

Report

Andrija Radočaj

May 5, 2025

Many thanks to my mentor, dr. Fesel-Kamenik, for the opportunity, his guidance and assistance throughout the traineeship.

Thank you to Lara, Iva, Nikiša, Peter, Bučka, Adrian and Fran for all the lunches and beautiful moments.

Special thanks to my partner in crime from the office, Lara.

Pika konc.

Contents

1	Introduction	4
2	Starting point	5
2.1	Careful with handling $\Delta\phi$	6
3	Preparing the environment	6
4	Submitting a job	8
5	Code parallelization	9
5.1	Function of one argument, one return value	10
5.2	Function of one argument, two return values	12
5.3	Function of two arguments, one return value	14
5.4	Function of two arguments, two return values	15
5.5	Implementation of shared memory usage	17
6	Reading from the files	23
6.1	JetClass data set	23
6.2	Rodem data set	27
6.3	Aspen Open Jets data set	30
7	JetClass data set: visualization of jets	32
8	Probability distribution of the radius	34
8.1	Contributions to the probability - 1D problem	34
8.2	Contributions to the probability density - 2D problem	34
8.3	Histogram and Poisson distribution relation - in general	35
8.4	Covariance matrix - in general	36
8.5	Covariance matrix of the probability (density) histogram	37
9	Explored modeling methods/tools	40
9.1	Bayesian machine scientist	40
9.2	Quantile transformation	43
9.3	Box-Cox transformation	60
9.4	Kernel density estimation	70
10	Results	72
10.1	Moments	72
10.2	Radius probability distribution	74
10.3	Radius probability density function	76

10.4	Bayesian machine scientist	79
------	--------------------------------------	----

1 Introduction

This report explores various data structures and analyzes data from multiple particle physics data sets. The second section introduces the relevant quantities. The third section is a guide through the preparation of the environment in which the codes can be run. The fourth section deals with running the codes on a cluster by submitting them as jobs. The fifth section is a tutorial on code parallelization and implementation of shared-memory usage between parallel processes. The sixth section shows how to read three different data sets. The seventh section helps with the visualization of jets. The eighth section goes in detail on how the normalization of the histogram introduces bin correlations. The ninth section is a tutorial on all the modeling methods and tools explored. The final tenth section presents the results. The end goal of this report was to find an analytic function that could describe the probability distribution or probability density function of the radius of QCD jets.

In this report, three files have been used.

- The file `WToQQ_070.root` was extracted from `JetClass.Pythia_train_100M_part7.tar` which was downloaded from <https://zenodo.org/records/6619768>.
This file was used to study the data structure in `.root` files and to compare the data with the other two files.
- The file `RunG.batch0.h5` was downloaded from <https://www.fdr.uni-hamburg.de/record/16505>.
This file contains a mixture of jets with diverse origins, so the assumption was made that jets with transverse momentum up to approximately 500 GeV are almost exclusively of QCD origin.
- The file `QCDjj_pT_450_1200_train01.h5` was extracted from `QCD.tar.gz` which was downloaded from <https://zenodo.org/records/12793616>.

Python codes were run on a remote computer cluster as jobs.

The Python distribution used is [Anaconda3/2023.03-1](#).

SSH (Secure Socket Shell) was run in WSL (Windows Subsystem for Linux) for remote connection.

The main codes `jet_dist.py`, `functions.py` and `data_root.py` are the ones responsible for the end results. They are not presented in this report directly because of their size.

2 Starting point

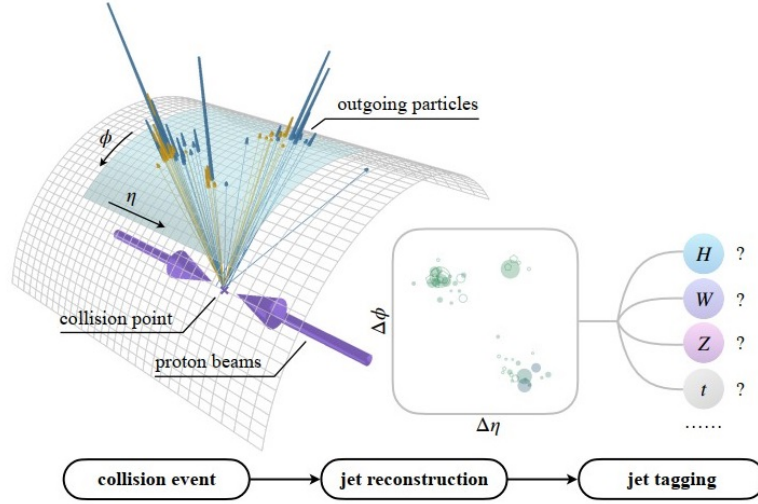


Figure 1: Illustration of jet tagging at the CERN's LHC [1]

The linear momentum of some jet's constituents is usually represented by $\vec{p} = (p_T, \eta, \phi)$ [2]. These quantities are best described by Figure 2 below and the equations (1)-(5).

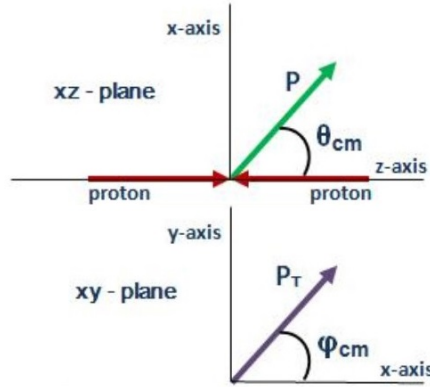


Figure 2: Coordinate system [2]

$$\eta = \frac{1}{2} \ln \left(\frac{p + p_z}{p - p_z} \right) = -\ln \left[\tan \left(\frac{\theta}{2} \right) \right] \quad (1)$$

$$p_x = p_T \cos \phi \quad (2)$$

$$p_y = p_T \sin \phi \quad (3)$$

$$p_z = p_T \sinh \eta \quad (4)$$

$$p = p_T \cosh \eta \quad (5)$$

Now, the quantity R is easily calculated with the help of equation (6) below:

$$R_{i(J)} = \sqrt{(\eta_J - \eta_{i(J)})^2 + (\phi_J - \phi_{i(J)})^2} \equiv \sqrt{(\Delta\eta_{i(J)})^2 + (\Delta\phi_{i(J)})^2}, \quad (6)$$

where $\eta_{i(J)}$ and $\phi_{i(J)}$ are from the i -th particle of the J -th jet, while η_J and ϕ_J are from the corresponding jet itself. The latter are calculated by first summing 3-momenta (p_x, p_y, p_z) of all particles in a jet and then calculating p_J and p_{J_z} .

2.1 Careful with handling $\Delta\phi$

Differences $\Delta\phi_{i(J)}$ should be treated with care because ϕ is an azimuthal angle. That is, it should be $\Delta\phi_{i(J)} \in [-\pi, \pi]$. While .root files already have $\Delta\phi_{i(J)}$ correctly calculated, in .h5 files $\Delta\phi_{i(J)}$ has to be carefully calculated with the help of ϕ_J and $\phi_{i(J)}$ as in the expressions (7) and (8) below:

$$\Delta\phi_{i(J)} \leftarrow \phi_J - \phi_{i(J)}; \quad (7)$$

$$\Delta\phi_{i(J)} \leftarrow \text{MOD}(\Delta\phi_{i(J)} + \pi, 2\pi) - \pi. \quad (8)$$

3 Preparing the environment

The first two steps are to be skipped if the environment wants to be set up on a local computer. The commands from the third step then have to be run in Anaconda Powershell Prompt.

The first step is to open WSL and remotely connect to the cluster via SSH.

The second step is to run the following commands.

WSL terminal:

```
module load Anaconda3/2023.03-1 #to load the module
conda init bash #just the first time
source ~/.bashrc #to be able to activate the environment
```

Figure 3: Commands for WSL once connected via SSH to the cluster

The third step is to run the following commands.

WSL terminal / Anaconda Powershell Prompt:

```
1 conda update conda
2 conda config --add channels anaconda
3 conda config --add channels conda-forge
4 conda info --envs #list all the environments
5 conda create -n jet scipy numpy matplotlib #create a new environment
```

```

6  conda activate jet #activate the new environment
7  python --version #check the python version
8  conda install uproot #for reading .root files
9  conda install h5py #for reading .h5 files
10 conda install vector #for managing 4-vectors
11 conda install awkward #works closely with vector package
12 conda install pandas #for use of machine scientist
13 conda install ipywidgets #for use of machine scientist
14 conda install sympy #for use of machine scientist
15 conda install scikit-learn #for quantile transformation and KDE
16 conda list package_name #check version of the installed package

```

Figure 4: Commands for setting up the environment

Additionally, if one is setting up the environment on a local computer, then, aside from running the mandatory commands from the third step in Anaconda Powershell Prompt, one can optionally run the following commands there too to make the environment appear as kernel in Jupyter Notebook.

Anaconda Powershell Prompt:

```

1  conda install ipykernel
2  #make the environment appear as kernel in Jupyter Notebook:
3  python -m ipykernel install --user --name jet --display-name "jet_project"
4  conda deactivate jet
5  jupyter notebook

```

Figure 5: Optional commands for Anaconda Powershell Prompt

After the environment has been set up, each time a new connection through SSH in WSL is made, the following commands have to be run.

WSL terminal:

```

module load Anaconda3/2023.03-1 #to load the module
source ~/.bashrc #to be able to activate the environment
conda activate jet #activate the environment

```

Figure 6: Commands for each time a new connection is made through SSH in WSL

4 Submitting a job

The first step in running some code on a cluster is to create a Python script and open it.

WSL terminal:

```
nano code.py
```

Once a Python script is created and opened, one can write some simple code in it to test things out.

code.py:

```
1 print("Hello_world!")
```

To exit the script, one should press Ctrl+X and then follow the instructions to save the changes or not. The second step is to create a shell script and open it.

WSL terminal:

```
nano job.sh
```

Once a shell script is created and opened, one should write the following instructions for SLURM to then read them. SLURM (Simple Linux Utility for Resource Management) is an open-source job scheduler used to manage and allocate resources on high-performance computing (HPC) clusters - it takes job requests, figures out when and where to run them, and makes sure resources are used efficiently.

job.sh:

```
# Use BASH (Bourne Again Shell) as a command-line interpreter:
#!/bin/bash

#SBATCH --job-name=steve_job

# In this file the output of a code will appear
#SBATCH --output=output.out #SLURM directive

# In this file potential errors and logging will appear
#SBATCH --error=error.err #SLURM directive

# To prevent accidental infinite loops and minimize disruption:
# - A job will be automatically killed when the time limit expires,
#   ensuring other users don't have to wait unnecessarily.
```

```
#SBATCH --time=00:00:10 # SLURM directive

# To position a job in the queue appropriately:
# - If the memory is too little for the code to run, the job will be
#   killed and may raise an error.
# - If the memory is too much, the job will wait longer in queue
#   than necessary, but once running, only the required memory
#   will be used.
#SBATCH --mem=100M # SLURM directive

#short(<3h)
#day(3h-1day)
#long(1day-30days)
#SBATCH --partition=short #SLURM directive

# Number of CPUs to use:
#SBATCH --cpus-per-task=1 #SLURM directive

python code.py
```

Again, to exit the script, one should press Ctrl+X and then follow the instructions to save the changes or not. The last step is to submit the job.

WSL terminal:

```
sbatch job.sh
```

To see the status of a job the following command can be run.

WSL terminal:

```
squeue -u username
```

One can open the output and error file with the help of nano editor, same as with the scripts.

5 Code parallelization

To speed up code execution, parallelization of loops, where each iteration is independent of the others, can be done. In this section a couple of simple codes with parallelized loops will be presented. For running the codes on a computer cluster, a shell script has to be created.

WSL terminal:

```
nano job.sh
```

job.sh:

```
#!/bin/bash

#SBATCH --job-name=steve_job

#SBATCH --output=output.out

#SBATCH --error=error.err

#SBATCH --time=00:00:30

#SBATCH --mem=100M

#SBATCH --partition=short

# For CPUs per task put the same number as the number
# of processes in a Python script:
#SBATCH --cpus-per-task=4

python code.py
```

5.1 Function of one argument, one return value

The first code works with a function that takes only one argument and returns only one value.

WSL terminal:

```
nano code.py
```

code.py:

```
1 import logging
2 from multiprocessing import Pool
3
4 def square(x):
5     logging.basicConfig(level=logging.INFO)
6     logging.info(f"argument_{x}_in_progress..")
7     return x * x
```

```

8
9 def main():
10     # List of arguments to pass to the function
11     input_data = [1, 2, 3, 4, 5]
12
13     # Create a pool of worker processes
14     with Pool(processes=4) as pool:
15         # Use map to pass each element of data
16         # as an argument to the function
17         squared = pool.map(square, input_data)
18
19     # Print results
20     print("Input:", input_data)
21     print("Squared:", squared)
22
23 if __name__ == "__main__":
24     main()

```

`if __name__ == "__main__":` acts as a guard to prevent certain parts of the code from being executed when child processes are created during multiprocessing.

