In [24]: min_weighting_func_val = ga.top_top_entity.get_phenotype()[0]
min_rule_weight = ga.top_top_entity.get_phenotype()[1]
samples_consequents = ga.top_top_entity.get_phenotype()[2]
max_antecedent_nums = ga.top_top_entity.get_phenotype()[3]
indices = ga.top_top_entity.get_phenotype()[4]
antecedent_ranges = ga.top_top_entity.get_phenotype()[5]
rules = ga.top_top_entity.get_phenotype()[6]

print(min_weighting_func_val)
print(min_rule_weight)
print(samples_consequents)
print(max_antecedent_nums)
print(indices)
print()
print("antecedent ranges:")
for dim in antecedent_ranges:
    print(dim)
print()
print(rules)
```



```
0.6135374006996508
0.0015296222391833407
```

```
2
[2, 5, 5]
[ 1 11  4  9]
```

antecedent ranges:

```
[(np.float64(11.9), np.float64(12.1)), (np.float64(18.9), np.float64(19.1))]
[(np.float64(3.9), np.float64(4.1)), (np.float64(17.9), np.float64(19.1))]
[(np.float64(6.9), np.float64(11.1)), (np.float64(100.9), np.float64(101.1))]
```

```
0 AND 0 AND 0 => 12.34 (weight: 0.1333)
1 AND 0 AND 0 => -22.53 (weight: 0.1333)
0 AND 0 AND 1 => 1.43 (weight: 0.0667)
0 AND 1 AND 0 => 275.26 (weight: 0.0667)
1 AND 1 AND 0 => 283.45 (weight: 0.0667)
0 AND 1 AND 1 => -116.47 (weight: 0.0000)
1 AND 0 AND 1 => 15.92 (weight: 0.0000)
1 AND 1 AND 1 => -127.69 (weight: 0.0000)
```

Total number of rules: 8

Sum of all weights: 0.4666666666666667

```
In [25]: data_reconstructed = infer_from_ruleset(rules, data[:, :-1], cutoff_weight=min_rule_weight)
r2_dataset = calc_R2(data[:, -1].reshape(-1,1), data_reconstructed)
print(f"R2 score: {r2_dataset}")
```

R2 score: 0.9210924276426304

```
In [26]: CSV_WITH_TP_RULES = "./app_data_demo/rules001_v2.csv"
```

```
In [27]: rules.to_csv(min_weight=0.0000001, filename=CSV_WITH_TP_RULES)
```

```
In [ ]:
```