The next step is to submit the job.

WSL terminal:

```
sbatch job.sh
```

After submitting the job, it should be checked if it is pending, running, or done.

WSL terminal:

```
squeue -u username
```

If the job does not appear in the list with the last command, that means that it is done and one should then first check the error file.

WSL terminal:

```
nano error.err
```

error.err:

```

INFO:root:argument 2 in progress..
INFO:root:argument 1 in progress..
INFO:root:argument 5 in progress..

```

```
INFO:root:argument 4 in progress..  
INFO:root:argument 3 in progress..
```

If the error file only contains log information, everything is fine, and one should proceed with opening the output file.

WSL terminal:

```
nano output.out
```

output.out:

```
Input: [1, 2, 3, 4, 5]  
Squared: [1, 4, 9, 16, 25]
```

5.2 Function of one argument, two return values

The second code works with a function that takes only one argument and returns two values.

WSL terminal:

```
nano code.py
```

code.py:

```
1 import logging  
2 from multiprocessing import Pool  
3  
4 def compute(x):  
5     logging.basicConfig(level=logging.INFO)  
6     logging.info(f"argument_{x}_in_progress..")  
7     return x * x, x + x  
8  
9 def main():  
10     # List of arguments to pass to the function  
11     input_data = [1, 2, 3, 4, 5]  
12  
13     # Create a pool of worker processes  
14     with Pool(processes=4) as pool:  
15         # Use map to pass each element of data as  
16         # an argument to the function  
17         results = pool.map(compute, input_data)  
18  
19     # Create empty lists to store the results separately
```

```

20     squares = []
21     sums = []
22
23     # Unpack the results manually using a for loop
24     for square,summation in results:
25         squares.append(square)
26         sums.append(summation)
27
28     # Print results
29     print("Input:", input_data)
30     print("Results:", results)
31     print("Squares:", squares)
32     print("Sums:", sums)
33
34 if __name__ == "__main__":
35     main()

```

WSL terminal:

```

sbatch job.sh

```

WSL terminal:

```

nano error.err

```

error.err:

```

INFO:root:argument 2 in progress..
INFO:root:argument 1 in progress..
INFO:root:argument 3 in progress..
INFO:root:argument 4 in progress..
INFO:root:argument 5 in progress..

```

WSL terminal:

```

nano output.out

```

output.out:

```

Input: [1, 2, 3, 4, 5]
Results: [(1, 2), (4, 4), (9, 6), (16, 8), (25, 10)]
Squares: [1, 4, 9, 16, 25]
Sums: [2, 4, 6, 8, 10]

```

5.3 Function of two arguments, one return value

The third code works with a function that takes two arguments and returns only one value.

WSL terminal:

```
nano code.py
```

code.py:

```
1 import logging
2 from multiprocessing import Pool
3
4 def add(x, y):
5     logging.basicConfig(level=logging.INFO)
6     logging.info(f"argument_{(x},{y})_in_progress..")
7     return x + y
8
9 def main():
10     # List of tuples, where each tuple contains two arguments of the function
11     input_data = [(1, 2), (3, 4), (5, 6), (7, 8), (9, 10)]
12
13     # Create a pool of worker processes
14     with Pool(processes=4) as pool:
15         # Use starmap to pass each tuple from data as
16         # arguments to the function
17         sums = pool.starmap(add, input_data)
18
19     # Print results
20     print("Input:", input_data)
21     print("Sums:", sums)
22
23 if __name__ == "__main__":
24     main()
```

WSL terminal:

```
sbatch job.sh
```

WSL terminal:

```
nano error.err
```

error.err:

```
INFO:root:argument (1,2) in progress..  
INFO:root:argument (3,4) in progress..  
INFO:root:argument (9,10) in progress..  
INFO:root:argument (5,6) in progress..  
INFO:root:argument (7,8) in progress..
```

WSL terminal:

```
nano output.out
```

output.out:

```
Input: [(1, 2), (3, 4), (5, 6), (7, 8), (9, 10)]  
Sums: [3, 7, 11, 15, 19]
```

5.4 Function of two arguments, two return values

The fourth code works with a function that takes two arguments and returns two values.

WSL terminal:

```
nano code.py
```

code.py:

```
1 import logging  
2 from multiprocessing import Pool  
3  
4 def compute(x, y):  
5     logging.basicConfig(level=logging.INFO)  
6     logging.info(f"argument_{(x,y)}_in_progress..  
7     return x * y, x + y  
8  
9 def main():  
10     # List of tuples, where each tuple contains two arguments of the function  
11     input_data = [(1, 2), (3, 4), (5, 6), (7, 8), (9, 10)]  
12  
13     # Create a pool of worker processes  
14     with Pool(processes=4) as pool:  
15         # Use starmap to pass each tuple from data  
16         # as arguments to the function  
17         results = pool.starmap(compute, input_data)  
18
```



```

19     # Create empty lists to store the results separately
20     products = []
21     sums = []
22
23     # Unpack the results manually using a for loop
24     for product,summation in results:
25         products.append(product)
26         sums.append(summation)
27
28     # Print results
29     print("Input:", input_data)
30     print("Results:", results)
31     print("Products:", products)
32     print("Sums:", sums)
33
34 if __name__ == "__main__":
35     main()

```

WSL terminal:

```

sbatch job.sh

```

WSL terminal:

```

nano error.err

```

error.err:

```

INFO:root:argument (3,4) in progress..
INFO:root:argument (1,2) in progress..
INFO:root:argument (5,6) in progress..
INFO:root:argument (7,8) in progress..
INFO:root:argument (9,10) in progress..

```

WSL terminal:

```

nano output.out

```

output.out:

```

Input: [(1, 2), (3, 4), (5, 6), (7, 8), (9, 10)]
Results: [(2, 3), (12, 7), (30, 11), (56, 15), (90, 19)]
Products: [2, 12, 30, 56, 90]
Sums: [3, 7, 11, 15, 19]

```

5.5 Implementation of shared memory usage

The fifth code shows inefficient usage of memory in code parallelization and why shared memory between processes is then useful. The goal of a code is to fill the blocks on the diagonal of some matrices. By first creating and filling small matrices and then copying them onto the diagonal blocks of a bigger matrix, extra memory is unnecessarily wasted.

WSL terminal:

```
nano code.py
```

code.py:

```
1 import logging
2 import numpy as np
3 from multiprocessing import Pool
4
5 # Function to compute the matrices
6 def compute_matrix(x, y, n):
7
8     logging.basicConfig(level=logging.INFO)
9     logging.info(f"argument_{(x,y)}_in_progress..")
10
11     # Create a (n x n) matrix filled with the same value x * y
12     product_matrix = np.full((n, n), x * y)
13
14     # Create a (n x n) matrix filled with the same value x + y
15     sum_matrix = np.full((n, n), x + y)
16
17     return product_matrix, sum_matrix
18
19 def main():
20     # List of tuples,
21     # where each tuple contains the first two arguments of the function
22     input_data = [(1, 2), (3, 4), (5, 6), (7, 8), (9, 10)]
23     input_data_len=len(input_data)
24
25     # Size of each matrix (2x2 matrix)
26     n = 2
27
28     # Prepare the list of full arguments (x, y, n) for each function call
29     args=[(*input_data[i],n) for i in range(input_data_len)]
30
31     # Create a pool of worker processes
32     with Pool(processes=4) as pool:
33         # Use starmap to distribute work across the pool of processes
34         results=pool.starmap(compute_matrix,args)
35
```

```

36     # Create a big matrix where each 'product_matrix'
37     # will be placed as a separate block on the diagonal
38     # (note: small matrices are already computed and stored in 'results')
39     # - allocating extra memory here for combining into one large matrix
40     final_product_matrix=np.zeros([input_data_len*n,input_data_len*n])
41
42     # Create a big matrix where each 'sum_matrix'
43     # will be placed as a separate block on the diagonal
44     # (note: small matrices are already computed and stored in 'results')
45     # - allocating extra memory here for combining into one large matrix
46     final_sum_matrix=np.zeros([input_data_len*n,input_data_len*n])
47
48     # Placing product_matrices and sum_matrices
49     # as blocks on the diagonal of a corresponding matrix
50     idx=0
51     for product_matrix,sum_matrix in results:
52         index1=idx*n
53         index2=(idx+1)*n
54
55         # Copy product_matrix into the diagonal block of final_product_matrix
56         final_product_matrix[index1:index2,index1:index2]=product_matrix
57
58         # Copy sum_matrix into the diagonal block of final_sum_matrix
59         final_sum_matrix[index1:index2,index1:index2]=sum_matrix
60         idx+=1
61
62     # Print the final matrices
63     print("Input:", input_data)
64     print("Product_block-diagonal_matrix:")
65     print(final_product_matrix)
66     print("Sum_block-diagonal_matrix:")
67     print(final_sum_matrix)
68
69 if __name__ == "__main__":
70     main()

```

WSL terminal:

```
sbatch job.sh
```

WSL terminal:

```
nano error.err
```

error.err:

```
INFO:root:argument (1,2) in progress..  
INFO:root:argument (3,4) in progress..  
INFO:root:argument (9,10) in progress..  
INFO:root:argument (5,6) in progress..  
INFO:root:argument (7,8) in progress..
```

WSL terminal:

```
nano output.out
```

output.out:

```
Input: [(1, 2), (3, 4), (5, 6), (7, 8), (9, 10)]  
Product block-diagonal matrix:  
[[ 2.  2.  0.  0.  0.  0.  0.  0.  0.  0.]  
 [ 2.  2.  0.  0.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0. 12. 12.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0. 12. 12.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0. 30. 30.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0. 30. 30.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0. 56. 56.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0. 56. 56.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0.  0.  0. 90. 90.]  
 [ 0.  0.  0.  0.  0.  0.  0.  0. 90. 90.]]  
Sum block-diagonal matrix:  
[[ 3.  3.  0.  0.  0.  0.  0.  0.  0.  0.]  
 [ 3.  3.  0.  0.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0.  7.  7.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0.  7.  7.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0. 11. 11.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0. 11. 11.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0. 15. 15.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0. 15. 15.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0.  0.  0. 19. 19.]  
 [ 0.  0.  0.  0.  0.  0.  0.  0. 19. 19.]]
```

The sixth code shows efficient usage of memory through memory sharing between processes. This time, only one big matrix is created, and its diagonal blocks are filled in directly without an intermediate step of first creating small matrices.

WSL terminal:

```
nano code.py
```

code.py:

```
1 import logging
2 import numpy as np
3 from multiprocessing import Pool
4 from multiprocessing.shared_memory import SharedMemory
5
6 # Function to compute the matrices
7 def compute_matrix(x, y, idx, n, input_data_len, shm_product_matrix_name,
8                   shm_sum_matrix_name):
9
10     logging.basicConfig(level=logging.INFO)
11     logging.info(f"argument_{x},{y} in progress..")
12
13     # Attach to the existing shared memory block
14     existing_shm_product_matrix = SharedMemory(name=shm_product_matrix_name)
15
16     # Create a NumPy array that wraps the shared memory buffer
17     product_matrix = np.ndarray((n*input_data_len,n*input_data_len),
18                                dtype=np.float64,buffer=existing_shm_product_matrix.buf)
19
20     # Attach to the existing shared memory block
21     existing_shm_sum_matrix = SharedMemory(name=shm_sum_matrix_name)
22
23     # Create a NumPy array that wraps the shared memory buffer
24     sum_matrix = np.ndarray((n*input_data_len,n*input_data_len),
25                              dtype=np.float64,buffer=existing_shm_sum_matrix.buf)
26
27     index1=idx*n
28     index2=(idx+1)*n
29
30     # Fill the block on the diagonal with the same value x * y
31     product_matrix[index1:index2,index1:index2] = np.full((n, n), x * y)
32
33     # Fill the block on the diagonal with the same value x + y
34     sum_matrix[index1:index2,index1:index2] = np.full((n, n), x + y)
35
36 def main():
37     # List of tuples,
38     # where each tuple contains the first two arguments of the function
39     input_data = [(1, 2), (3, 4), (5, 6), (7, 8), (9, 10)]
40     input_data_len=len(input_data)
41
42     # Size of each matrix (2x2 matrix)
43     n = 2
44
45     # Create shared memory for the big matrix
46     # Each float64 takes 8 bytes
47     shm_product_matrix = SharedMemory(
48         create=True,size=n*input_data_len*n*input_data_len*8)
```

```

47
48     # View into the shared memory
49     product_matrix = np.ndarray((n*input_data_len,n*input_data_len),
50                                 dtype=np.float64,buffer=shm_product_matrix.buf)
51
52     # Create shared memory for the big matrix
53     # Each float64 takes 8 bytes
54     shm_sum_matrix = SharedMemory(
55         create=True,size=n*input_data_len*n*input_data_len*8)
56
57     # View into the shared memory
58     sum_matrix = np.ndarray((n*input_data_len, n*input_data_len),
59                             dtype=np.float64,buffer=shm_sum_matrix.buf)
60
61     # Prepare the list of full arguments
62     # (x,y,idx,n,input_data_len,shm_product_matrix.name,shm_sum_matrix.name)
63     # for each function call:
64     args=[(*input_data[idx],idx,n,input_data_len,
65           shm_product_matrix.name,shm_sum_matrix.name)
66           for idx in range(input_data_len)]
67
68     # Create a Pool of processes
69     with Pool(processes=4) as pool:
70         # Use starmap to distribute work across the pool of processes
71         pool.starmap(compute_matrix,args)
72
73     # Print the final matrices
74     print("Input:", input_data)
75     print("Product_block-diagonal_matrix:")
76     print(product_matrix)
77     print("Sum_block-diagonal_matrix:")
78     print(sum_matrix)
79
80     # Clean up the shared memory once everything is done
81     shm_product_matrix.close() # Detach from the shared memory
82     shm_product_matrix.unlink() # Release the shared memory block
83
84     # Clean up the shared memory once everything is done
85     shm_sum_matrix.close() # Detach from the shared memory
86     shm_sum_matrix.unlink() # Release the shared memory block
87
88 if __name__ == "__main__":
89     main()

```

WSL terminal:

```
sbatch job.sh
```

WSL terminal:

```
nano error.err
```

error.err:

```
INFO:root:argument (1,2) in progress..  
INFO:root:argument (3,4) in progress..  
INFO:root:argument (9,10) in progress..  
INFO:root:argument (5,6) in progress..  
INFO:root:argument (7,8) in progress..
```

WSL terminal:

```
nano output.out
```

output.out:

```
Input: [(1, 2), (3, 4), (5, 6), (7, 8), (9, 10)]  
Product block-diagonal matrix:  
[[ 2.  2.  0.  0.  0.  0.  0.  0.  0.  0.]  
 [ 2.  2.  0.  0.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0. 12. 12.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0. 12. 12.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0. 30. 30.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0. 30. 30.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0. 56. 56.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0. 56. 56.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0.  0.  0. 90. 90.]  
 [ 0.  0.  0.  0.  0.  0.  0.  0. 90. 90.]]  
Sum block-diagonal matrix:  
[[ 3.  3.  0.  0.  0.  0.  0.  0.  0.  0.]  
 [ 3.  3.  0.  0.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0.  7.  7.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0.  7.  7.  0.  0.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0. 11. 11.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0. 11. 11.  0.  0.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0. 15. 15.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0. 15. 15.  0.  0.]  
 [ 0.  0.  0.  0.  0.  0.  0.  0. 19. 19.]  
 [ 0.  0.  0.  0.  0.  0.  0.  0. 19. 19.]]
```

6 Reading from the files

6.1 JetClass data set

Each data set mentioned consists of clustered jets. Each jet has its own data components, as well as each particle within each jet. Figure 7 and Figure 8 help in understanding the data structure of .root files. Each jet has its own dictionary for which keys represent different entries or, in other words, hold different information about the jet.

code.ipynb:

```
1 import uproot
2 filepath = "WToQQ_070.root" #change if needed
3 tree = uproot.open(filepath)['tree']
4 tree.arrays()
```

Figure 7: Jupyter Notebook cell: inspection of data structure in .root file

output:

```
{part_px: [62.5, 68.2, ..., 0.946, 1.21], part_py: [215, ...], ...},
 {part_px: [-109, -82, ..., -1.39, -0.967], part_py: [...], part_pz: ..., ...},
 ...,
 {part_px: [-159, -81.3, ..., -0.957], part_py: [-74.2, ...], ...}]
```

```
-----
type: 100000 * {
  part_px: var * float32,
  part_py: var * float32,
  part_pz: var * float32,
  part_energy: var * float32,
  part_deta: var * float32,
  part_dphi: var * float32,
  part_d0val: var * float32,
  part_d0err: var * float32,
  part_dzval: var * float32,
  part_dzerr: var * float32,
  part_charge: var * float32,
  part_isChargedHadron: var * int32,
  part_isNeutralHadron: var * int32,
  part_isPhoton: var * int32,
  part_isElectron: var * int32,
  part_isMuon: var * int32,
  label_QCD: float32,
  label_Hbb: bool,
  label_Hcc: bool,
  label_Hgg: bool,
  label_H4q: bool,
```



```

label_Hqq1: bool,
label_Zqq: int32,
label_Wqq: int32,
label_Tbqq: int32,
label_Tbl: int32,
jet_pt: float32,
jet_eta: float32,
jet_phi: float32,
jet_energy: float32,
jet_nparticles: float32,
jet_sdmass: float32,
jet_tau1: float32,
jet_tau2: float32,
jet_tau3: float32,
jet_tau4: float32,
aux_genpart_eta: float32,
aux_genpart_phi: float32,
aux_genpart_pid: float32,
aux_genpart_pt: float32,
aux_truth_match: float32 }

```

Figure 8: Data structure in .root file

In terms of reading the JetClass data set's .root files, the procedure in Figure 9 can be used. This procedure utilizes Python lists as data structures with flexible shapes, which then point to the fact that not each jet has the same number of particles.

code.py:

```

1 import uproot, vector
2 filepath = "WToQQ_070.root" #change if needed
3 tree = uproot.open(filepath)['tree']
4 jets = {}
5 components = ['px', 'py', 'pz', 'energy', 'dphi', 'deta', 'isChargedHadron', '
    isNeutralHadron', 'isPhoton', 'isElectron', 'isMuon']
6 for component in components:
7     label = 'part_' + component
8     branch = tree[label]
9     jets[label] = [jet for jet in branch.array().tolist()]
10 p4 = vector.zip({'px': jets['part_px'],
11                 'py': jets['part_py'],
12                 'pz': jets['part_pz'],
13                 'energy': jets['part_energy']})
14 jets['part_pt']=p4.pt.tolist()
15 jets['part_eta']=p4.eta.tolist()
16 jets['part_phi']=p4.phi.tolist()
17 components = ['pt', 'eta', 'phi', 'nparticles', 'energy']

```

```

18 for component in components:
19     label = 'jet_' + component
20     branch = tree[label]
21     jets[label] = [jet for jet in branch.array().tolist()]

```

Figure 9: Reading from .root file [3]

Another way of reading the data from JetClass data set's .root files is by using `read_file` function of `dataloader.py` which can be found [here](#). This procedure involves zero-padding in order for each jet to have the same number of particles as to then enable using numpy arrays of fixed shapes for better code efficiency. The maximum number of particles can be set and each jet that has fewer particles than the set maximum is zero-padded while each jet that has more particles than the set maximum is truncated. A byproduct of zero-padding are dummy constituents which have all their features equal to zero. Each particle with its transverse momentum being zero will be considered a dummy constituent because, while technically a real particle can have zero transverse momentum, it is extremely rare. Moreover, in this report, as will be seen, particles with zero transverse momentum do not contribute to any calculations. **Important note: for the purposes of this report, the calculation and return of the one-hot encoded numpy array holding the information about the origin of a jet in `read_file` function are commented out.**

code.py:

```

1 import sys
2 sys.path.append('./')
3
4 import numpy as np
5 from dataloader import read_file
6
7 input="WToQQ_070.root" #change if needed
8
9 particle_features=['part_pt','part_eta','part_phi',
10                  'part_deta','part_dphi','part_energy',
11                  'part_isChargedHadron','part_isNeutralHadron',
12                  'part_isElectron','part_isMuon','part_isPhoton']
13
14 jet_features=['jet_pt','jet_eta','jet_phi']
15
16 x_particles, x_jets = read_file(
17     filepath=input,particle_features=particle_features,
18     jet_features=jet_features)
19
20 mask = x_particles[:, 0, :] == 0
21
22 full_mask = np.repeat(mask[:, np.newaxis, :], x_particles.shape[1], axis=1)

```

```

23 x_particles=np.ma.masked_where(full_mask,x_particles)
24
25 print(f"x_particles.shape={x_particles.shape}")
26 #(num_jets,num_particle_features,num_particles)
27 print(f"\n",flush=True)
28
29 print(f"x_jets.shape={x_jets.shape}")
30 #(num_jets,num_jet_features)
31 print(f"\n",flush=True)
32
33 jets={}
34
35 for pf,label in enumerate(particle_features):
36     jets[label]=x_particles[:,pf,:]
37
38 for jf,label in enumerate(jet_features):
39     jets[label]=x_jets[:,jf]
40
41 jets['jet_nparticles']=jets['part_pt'].count(axis=1)

```

Figure 10: Reading from any JetClass data set's .root file [1]

As commented in `dataloader.py` the numpy array *x_particles*, which holds information about particle level features, is of shape $(num_jets, num_particles_features, max_num_particles)$ while the numpy array *x_jets*, which holds the information about jet level features, is of shape $(num_jets, num_jet_features)$. The purpose of masking is to mark each dummy constituent as irrelevant. In order to see how masking in the above figure works, one can run the following code in Figure 11 the output of which can be seen in Figure 12.

code.ipynb:

```

1 import numpy as np
2 #two matrices in c represent two jets
3 #each column represents each particle
4 #each row represents different feature
5 c=np.array([[1,0,5,7],[0,6,0,1],[0,2,3,2]],[[0,4,9,0],[0,7,5,2],[6,3,8,4]])
6 print("c:")
7 print(c)
8 print("-----")
9 mask=c[:,0,:]==0
10 print("mask:")
11 print(mask)
12 print("-----")
13 full_mask=np.repeat(mask[:,np.newaxis,:],c.shape[1],axis =1)
14 print("full_mask:")
15 print(full_mask)
16 print("-----")

```

```

17 c_masked=np.ma.masked_where(full_mask,c)
18 print("c_masked:")
19 print(c_masked)

```

Figure 11: Helper cell for illustration of masking

output:

```

c:
[[[1 0 5 7]
  [0 6 0 1]
  [0 2 3 2]]

 [[0 4 9 0]
  [0 7 5 2]
  [6 3 8 4]]]

-----
mask:
[[False  True False False]
 [ True False False  True]]

-----
full_mask:
[[[False  True False False]
  [False  True False False]
  [False  True False False]]

 [[ True False False  True]
 [ True False False  True]
 [ True False False  True]]]

-----
c_masked:
[[[1 -- 5 7]
  [0 -- 0 1]
  [0 -- 3 2]]

 [[-- 4 9 --]
  [-- 7 5 --]
  [-- 3 8 --]]]

```

Figure 12: Output of the helper cell

6.2 Rodem data set

Reading from Rodem data set's .h5 files is described in Figure 13 below and the data is also zero-padded.

code.py:

```
1 import h5py
2 import numpy as np
3
4 def load_cnsts_rodem(ifile: str, mask_pt: list):
5     """Load constituents from an HDF5 file."""
6     with h5py.File(ifile, "r") as f:
7
8         cnsts=f["objects/jets/jet1_cnsts"][mask_pt,:,:3] #first three features
9
10        mask = np.repeat(cnsts[:, :, 0] == 0, cnsts.shape[2])
11        mask = mask.reshape(-1, cnsts.shape[1], cnsts.shape[2])
12        cnsts = np.ma.masked_where(mask, cnsts)
13
14        print(f"cnsts.shape={cnsts.shape}")
15        #(num_jets,num_particles,num_particle_features)
16        print(f"\n",flush=True)
17
18        return cnsts[:, :, 0], cnsts[:, :, 1], cnsts[:, :, 2] #pt,eta,phi
19
20 def load_jets_rodem(ifile: str, mask_pt: list):
21     """Load jets from an HDF5 file."""
22     with h5py.File(ifile, "r") as f:
23
24         jets = f["objects/jets/jet1_obs"][mask_pt,:3] #first three features
25
26         print(f"jets.shape={jets.shape}")
27         #(num_jets,num_jet_features)
28         print(f"\n",flush=True)
29
30         return jets[:,0], jets[:,1], jets[:,2] #pt,eta,phi
31
32 #the input HDF5 file containing the relevant jets:
33 input = "QCDjj_pT_450_1200_train01.h5" #change if needed
34
35 #the number of jets to load:
36 n_jets = 400_000 #change if needed
37
38 #jet pt range:
39 min_pt=450 #change if needed
40 max_pt=1000 #change if needed
41
42 with h5py.File(input, "r") as f:
43     mask_pt=[index for index,value in enumerate(f["objects/jets/jet1_obs"][:
44         n_jets,0]) if ((value>=min_pt) and (value<=max_pt))]
45
46 jets={}
47
```

```

47 #rows -> jets ; columns -> particles
48 jets['part_pt'],jets['part_eta'],jets['part_phi']=load_cnsts_rodem(input,mask_pt
    )
49
50 #flat arrays
51 jets['jet_pt'],jets['jet_eta'],jets['jet_phi']=load_jets_rodem(input,mask_pt)
52
53 jets['jet_nparticles']=jets['part_pt'].count(axis=1)
54
55 print(f"Smallest_jet_pt_is:{np.min(jets['jet_pt'])}")
56 print(f"Biggest_jet_pt_is:{np.max(jets['jet_pt'])}",flush=True)

```

Figure 13: Reading from any Rodem data set's .h5 file [4]

To understand and visualize the masking done in the figure above it is best to run the cell given in Figure 14 of which the output can be seen in Figure 15.

code.ipynb:

```

1 import numpy as np
2 #two matrices in c represent two jets
3 #each row represents each particle
4 #each column represents different feature
5 c=np.array([[1,2,5],[0,6,0],[0,2,3]],[[6,4,9],[0,7,5],[6,3,8]])
6 print("c:")
7 print(c)
8 print("-----")
9 mask=np.repeat(c[:,0]==0, c.shape[2])
10 print("mask:")
11 print(mask)
12 print("-----")
13 mask_shaped=mask.reshape(-1, c.shape[1], c.shape[2])
14 print("mask_shaped:")
15 print(mask_shaped)
16 print("-----")
17 c_masked=np.ma.masked_where(mask_shaped,c)
18 print("c_masked:")
19 print(c_masked)
20 print("-----")
21 count1=c_masked.count(axis=1)
22 print("count1:")
23 print(count1)
24 print("-----")
25 count2=c_masked.count(axis=1)[:0]
26 print("count2:")
27 print(count2)

```

Figure 14: Helper cell for illustration

```
output:

c:
[[[1 2 5]
  [0 6 0]
  [0 2 3]]

 [[6 4 9]
  [0 7 5]
  [6 3 8]]]

-----

mask:
[False False False  True  True  True  True  True  True  False False False
  True  True  True False False False]

-----

mask_shaped:
[[[False False False]
  [ True  True  True]
  [ True  True  True]]

 [[False False False]
  [ True  True  True]
  [False False False]]]

-----

c_masked:
[[[1 2 5]
  [-- -- --]
  [-- -- --]]

 [[6 4 9]
  [-- -- --]
  [6 3 8]]]

-----

count1:
[[1 1 1]
 [2 2 2]]

-----

count2:
[1 2]
```

Figure 15: Output of the helper cell

6.3 Aspen Open Jets data set

Reading from Aspen Open Jets data set's .h5 files is described in Figure 16 below and, like the previous data set, contains zero-padded data.

code.py:

```
1 import h5py
2 import vector
3 import awkward as ak
4 import numpy as np
5
6 def load_cnsts_aspen(iframe: str, mask_pt: list):
7     """Load constituents from an HDF5 file."""
8     with h5py.File(iframe, "r") as f:
9
10         PFCands = f["PFCands"][mask_pt, :, :4] #px,py,pz,E
11
12         p4_cnsts=vector.zip({'px':PFCands[:, :, 0],
13                             'py':PFCands[:, :, 1],
14                             'pz':PFCands[:, :, 2],
15                             'energy':PFCands[:, :, 3]})
16
17         cnsts=np.zeros([PFCands.shape[0],PFCands.shape[1],3])
18         cnsts[:, :, 0]=ak.to_numpy(p4_cnsts.pt, allow_missing=True)
19         cnsts[:, :, 1]=ak.to_numpy(p4_cnsts.eta, allow_missing=True)
20         cnsts[:, :, 2]=ak.to_numpy(p4_cnsts.phi, allow_missing=True)
21
22         #masking now done the same way as for Rodem data set
23         mask = np.repeat(cnsts[:, :, 0] == 0, cnsts.shape[2])
24         mask = mask.reshape(-1, cnsts.shape[1], cnsts.shape[2])
25         cnsts = np.ma.masked_where(mask, cnsts)
26
27         print(f"cnsts.shape={cnsts.shape}")
28         #(num_jets,num_particles,num_particle_features)
29         print(f"\n", flush=True)
30
31         return cnsts[:, :, 0], cnsts[:, :, 1], cnsts[:, :, 2] #pt,eta,phi
32
33 def load_jets_aspen(iframe: str, mask_pt: list):
34     """Load jets from an HDF5 file."""
35     with h5py.File(iframe, "r") as f:
36
37         jets = f["jet_kinematics"][mask_pt, :3]
38
39         print(f"jets.shape={jets.shape}")
40         #(num_jets,num_jet_features)
41         print(f"\n", flush=True)
42
43         return jets[:, 0], jets[:, 1], jets[:, 2] #pt,eta,phi
44
45 #the input HDF5 file:
46 input = "RunG_batch0.h5"
47
```



```

48 #the number of jets to load:
49 n_jets = 400_000
50
51 #jet pt range:
52 min_pt=300 #change if needed
53 max_pt=500 #change if needed
54
55 jets={}
56
57 with h5py.File(input, "r") as f:
58
59     mask_pt=[index for index,value in enumerate(f["jet_kinematics"][:n_jets,0])
60              if ((value>=min_pt) and (value<=max_pt))]
61
62     #rows -> jets ; columns -> particles
63     jets['part_pt'],jets['part_eta'],jets['part_phi']=load_cnsts_aspen(input,
64                                mask_pt)
65
66     #flat arrays
67     jets['jet_pt'],jets['jet_eta'],jets['jet_phi']=load_jets_aspen(input,mask_pt
68                                )

```

Figure 16: Reading from any Aspen Open Jets data set's .h5 file [5]

7 JetClass data set: visualization of jets

After reading a file, with the help of the code in Figure 17, jets can be visualized as in Figure 18 below. Every constituent of each jet in JetClass data set's .root files is known to be a specific kind of particle through one-hot encoding.

Important: The code snippet below is a continuation of the code used to read the file.

code.py:

```

1 import os
2 import matplotlib.pyplot as plt
3
4 if not os.path.exists("./Plots"):
5     os.makedirs("./Plots")
6
7 arg=np.argmax(jets['jet_pt'])
8 print(f"index_of_the_jet_with_the_highest_pT:{arg}")
9 print(f"\n",flush=True)
10
11 fig,ax=plt.subplots(figsize=(20,12))
12
13 for i in range(jets['jet_nparticles'][arg]):

```

```

14     if jets['part_isChargedHadron'][arg,i]==1:
15         m='o'
16         fcs='lightgreen'
17     elif jets['part_isNeutralHadron'][arg,i]==1:
18         m='o'
19         fcs='none'
20     elif jets['part_isElectron'][arg,i]==1:
21         m='^'
22         fcs='lightgreen'
23     elif jets['part_isMuon'][arg,i]==1:
24         m='v'
25         fcs='lightgreen'
26     elif jets['part_isPhoton'][arg,i]==1:
27         m='p'
28         fcs='none'
29     ax.scatter(jets['part_deta'][arg,i],jets['part_dphi'][arg,i],color='
        lightgreen',marker=m,facecolors=fcs,s=jets['part_energy'][arg,i])
30
31 ax.set_xlabel(r"$\Delta\eta$")
32 ax.set_ylabel(r"$\Delta\phi$")
33 fig.savefig("./Plots/jet.jpg")

```

Figure 17: Plotting for visualization of a jet

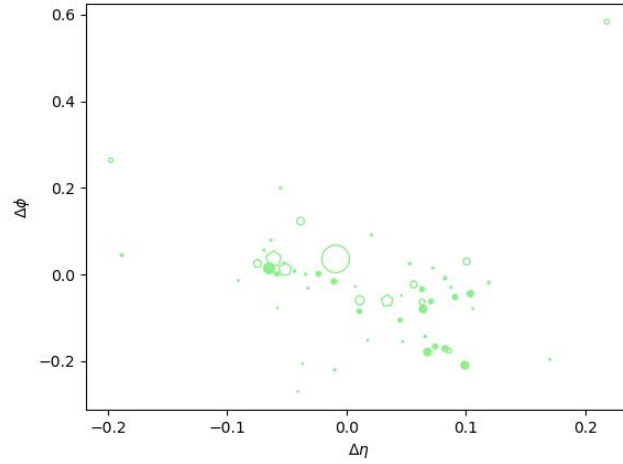


Figure 18: Jet with the highest p_T from WToQQ_070.root file. The circles, triangles (upward- or downward-directed), and pentagons represent the particle types, which are hadrons, leptons (electrons or muons), and photons, respectively. The solid (hollow) markers stand for electrically charged (neutral) particles [1].

8 Probability distribution of the radius

8.1 Contributions to the probability - 1D problem

The contribution of a particle that falls within some radius interval $[R_i - \Delta R, R_i + \Delta R]$ to the probability of the given interval is proportional to its transverse momentum because particles with higher transverse momentum contribute more significantly to the total momentum of the jet. These higher-energy particles are more likely to be involved in the jet structure and are therefore more likely to contribute to the probability in the specified radius interval. To obtain the probability plot of a certain jet energy or jet p_T range, numerous jets have to be considered to make a histogram and obtain the appropriate statistics. The jet's contribution to the probability of the given radius interval is then proportional to the sum of the transverse momenta of its particles that fall within the interval:

$$w_{ik(J)}^p = p_{T_{ik(J)}} \quad (9)$$

$$H_{i(J)}^p = \sum_{k(counts)} w_{ik(J)}^p, \quad (10)$$

where $p_{T_{ik(J)}}$ is the transverse momentum value of a particle that belongs to jet J and whose value R falls in the interval $[R_i - \Delta R, R_i + \Delta R]$. **For jet J , there are n_1 particles in total, and n_2 of these particles have the radius that falls within the interval $[R_i - \Delta R, R_i + \Delta R]$. The index k then runs from 1 to n_2 , where each value of k corresponds to a particular particle within the jet and the interval.** In other words, $H_{i(J)}^p$ represents the unscaled contribution of a jet J to the probability of the radius interval $[R_i - \Delta R, R_i + \Delta R]$. The contributions to the probability of each jet to some radius interval have to be summed up in a certain way to ensure that the sum of all probability contributions from the whole radius range is equal to one.

$$h_i^p = \frac{1}{N(jets)} \sum_J \frac{H_{i(J)}^p}{\sum_{m(bins)} H_{m(J)}^p}. \quad (11)$$

With the applied normalization, all of the particle contributions to probability will sum up to one no matter how many jets are taken into account, while the contributions from dissimilar jets are kept uncorrelated at the same time.

8.2 Contributions to the probability density - 2D problem

Taking into account what was said in the previous subsection, the probability density at a certain value of R also must account for the fact that R is a radius in $(\Delta\eta, \Delta\phi)$ plane. The particle's contribution to the probability density of the interval $[R_i - dR, R_i + dR]$ has to stay proportional to its transverse momentum, but it has to be divided by the value of R_i because the infinitesimal surface area in a two-dimensional plane is calculated as $R \cdot dR \cdot d\alpha$, where α is just the symbol taken for the azimuthal angle in $(\Delta\eta, \Delta\phi)$ plane. The probability density of interest, called $\rho(R)$, is the one already integrated over the azimuthal angle.

$$\int_R \int_\alpha \rho(R, \alpha) \cdot R \cdot dR \cdot d\alpha = \int_R \rho(R) \cdot R \cdot dR \quad (12)$$

As with probability, to obtain the probability density plot of a certain jet energy or jet p_T range, numerous jets have to be considered in order to make a histogram and get appropriate statistics. As the contribution of a jet to the probability of a given radius, the contribution of a jet to the probability density in the interval $[R_i - \Delta R, R_i + \Delta R]$ must be proportional to the sum of transverse momenta of all its particles whose radius falls within the interval and, as described above, must be divided by R_i :

$$w_{ik(J)}^d = \frac{p_{Tik(J)}}{R_i} \quad (13)$$

$$H_{i(J)}^d = \sum_{k(counts)} w_{ik(J)}^d. \quad (14)$$

In other words, $H_{i(J)}^d$ represents the unscaled contribution of a jet J to the probability density of the interval $[R_i - \Delta R, R_i + \Delta R]$.

In order to ensure that the expression (12) is equal to one, the histogram must be normalized in a certain way:

$$h_i^d = \frac{1}{N(jets)} \sum_J \frac{H_{i(J)}^d}{\Delta R \cdot \sum_{m(bins)} R_m H_{m(J)}^d}, \quad (15)$$

so that in the end the result is as follows:

$$\int_R \rho(R) \cdot R \cdot dR \rightarrow \sum_{i(bins)} h_i^d \cdot R_i \cdot \Delta R = \quad (16)$$

$$= \sum_J \frac{1}{N(jets)} \sum_{i(bins)} \frac{H_{i(J)}^d}{\Delta R \cdot \sum_{m(bins)} R_m H_{m(J)}^d} \cdot R_i \cdot \Delta R = \sum_J \frac{1}{N(jets)} = 1. \quad (17)$$

Observe again that the contributions from dissimilar jets are kept uncorrelated.

Moreover, it is now clear that the expression (15) is directly connected to the expression (11):

$$h_i^d \cdot R_i \cdot \Delta R = \quad (18)$$

$$= \frac{1}{N(jets)} \sum_J \frac{H_{i(J)}^d}{\Delta R \cdot \sum_{m(bins)} R_m H_{m(J)}^d} \cdot R_i \cdot \Delta R \quad (19)$$

$$= \frac{1}{N(jets)} \sum_J \frac{H_{i(J)}^p}{\sum_{m(bins)} H_{m(J)}^p} \quad (20)$$

$$= h_i^p. \quad (21)$$

8.3 Histogram and Poisson distribution relation - in general

If the number of events is huge ($n \gg$) and the probability for some event to happen is very small ($p \ll$) then one is talking about a Poisson distribution. When falling into a bin of non-weighted histogram with a lot of counts ($n \gg$) is considered an event with small probability ($p \ll$) then the variance of this bin's value follows the expression for Poisson distribution and is equal to $Var = n \cdot p =$

$n(i)$, where $n(i)$ is the number of samples that fall into that specific bin. On the other hand, in the case where there is a histogram with weights, then one has a sum of weighted Poisson distributions. Here, $n \gg$ is the number of samples with the specific weight, while $p \ll$ is the probability that these samples will fall in the specific bin. In order to make the connection between weighted histograms and a linear combination of random variables, it is necessary to introduce $X_k^{(i)}$ as a random variable which for samples with weight value w_k gives 1 when they fall into the i -th bin and 0 otherwise. For the purpose of this section, all the weights $w_k^{(i)}$ whose counts fall into the i -th bin are arranged so that only the first $m(i)$ weights have dissimilar values:

$$\{w_1^{(i)}, \dots, w_{m(i)}^{(i)}, \dots, w_{n(i)}^{(i)}; m(i) \leq n(i)\}. \quad (22)$$

It is now possible to define a new random variable for each bin that represents a final histogram value:

$$X^{(i)} = \sum_{k=1}^{m(i)} w_k^{(i)} X_k^{(i)}, \quad (23)$$

For a linear combination of random uncorrelated variables $(X_1^{(i)}, \dots, X_{m(i)}^{(i)})$ with correspondingly different weights $(w_1^{(i)}, \dots, w_{m(i)}^{(i)})$ the variance is calculated in this manner:

$$Var(X^{(i)}) = \sum_{k=1}^{m(i)} w_k^{(i)2} Var(X_k^{(i)}), \quad (24)$$

where each $X_k^{(i)}$ is a Poisson distribution with variance equal to the number of samples in the i -th bin which have the same weight $w_k^{(i)}$ so essentially one is summing all of the squared weights whose counts fall into the i -th bin. Taking the latter into account, for the specific i -th bin with $n(i)$ counts, the variance is then equal to:

$$Var(X^{(i)}) = \sum_{k=1}^{n(i)} w_k^{(i)2}. \quad (25)$$

Normalization of a histogram, or analogously the weights, changes the situation and introduces bin correlations. The covariance matrix must be calculated to fit the data with the complete information. **Standard deviation of the probability for the specific bin in a weighted histogram is then equal to the square root of the sum of squared weights of all samples that fall into that bin.**

8.4 Covariance matrix - in general

The propagation formula for the covariance matrix will be of much use in this report, so here is the short proof.

$$Cov(X, Y) = E[(X - E[X])(Y - E[Y])] \quad (26)$$

For $f(X_1, \dots, X_n)$ the first-order Taylor expansion around $(E[X_1], \dots, E[X_n])$ follows:

$$f(X_1, \dots, X_n) \approx f(E[X_1], \dots, E[X_n]) + \sum_i \left. \frac{\partial f}{\partial X_i} \right|_{E[X_i]} (X_i - E[X_i]) \quad (27)$$

$$E[f] \approx f(E[X_1], \dots, E[X_n]) + \sum_i \left. \frac{\partial f}{\partial X_i} \right|_{E[X_i]} (E[X_i] - E[X_i]) = f(E[X_1], \dots, E[X_n]) \quad (28)$$

$$f(X_1, \dots, X_n) \approx E[f] + \sum_i \left. \frac{\partial f}{\partial X_i} \right|_{E[X_i]} (X_i - E[X_i]) \quad (29)$$

Now, for $f(X_1, \dots, X_n)$ and $g(X_1, \dots, X_n)$ in the Taylor expansion of first order, it stands:

$$Cov(f, g) = E[(f - E[f])(g - E[g])] = \quad (30)$$

$$= \sum_i \sum_j \left. \frac{\partial f}{\partial X_i} \right|_{E[X_i]} \left. \frac{\partial g}{\partial X_j} \right|_{E[X_j]} E[(X_i - E[X_i])(X_j - E[X_j])] \quad (31)$$

$$= \sum_i \sum_j \left. \frac{\partial f}{\partial X_i} \right|_{E[X_i]} \left. \frac{\partial g}{\partial X_j} \right|_{E[X_j]} Cov(X_i, X_j) \quad (32)$$

8.5 Covariance matrix of the probability (density) histogram

Now, the covariance matrix propagation formula can be utilized to calculate the covariance matrix of the probability (density) histogram where the correlations between bins are introduced due to normalization.

$$Cov(h_i, h_j) = \sum_{J'} \sum_{J''} \sum_{k \text{ (bins)}} \sum_{l \text{ (bins)}} \frac{\partial h_i}{\partial H_{k(J')}} \frac{\partial h_j}{\partial H_{l(J'')}} Cov(H_{k(J')}, H_{l(J'')}) = \quad (33)$$

$$= \sum_{J'} \sum_{k \text{ (bins)}} \frac{\partial h_i}{\partial H_{k(J')}} \frac{\partial h_j}{\partial H_{k(J')}} Var(H_{k(J')}) \quad (34)$$

Now, if one is calculating the covariances between bins of the probability histogram, the following expressions need to be followed.

$$\frac{\partial h_i^p}{\partial H_{k(J')}} = \begin{cases} \frac{1}{N(jets)} \frac{\sum_{l \text{ (bins)}} H_{l(J')}^p - H_{i(J')}^p}{\left(\sum_{l \text{ (bins)}} H_{l(J')}^p\right)^2} & \text{if } i=k \\ \frac{1}{N(jets)} \frac{-H_{i(J')}^p}{\left(\sum_{l \text{ (bins)}} H_{l(J')}^p\right)^2} & \text{if } i \neq k \end{cases}$$

On the other hand, if one is calculating the covariances between bins of the probability density histogram, then the following expressions need to be followed.

$$\frac{\partial h_i^d}{\partial H_{k(J')}} = \begin{cases} \frac{1}{N(jets) \cdot \Delta R} \frac{\sum_{l \text{ (bins)}} R_l H_{l(J')}^d - R_k H_{i(J')}^d}{\left(\sum_{l \text{ (bins)}} R_l H_{l(J')}^d\right)^2} & \text{if } i=k \\ \frac{1}{N(jets) \cdot \Delta R} \frac{-R_k H_{i(J')}^d}{\left(\sum_{l \text{ (bins)}} R_l H_{l(J')}^d\right)^2} & \text{if } i \neq k \end{cases}$$

Because there may not be that many particles in some J' jet, the variance of $H_{k(J')}$ whose bin itself can have only one count or none, needs to be looked at closely; it cannot just be taken naively as

a sum of squared weights in that bin. There is a trick, because essentially one is dealing with an ensemble of histograms or, in other words, with a lot of jet-wise partial histograms.

$$H_k = \sum_J H_{k(J)}; \quad k = 1, \dots, N \quad (35)$$

$$S_{k \setminus J'} = \sum_{J \neq J'} H_{k(J)}; \quad k = 1, \dots, N \quad (36)$$

$$H_k = H_{k(J')} + S_{k \setminus J'} \rightarrow \text{Cov}(H_k, S_{k \setminus J'}) = \text{Var}(S_{k \setminus J'}) \quad (37)$$

$$\text{Var}(H_{k(J')}) = \text{Var}(H_k - S_{k \setminus J'}) = \text{Var}(H_k) + \text{Var}(S_{k \setminus J'}) - 2 \cdot \text{Cov}(H_k, S_{k \setminus J'}) = \quad (38)$$

$$= \text{Var}(H_k) + \text{Var}(S_{k \setminus J'}) - 2 \cdot \text{Var}(S_{k \setminus J'}) = \text{Var}(H_k) - \text{Var}(S_{k \setminus J'}) \quad (39)$$

Both variances from the result can now be calculated as the sum of squared weights because they are essentially the bins of two weighted histograms ($\{H_k; k = 1, \dots, N\}$ and $\{S_{k \setminus J'}; k = 1, \dots, N\}$) with many counts. As can be seen, regardless of the fact that there may be just a few counts per bin for $\{H_{k(J')}; k = 1, \dots, N\}$ histogram, variance of its k -th bin is still the sum of squared weights even if there is only one count or none.

In the case of probability, the variance of $H_{k(J')}^p$ is then calculated as follows.

$$w_{km(J')}^p = p_{Tkm(J')} \quad (40)$$

$$H_{k(J')}^p = \sum_{m(\text{counts})} w_{km(J')}^p \quad (41)$$

$$\text{Var}(H_{k(J')}^p) = \sum_{m(\text{counts})} [w_{km(J')}^p]^2. \quad (42)$$

On the other hand, in the case of the probability density, the variance of $H_{k(J')}^d$ is calculated as follows:

$$w_{km(J')}^d = \frac{p_{Tkm(J')}}{R_k} \quad (43)$$

$$H_{k(J')}^d = \sum_{m(\text{counts})} w_{km(J')}^d \quad (44)$$

$$\text{Var}(H_{k(J')}^d) = \sum_{m(\text{counts})} [w_{km(J')}^d]^2. \quad (45)$$

It can be seen now that if the weights of all the counts of a jet J' are multiplied by the same constant, the result in (34) would not change:

$$w_{km(J')} \rightarrow \text{const}(J') \cdot w_{km(J')} \quad (46)$$

$$\frac{\partial h_i}{\partial H_{k(J')}} \rightarrow \frac{1}{\text{const}(J')} \cdot \frac{\partial h_i}{\partial H_{k(J')}} \quad (47)$$

$$\frac{\partial h_j}{\partial H_{k(J')}} \rightarrow \frac{1}{\text{const}(J')} \cdot \frac{\partial h_j}{\partial H_{k(J')}} \quad (48)$$

$$\text{Var}(H_{k(J')}) \rightarrow \text{const}(J')^2 \cdot \text{Var}(H_{k(J')}) \quad (49)$$

$$\frac{\partial h_i}{\partial H_{k(J')}} \frac{\partial h_j}{\partial H_{k(J')}} \text{Var}(H_{k(J')}) \rightarrow \frac{\partial h_i}{\partial H_{k(J')}} \frac{\partial h_j}{\partial H_{k(J')}} \text{Var}(H_{k(J')}) \quad (50)$$

Namely, if one chooses this constant in a smart way, it can make numerical calculations more stable and avoid dealing with summing big weights or their squares. The constant is chosen in such a way that the new weights are now:

$$[w_{km(J')}^p]' = \frac{w_{km(J')}^p}{N(\text{jets}) \cdot \sum_{l(\text{bins})} \sum_{n(\text{counts})} w_{ln(J')}^p}, \quad (51)$$

$$[w_{km(J')}^d]' = \frac{w_{km(J')}^d}{N(\text{jets}) \cdot \Delta R \cdot \sum_{l(\text{bins})} R_l \sum_{n(\text{counts})} w_{ln(J')}^d}. \quad (52)$$

It can also be seen that with these new weights in place, the calculations of h_k^p and h_k^d is now the straightforward sum of all weights whose counts fall into the interval $[R_k - \Delta R, R_k + \Delta R]$:

$$h_k^p = \sum_{J'} \sum_m [w_{km(J')}^p]' = \quad (53)$$

$$= \frac{1}{N(\text{jets})} \sum_{J'} \frac{\sum_m w_{km(J')}^p}{\sum_{l(\text{bins})} \sum_{n(\text{counts})} w_{ln(J')}^p} \quad (54)$$

$$= \frac{1}{N(\text{jets})} \sum_{J'} \frac{H_{k(J')}^p}{\sum_{l(\text{bins})} H_{l(J')}^p}, \quad (55)$$

$$h_k^d = \sum_{J'} \sum_m [w_{km(J')}^d]' = \quad (56)$$

$$= \frac{1}{N(\text{jets})} \sum_{J'} \frac{\sum_m w_{km(J')}^d}{\Delta R \cdot \sum_{l(\text{bins})} R_l \sum_{n(\text{counts})} w_{ln(J')}^d} \quad (57)$$

$$= \frac{1}{N(\text{jets})} \sum_{J'} \frac{H_{k(J')}^d}{\Delta R \cdot \sum_{l(\text{bins})} R_l H_{l(J')}^d}. \quad (58)$$

9 Explored modeling methods/tools

The objective was to identify an analytic function that could describe the probability distribution or probability density function of the radius of QCD jets. Furthermore, the aim was to develop a parametric fit: a single functional form that could model the R distributions across narrow jet p_T bins, with parameters that would also need to be fitted for each case. This chapter provides an overview of the methods and tools that were tested.

9.1 Bayesian machine scientist

Bayesian machine scientist is a powerful package for finding analytical fit functions for complex data [6]. Python implementation can be found [here](#), where a tutorial is also presented. In this subsection, a simple example of usage is given. As always, one has to first create a shell script.

WSL terminal:

```
nano job.sh
```

job.sh:

```
#!/bin/bash

#SBATCH --job-name=steve_job

#SBATCH --output=output.out

#SBATCH --error=error.err

#SBATCH --time=00:10:00

#SBATCH --mem=2G

#SBATCH --partition=short

#SBATCH --cpus-per-task=1

python code.py
```

WSL terminal:

```
nano code.py
```

code.py:

```
1 import os
2 import sys
3 sys.path.append('./machine/') #downloaded files are here
4
5 import warnings
6 warnings.filterwarnings('ignore')
7 import numpy as np
8 import pandas as pd
9 import matplotlib.pyplot as plt
10
11 from copy import deepcopy
12 from mcmc import * #downloaded mcmc.py file
13 from parallel import * #downloaded parallel.py file
14 from fit_prior import read_prior_par #downloaded read_prior_par.py file
15
16 print(sys.version) #3.12.3
17 print(f"\n",flush=True)
18
19 num_points = 20
20 x1 = np.random.uniform(-5, 5, num_points)
21 x2 = np.random.uniform(-5, 5, num_points)
22
23 # True function
24 y_true = 0.4 * x1**2 - 3.2167 * x2**3
25
26 # Add noise (Gaussian noise with mean 0 and standard deviation 1)
27 noise = np.random.normal(0, 1, size=num_points)
28 y_noisy = y_true + noise
29
30 x=np.zeros([num_points,2])
31 x[:,0]=x1
32 x[:,1]=x2
33
34 x=pd.DataFrame(data=x,columns=['x1','x2'])
35 y=pd.Series(data=y_noisy)
36 XLABS=['x1','x2']
37
38 # Read the hyperparameters for the prior from downloaded file
39 file='./machine/final_prior_param_sq.named_equations.nv13.np13.2016-09-01_17
    _05_57.196882.dat'
40 prior_par = read_prior_par(file)
41
42 # Set the temperatures for the parallel tempering
43 Ts = [1] + [1.04**k for k in range(1, 20)]
44
45 # Initialize the parallel machine scientist
46 pms = Parallel(
```

```

47     Ts,
48     variables=XLABS,
49     parameters=['a%d' % i for i in range(13)],
50     x=x, y=y,
51     prior_par=prior_par,
52 )
53
54 # Number of MCMC steps
55 nstep = 3000
56
57 # MCMC
58 description_lengths, mdl, mdl_model = [], np.inf, None
59 for i in range(nstep):
60     # MCMC update
61     pms.mcmc_step() # MCMC step within each T
62     pms.tree_swap() # Attempt to swap two randomly selected consecutive temps
63     # Add the description length to the trace
64     description_lengths.append(pms.tl.E)
65     # Check if this is the MDL expression so far
66     if pms.tl.E < mdl:
67         mdl, mdl_model = pms.tl.E, deepcopy(pms.tl)
68     # Keep track of the progress
69     #if (i+1) % 100 == 0:
70     #    print(f"Step {i+1}/{nstep}")
71     #    print(f"\n",flush=True)
72
73 print("Best_model:\t", mdl_model)
74 print("Desc._length:\t", mdl)
75 print("Parameter_values:")
76 print(mdl_model.par_values)
77 print(f"\n",flush=True)
78
79 if not os.path.exists("./Plots"):
80     os.makedirs("./Plots")
81
82 plt.figure(figsize=(6, 6))
83 plt.scatter(mdl_model.predict(x), y)
84 plt.plot((y_noisy.min(),y_noisy.max()),(y_noisy.min(),y_noisy.max()))
85 plt.xlabel('MDL_model_predictions', fontsize=14)
86 plt.ylabel('Actual_values', fontsize=14)
87 plt.savefig("./Plots/model.jpg")

```

WSL terminal:

```
sbatch job.sh
```

WSL terminal:

```
nano output.out
```

output.out:

```
3.12.3 | packaged by conda-forge | (main, Apr 15 2024, 18:38:13) [GCC 12.3.0]

Best model:      (((x2 * _a9_) ** 3) + ((x1 ** 2) * _a3_))
Desc. length:    65.4369873882440
Parameter values:
{'d0': {'_a3_': np.float64(0.39042008279477214), '_a9_': np.float64
(-1.4766385375672322), '_a0_': 1.0, '_a1_': 1.0, '_a2_': 1.0, '_a4_': 1.0, '
_a5_': 1.0, '_a6_': 1.0, '_a7_': 1.0, '_a8_': 1.0, '_a10_': 1.0, '_a11_':
1.0, '_a12_': 1.0}}
```

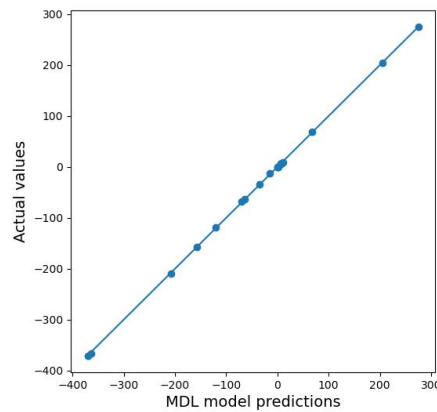


Figure 19: Quality of the returned model: model.jpg

However, as will be shown, machine scientist can sometimes return a rather unwieldy function. Moreover, since it does not account for the data uncertainties, the resulting function does not guarantee stable behavior when used later in error-weighted fitting.

9.2 Quantile transformation

Quantile transformation in Python is a technique that transforms features to follow a uniform or normal distribution by mapping the original data to quantiles. Quantiles are cut points that divide the range of a probability distribution into continuous intervals, each containing an equal probability mass. The following script helps with visualization of what the procedure can do. The scripts in this subsection were run in the Jupyter Notebook on a local computer for simplicity.

quantiles1.ipynb:

```
1 import os
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from numpy import exp
5 from sklearn.preprocessing import QuantileTransformer
6
7 # **GENERATE SOME DATA:**
8
9 # Set a seed for reproducibility
10 np.random.seed(10)
11
12 n_samples=1000
13 # Generate Gaussian data sample
14 data = np.random.randn(n_samples)
15
16 # Add a skew to the data distribution
17 data = exp(data)
18
19 #-----
20
21 # **PLOTting THE HISTOGRAM OF THE ORIGINAL DATA:**
22
23 if not os.path.exists("./Plots"):
24     os.makedirs("./Plots")
25
26 # Histogram of the raw data with a skew
27 plt.figure(figsize=(6,6))
28 plt.hist(data, bins=25, ec='black')
29 plt.xlabel("data")
30 plt.ylabel("counts")
31 plt.savefig("./Plots/data.jpg")
32
33 #-----
34
35 # **PLOTting THE HISTOGRAM OF THE TRANSFORMED DATA:**
36
37 # Reshape data to have rows and columns
38 data = data.reshape(-1,1)
39
40 # Quantile transform the raw data
41 qt = QuantileTransformer(output_distribution='normal')
42 data_trans = qt.fit_transform(data)
43
44 # Histogram of the transformed data
45 plt.figure(figsize=(6,6))
46 plt.hist(data_trans, bins=25, ec='black')
47 plt.xlabel("data_trans")
```

```

48 plt.ylabel("counts")
49 plt.savefig("./Plots/data_trans.jpg")

```

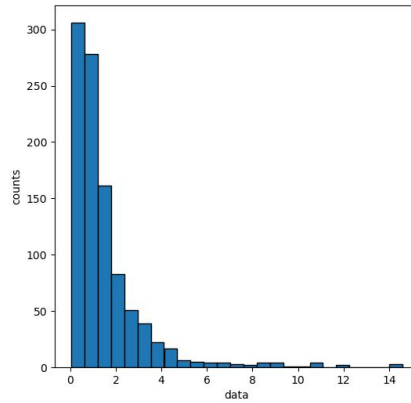


Figure 20: Original skewed data: data.jpg

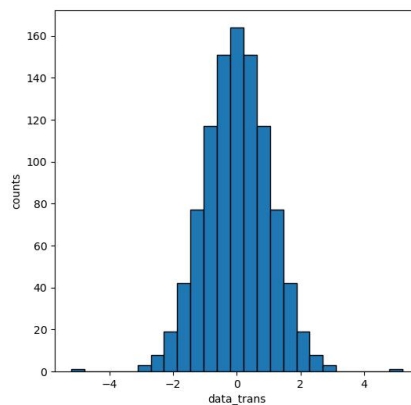


Figure 21: Transformed data: data_trans.jpg

The way this method calculates quantiles is described in detail in the following explanation and then checked in the code.

Quantiles Calculation

1. qt.fit procedure sorts the data \rightarrow sorted_data
2. splitting indices of sorted_data in (n-quantiles-1) parts
 \rightarrow step between two subsequent indices: $(n_samples-1)/(n_quantiles-1)$
3. get the indices of elements in sorted_data which will go in qt.quantiles_
 $\rightarrow \{0, 1, 2, \dots, k, \dots, n_quantiles-1\} \cdot (n_samples-1)/(n_quantiles-1) =$
 $= \{idx(0), idx(1), \dots, idx(k), \dots, idx(n_quantiles-1)\}$
 \rightarrow observe: $idx(0)=0$ and $idx(n_quantiles-1)=n_samples-1$
 \rightarrow observe: if $n_quantiles=n_samples \rightarrow idx(k)=k \ (\forall k)$

4. if for example $\text{idx}(k)=3.7$ then:

$\text{qt.quantiles_}[k,0]=\text{sorted_data}[3]+0.7\cdot(\text{sorted_data}[4]-\text{sorted_data}[3])$

→ observe: if $n_{\text{quantiles}}=n_{\text{samples}} \rightarrow \text{idx}(k)=k \ (\forall k) \rightarrow \text{qt.quantiles_}[:,0]=\text{sorted_data}$

In addition to helping one understand how quantiles are calculated, the following script shows how the transformation of data is then performed.

quantiles2.ipynb:

```
1 import random
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from numpy import exp
5 from scipy.stats import norm
6 from sklearn.preprocessing import QuantileTransformer
7
8 # **GENERATE SOME DATA:**
9
10 # Set a seed for reproducibility
11 np.random.seed(10)
12
13 n_samples=1000
14 # Generate Gaussian data sample
15 data = np.random.randn(n_samples)
16 # Add a skew to the data distribution
17 data = exp(data)
18 # Reshape data to have rows and columns
19 data = data.reshape(-1,1)
20
21 #-----
22
23 # **CALCULATING THE QUANTILES AND PERCENTILES:**
24
25 # Setting the transformer
26 n_quantiles=6
27 qt = QuantileTransformer(
28     output_distribution='normal', subsample=None, n_quantiles=n_quantiles)
29
30 # Finding the quantiles
31 qt.fit(X=data.reshape(-1,1))
32
33 quantiles=qt.quantiles_.flatten()
34 percentiles=np.zeros([n_quantiles])
35
36 idx=[((n_samples-1)/(n_quantiles-1))*k for k in range(n_quantiles)]
37 quantiles_check=np.zeros([n_quantiles])
38 sorted_data=np.sort(data,axis=0).flatten()
39
```

```

40 for k in range(n_quantiles):
41
42     percentiles[k]=k/(n_quantiles-1)
43
44     print(f"{k}. quantile: {quantiles[k]}-> percentile: {percentiles[k]*100}%")
45
46     if idx[k]!=int(idx[k]):
47
48         # Manual check:
49         percentage=(idx[k]-int(idx[k]))
50         step=sorted_data[int(idx[k])+1]-sorted_data[int(idx[k])]
51         addition=percentage*step
52         quantiles_check[k]=sorted_data[int(idx[k])]+addition
53
54     else:
55
56         # Manual check:
57         quantiles_check[k]=sorted_data[int(idx[k])]
58
59     print(f"quantile manual check: {quantiles_check[k]}")
60
61 print(f"\n")
62
63 #-----
64
65 # **TRANSFORMING THE QUANTILES:**
66
67 for k in range(n_quantiles):
68
69     print(f"Percentile point function check for {k}. quantile:")
70     print("qt.transform(quantile):")
71     print(qt.transform(np.array([quantiles[k]]).reshape(-1,1)))
72
73     if k==0:
74         print(f"norm.ppf(percentile):")
75         print(norm.ppf(percentiles[k]+1e-7,loc=0,scale=1.0))
76     elif k==n_quantiles-1:
77         print(f"norm.ppf(percentile):")
78         print(norm.ppf(percentiles[k]-1e-7,loc=0,scale=1.0))
79     else:
80         print(f"norm.ppf(percentile):")
81         print(norm.ppf(percentiles[k],loc=0,scale=1.0))
82
83     print(f"\n")
84
85 #-----
86
87 # **PLAYING WITH A POINT THAT IS POSITIONED BETWEEN TWO QUANTILES:**

```



```

88
89 random.seed(10)
90
91 # Point that needs to be transformed
92 # is somewhere between quantiles[k] and quantiles[k+1]:
93 k=random.randint(0,n_quantiles-2)
94 num=random.random()
95 x=quantiles[k]+num*(quantiles[k+1]-quantiles[k])
96
97 # Percentile is calculated with linear interpolation
98 a=percentiles[k]
99 b=((percentiles[k+1]-percentiles[k])/(quantiles[k+1]-quantiles[k]))
100 x0=quantiles[k]
101 percentile=a+b*(x-x0)
102 x_transformed=norm.ppf(percentile)
103
104 print(f"x={x} -> percentile(x)={percentile}")
105 print(f"\n")
106 print(f"qt.transform(x):")
107 print(qt.transform(np.array(x).reshape(-1,1)))
108 print("x_transformed:")
109 print(x_transformed)
110 print(f"\n")
111
112 # Transform the point back to the original space
113 percentile=norm.cdf(x_transformed)
114 x_back=(percentile-a)/b+x0
115
116 print("qt.inverse_transform(qt.transform(x)):")
117 print(qt.inverse_transform(qt.transform(np.array(x).reshape(-1,1))))
118 print("x_back:")
119 print(x_back)

```

output:

```

0. quantile: 0.04058318972369623 -> percentile: 0.0%
quantile manual check: 0.04058318972369623
1. quantile: 0.45261614214031254 -> percentile: 20.0%
quantile manual check: 0.45261614214031254
2. quantile: 0.7880576872337507 -> percentile: 40.0%
quantile manual check: 0.7880576872337507
3. quantile: 1.2357563673533 -> percentile: 60.0%
quantile manual check: 1.2357563673533
4. quantile: 2.1275282556894264 -> percentile: 80.0%
quantile manual check: 2.1275282556894264
5. quantile: 14.583785188055503 -> percentile: 100.0%
quantile manual check: 14.583785188055503

```

```

Percentile point function check for 0. quantile:
qt.transform(quantile):
[[-5.19933758]]
norm.ppf(percentile):
-5.1993375821928165

Percentile point function check for 1. quantile:
qt.transform(quantile):
[[-0.84162123]]
norm.ppf(percentile):
-0.8416212335729142

Percentile point function check for 2. quantile:
qt.transform(quantile):
[[-0.2533471]]
norm.ppf(percentile):
-0.2533471031357997

Percentile point function check for 3. quantile:
qt.transform(quantile):
[[0.2533471]]
norm.ppf(percentile):
0.2533471031357997

Percentile point function check for 4. quantile:
qt.transform(quantile):
[[0.84162123]]
norm.ppf(percentile):
0.8416212335729143

Percentile point function check for 5. quantile:
qt.transform(quantile):
[[5.19933758]]
norm.ppf(percentile):
5.199337582290661

x=2.5334162010006205 -> percentile(x)=0.806517013056411

qt.transform(x):
[[0.86513267]]

```

```

x_transformed:
0.8651326722548073

qt.inverse_transform(qt.transform(x)):
[[2.5334162]]
x_back:
2.533416201000622

```

The next presented script aims to try to use this technique in a non-conventional way to fit the original data, to be precise, a weighted histogram. Since the method only works with non-weighted data, sampling has to be done first. Once the sampling is done, the quantile transformation can be applied, and the transformed data can be fitted. After fitting the transformed data with a known distribution, one can then sample from it and transfer these newly created samples back to the original space. The script uses the module `data_inspect.py` to read the data and create a histogram that represents the radius probability distribution of a jet p_T bin.

`data_inspect.py:`

```

1 import h5py
2 import numpy as np
3
4 # **READING FROM A FILE**
5
6 def load_cnsts_rodem(ifile: str, mask_pt: list):
7     """Load constituents from an HDF5 file."""
8     with h5py.File(ifile, "r") as f:
9
10         cnsts=f["objects/jets/jet1_cnsts"][mask_pt,:,:3] #first three features
11
12         mask = np.repeat(cnsts[:, :, 0] == 0, cnsts.shape[2])
13         mask = mask.reshape(-1, cnsts.shape[1], cnsts.shape[2])
14         cnsts = np.ma.masked_where(mask, cnsts)
15
16         print(f"cnsts.shape={cnsts.shape}")
17         #(num_jets,num_particles,num_particle_features)
18         print(f"\n")
19
20         return cnsts[:, :, 0], cnsts[:, :, 1], cnsts[:, :, 2] #pt,eta,phi
21
22 def load_jets_rodem(ifile: str, mask_pt: list):
23     with h5py.File(ifile, "r") as f:
24
25         jets = f["objects/jets/jet1_obs"][mask_pt,:,:3] #first three features
26
27         print(f"jets.shape={jets.shape}")

```

```

28     #(num_jets,num_jet_features)
29     print(f"\n")
30
31     return jets[:,0],jets[:,1],jets[:,2] #pt,eta,phi
32
33     #-----
34
35     # **ARRANGE DATA IN A DICTIONARY:**
36
37     def read_data(input,mask_pt):
38
39         jets={}
40
41         #rows -> jets ; columns -> particles
42         jets['part_pt'],jets['part_eta'],jets['part_phi']=load_cnsts_rodem(input,
43                                     mask_pt)
44
45         #flat arrays
46         jets['jet_pt'],jets['jet_eta'],jets['jet_phi']=load_jets_rodem(input,mask_pt
47                                     )
48
49         jets['jet_nparticles']=jets['part_pt'].count(axis=1)
50
51         print(f"The smallest jet_pt is: {np.min(jets['jet_pt'])}")
52         print(f"The biggest jet_pt is: {np.max(jets['jet_pt'])}")
53         print(f"\n")
54
55         return jets
56
57     #-----
58
59     # **CREATE A HISTOGRAM AND CALCULATE ITS COVARIANCE MATRIX:**
60
61     def hist_cov(R,bins_r,weights):
62
63         # **MAKING A HISTOGRAM:**
64
65         n_jets=len(R)
66
67         R_middles = (bins_r[1:]+bins_r[:-1])/2.
68         num_R_bins=len(R_middles)
69
70         weights=np.ma.array(weights/(n_jets*np.sum(weights,axis=1).reshape(-1,1)),
71                               dtype=np.float64)
72
73         W=np.sum(weights,axis=1)
74
75         hist=np.zeros([num_R_bins],dtype=np.float64)

```

```

74 hist, _ = np.histogram(R.compressed(),
75                          bins=bins_r, weights=weights.compressed())
76
77 #-----
78
79 # **COMPUTE THE COVARIANCE MATRIX**
80
81 H=np.zeros([n_jets,num_R_bins],dtype=np.float64)
82 H2=np.zeros([n_jets,num_R_bins],dtype=np.float64)
83
84 for J in range(n_jets):
85
86     H[J,:], _ = np.histogram(R[J].compressed(),
87                              bins=bins_r, weights=weights[J].compressed())
88     H2[J,:], _ = np.histogram(R[J].compressed(),
89                               bins=bins_r, weights=weights[J].compressed()**2)
90
91 a=-H[:, :, np.newaxis]*np.ones(num_R_bins)
92
93 mask=np.eye(num_R_bins, dtype=bool)
94
95 a[:, mask] += np.repeat(W, num_R_bins).reshape(a.shape[0], a.shape[1])
96
97 a=a/(n_jets*(W ** 2)[ :, np.newaxis, np.newaxis])
98
99 b=a
100
101 cov=np.zeros([num_R_bins,num_R_bins],dtype=np.float64)
102 cov=np.einsum('mik,mjk,mk->ij', a, b, H2)
103
104 """
105 for i in range(num_R_bins):
106
107     for j in range(num_R_bins):
108
109         for J in range(n_jets):
110
111             for k in range(num_R_bins):
112
113                 if k==i:
114                     a=(W[J]-H[J,i])/(n_jets*W[J]**2)
115                 else:
116                     a=-H[J,i]/(n_jets*W[J]**2)
117
118                 if k==j:
119                     b=(W[J]-H[J,j])/(n_jets*W[J]**2)
120                 else:
121                     b=-H[J,j]/(n_jets*W[J]**2)

```

```

122
123         cov[i,j]+=a*b*H2[J,k]
124     """
125
126     #-----
127
128     return hist,cov

```

quantiles3.ipynb:

```

1  import sys
2  sys.path.append("./")
3
4  import os
5  import h5py
6  import numpy as np
7  import matplotlib.pyplot as plt
8  from data_inspect import *
9  from scipy.stats import norm
10 from sklearn.preprocessing import QuantileTransformer
11
12 # **SETTING THE DATA STRUCTURE:**
13
14 # The input HDF5 file containing the relevant jets:
15 input = "QCDjj_pT_450_1200_train01.h5" #change if needed
16
17 # The number of jets to load:
18 n_jets = 100_000 #change if needed
19
20 # Jet pt range:
21 min_pt=450 #change if needed
22 max_pt=1000 #change if needed
23 width_pt=10
24
25 # Choosing only the jets from the set jet pt range
26 with h5py.File(input, "r") as f:
27
28     mask_pt=[index for index,value in enumerate(
29         f["objects/jets/jet1_obs"][:n_jets,0])
30         if ((value>=min_pt) and (value<=max_pt))]
31
32 R_max=1.0
33 width=0.01
34
35 jets=read_data(input,mask_pt)
36
37 #-----
38

```

```

39 # **R DISTRIBUTION OF JETS:**
40
41 #phi angle differences between the constituents and the associated jet:
42 #rows -> jets ; columns -> particles
43 delta_phi=jets['jet_phi'].reshape(-1,1)-jets['part_phi']
44
45 #because the phi angle difference should be inside [-np.pi,+np.pi]:
46 delta_phi_new=np.ma.mod(delta_phi+np.pi,2*np.pi)-np.pi
47
48 #eta angle differences between the constituents and the associated jet:
49 #rows -> jets ; columns -> particles
50 delta_eta=jets['jet_eta'].reshape(-1,1)-jets['part_eta']
51
52 #rows -> jets ; columns -> particles
53 R=np.ma.sqrt(delta_eta**2+delta_phi_new**2)
54
55 print(f"number_of_particles_loaded:_\
56 {len(R.compressed())}")
57
58 print(f"number_of_particles_with_R=0:_\
59 {len(np.where(R.compressed()==0.0)[0])}")
60
61 print(f"percentage_of_particles_with_R=0:_\
62 {np.round(len(np.where(R.compressed()==0.0)[0])*100/len(R.compressed()),3)}%")
63
64 print(f"number_of_particles_with_R>={R_max}:_\
65 {len(np.where(R.compressed()>=R_max)[0])}")
66
67 print(f"percentage_of_particles_with_R>={R_max}:_\
68 {np.round(len(np.where(R.compressed()>=R_max)[0])*100/len(R.compressed()),3)}%")
69
70 print(f"\n")
71
72 jets['part_pt']=np.ma.masked_where(R<0,jets['part_pt'])
73 jets['part_pt']=np.ma.masked_where(R>=R_max,jets['part_pt'])
74
75 R=np.ma.masked_where(R<0,R)
76 R=np.ma.masked_where(R>=R_max,R)
77
78 #-----
79
80 # **SETTING pT BINS:**
81
82 bins_pt=np.arange(min_pt,max_pt+width_pt,width_pt)
83 pt_middles = (bins_pt[1:]+bins_pt[:-1])/2.
84 num_pt_bins=len(pt_middles)
85
86 print(f"num_pt_bins={num_pt_bins}")

```

```

87 print(f"\n")
88
89 #-----
90
91 # **SETTING R BINS:**
92
93 bins_r=np.arange(0,R_max,width)
94 R_middles = (bins_r[1:]+bins_r[:-1])/2.
95 num_R_bins=len(R_middles)
96
97 print(f"num_R_bins={num_R_bins}")
98 print(f"\n",flush=True)
99
100 #-----
101
102 mask_binpt=[]
103
104 for k in range(num_pt_bins):
105
106     mask_binpt.append([index for index,value in enumerate(
107         jets['jet_pt']) if ((value>bins_pt[k])
108             and (value<=bins_pt[k+1]))])
109
110 num_jets=[len(sublist) for sublist in mask_binpt]
111
112 #-----
113
114
115 # **CREATING THE HISTOGRAM AND CALCULATING ITS COVARIANCE MATRIX**
116
117 idx=2
118
119 weights=jets['part_pt'][mask_binpt[idx]]
120
121 hist,cov_matrix=hist_cov(R[mask_binpt[idx]],bins_r,weights)
122
123 #-----
124
125 # **SAMPLING:**
126
127 weights=weights/(num_jets[idx]*np.sum(weights,axis=1).reshape(-1,1))
128
129 n_samples = 500000
130 data = np.random.choice(R[mask_binpt[idx],:].compressed(),
131     size=n_samples,p=weights.compressed())
132
133 #-----
134

```



```

135 # **SAMPLING VISUALIZED:**
136
137 if not os.path.exists("./Plots"):
138     os.makedirs("./Plots")
139
140 fig,ax=plt.subplots(3,1,figsize=(6,15))
141
142 ax[0].set_title(f"{num_jets[idx]}_jets_with_pT_in_range_"
143                f"{bins_pt[idx]}-{bins_pt[idx+1]}_GeV")
144 ax[0].bar(R_middles,hist,width=width,ec='black')
145 ax[0].set_xlabel("R_original")
146 ax[0].set_ylabel("original_weights")
147
148 ax[1].set_title(f"{num_jets[idx]}_jets_with_pT_in_range_"
149                f"{bins_pt[idx]}-{bins_pt[idx+1]}_GeV")
150 ax[1].hist(data,bins=bins_r,ec='black')
151 ax[1].set_xlabel("R_samples")
152 ax[1].set_ylabel("counts")
153
154 a=1/n_samples*np.ones(shape=n_samples)
155 ax[2].set_title(f"{num_jets[idx]}_jets_with_pT_in_range_"
156                f"{bins_pt[idx]}-{bins_pt[idx+1]}_GeV")
157 ax[2].hist(data,bins=bins_r,weights=a,ec='black')
158 ax[2].set_xlabel("R_samples")
159 ax[2].set_ylabel("weights_counts/n_samples")
160
161 fig.savefig("./Plots/sampling.jpg")
162
163 #-----
164
165 # **TRANSFORMING SAMPLED R, PLOTTING IT AND FITTING**
166
167 # Choose either "uniform" or "normal":
168 choice="normal"
169
170 # Setting the transformer
171 n_quantiles=100
172 qt = QuantileTransformer(
173     output_distribution=choice,subsample=None,n_quantiles=n_quantiles)
174
175 # Transforming:
176 data_trans = qt.fit_transform(X=data.reshape(-1,1))
177
178 fig,ax=plt.subplots(1,2,figsize=(12,5))
179
180 ax[0].set_title(f"{num_jets[idx]}_jets_with_pT_in_range_"
181                f"{bins_pt[idx]}-{bins_pt[idx+1]}_GeV")
182 ax[0].hist(data_trans,bins=len(bins_r),ec='black')

```

```

183 ax[0].set_xlabel("R-transformedsamples")
184 ax[0].set_ylabel("counts")
185
186 if choice=="uniform":
187     # Fitting uniform:
188     x=np.linspace(min(data_trans.flatten()),max(data_trans.flatten()),10000)
189     pdf=np.ones(len(x))*1/(max(data_trans.flatten())-min(data_trans.flatten()))
190
191 elif choice=="normal":
192     # Fitting a Gaussian:
193     x=np.linspace(min(data_trans.flatten()),max(data_trans.flatten()),10000)
194     mu,sigma = norm.fit(data_trans)
195     pdf = norm.pdf(x,loc=mu,scale=sigma)
196
197 ax[1].set_title(f"{num_jets[idx]}jetswithTinrange"
198               f"{bins_pt[idx]}-{bins_pt[idx+1]}GeV")
199 ax[1].plot(x, pdf, 'r-')
200 ax[1].hist(data_trans,bins=len(bins_r),ec='black',density=True)
201 ax[1].set_xlabel("R-transformedsamples")
202 ax[1].set_ylabel("weights=counts/[sum(counts)*width]")
203
204 fig.savefig(f"./Plots/trans_{choice}.jpg")
205
206 # -----
207
208 # **SAMPLING FROM THE FITTED DISTRIBUTION:**
209
210 prob=pdf/np.sum(pdf)
211 samples = np.random.choice(x,size=n_samples,p=prob)
212
213 plt.figure(figsize=(6,5))
214 plt.hist(samples,bins=len(bins_r),density=True,ec='black')
215 plt.plot(x,pdf,'r-')
216 plt.xlabel(f"samplesfromthe{choice}distributionfit")
217 plt.ylabel("weights=counts/[sum(counts)*width]")
218 plt.savefig(f"./Plots/samples_{choice}.jpg")
219
220 # -----
221
222 # **TRASFORMING THE FIT BACK TO THE ORIGINAL SPACE**
223
224 samples_back=qt.inverse_transform(samples.reshape(-1,1))
225
226 hist_fit=[]
227 middles_fit=[]
228
229 # Making the fit as smooth as possible
230 dif=0.001

```

```

231 for i in range(0,70):
232     h, be = np.histogram(samples_back.flatten(),
233                           bins=len(bins_r), weights=a, range=(0+dif*i, R_max+dif*i))
234     hist_fit.extend(h)
235     bm = (be[1:]+be[:-1])/2.
236     middles_fit.extend(bm)
237
238 sorting=np.argsort(np.argsort(middles_fit))
239 hist_fit = [val for _, val in sorted(zip(sorting,hist_fit))]
240 middles_fit=[val for _, val in sorted(zip(sorting,middles_fit))]
241
242 #-----
243
244 # **PLOT OF THE FIT IN THE ORIGINAL SPACE**
245
246 plt.figure(figsize=(6,5))
247 plt.title(f"{num_jets[idx]} jets with pT in range "
248           fr"{bins_pt[idx]}-{bins_pt[idx+1]} GeV $\rightarrow$ {choice}, "
249           f"n_quantiles={n_quantiles}")
250 plt.scatter(R_middles,hist,color='b',marker='x')
251 plt.plot(middles_fit,hist_fit,'r-')
252 plt.xlabel("R-original")
253 plt.ylabel("original weights")
254 plt.savefig(f"./Plots/fit_{choice}.jpg")

```

output:

```

cnsts.shape=(98758, 100, 3)

jets.shape=(98758, 3)

The smallest jet_pt is: 450.00129118847076
The biggest jet_pt is: 999.8620172940161

number of particles loaded: 5365439
number of particles with R=0: 1
percentage of particles with R=0: 0.0%
number of particles with R>=1.0: 14630
percentage of particles with R>=1.0: 0.273%

num_pt_bins=55

num_R_bins=99

```

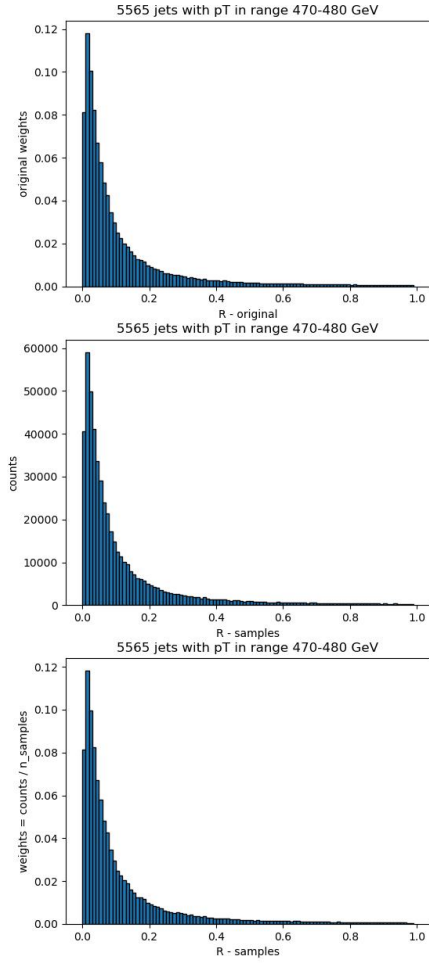


Figure 22: Sampling of the weighted data: sampling.jpg

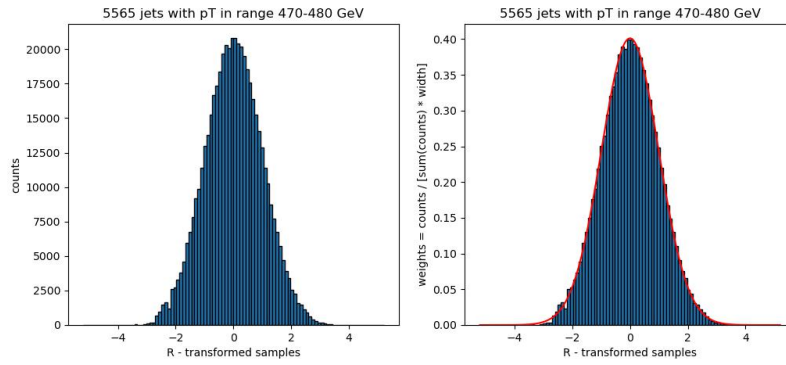


Figure 23: Fitting of transformed sampled data: trans_normal.jpg

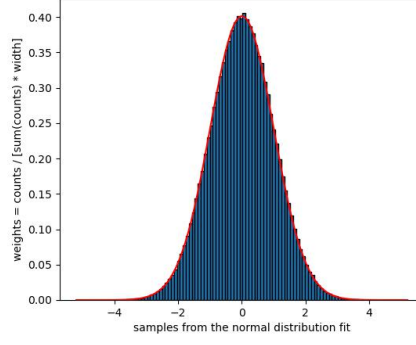


Figure 24: Sampling from the fitted distribution: samples_normal.jpg

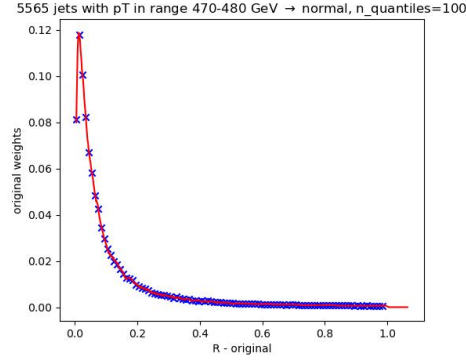


Figure 25: Fitted distribution transformed back to the original space: fit_normal.jpg

As one can see from the figures above, the fit of the known distribution transformed back to the original space is quite good for the chosen jet p_T bin. However, the problem is that this method is nonparametric, so one does not end up with some expression for the fit function.

9.3 Box-Cox transformation

In an attempt to find some parametric method of transforming the data, fitting it and then transforming the fit back to the original space, the Box-Cox transformation was found and explored. The transformation is given by the following expression.

$$y = \begin{cases} \frac{(x^\lambda - 1)}{\lambda}, & \lambda \neq 0 \\ \log(x), & \lambda = 0 \end{cases}$$

The algorithm finds the best parameter λ when the data is loaded. The goal of the transformation is to end up with the normal distribution of the data. Compared to the quantile transformation, this transformation is parametric. Moreover, unlike the quantile transformation, it preserves the data structure - there is no sorting of the data before the transformation is applied. However, the quantile transformation technically always works, while the Box-Cox transformation does not.

The script in this subsection was also run in the Jupyter Notebook on a local computer for simplicity. The first part of the script mirrors the procedure introduced in `quantiles3.ipynb`.

boxcox.ipynb:

```
1 import sys
2 sys.path.append("./")
3
4 import os
5 import h5py
6 import numpy as np
7 import matplotlib.pyplot as plt
8 from data_inspect import *
9 from inspect import signature
10 from scipy.stats import norm, boxcox
11 from scipy.special import inv_boxcox
12 from scipy.optimize import curve_fit
13
14 # **SETTING THE DATA STRUCTURE:**
15
16 # The input HDF5 file containing the relevant jets:
17 input = "QCDjj_pT_450_1200_train01.h5" #change if needed
18
19 # The number of jets to load:
20 n_jets = 100_000 #change if needed
21
22 # Jet pt range:
23 min_pt=450 #change if needed
24 max_pt=1000 #change if needed
25 width_pt=10
26
27 # Choosing only the jets from the set jet pt range
28 with h5py.File(input, "r") as f:
29
30     mask_pt=[index for index,value in enumerate(
31         f["objects/jets/jet1_obs"][:n_jets,0])
32             if ((value>=min_pt) and (value<=max_pt))]
33
34 R_max=1.0
35 width=0.01
36
37 jets=read_data(input,mask_pt)
38
39 #-----
40
41 # **R DISTRIBUTION OF JETS:**
42
43 #phi angle differences between the constituents and the associated jet:
44 #rows -> jets ; columns -> particles
45 delta_phi=jets['jet_phi'].reshape(-1,1)-jets['part_phi']
46
47 #because the phi angle difference should be inside [-np.pi,+np.pi]:
```

```

48 delta_phi_new=np.ma.mod(delta_phi+np.pi,2*np.pi)-np.pi
49
50 #eta angle differences between the constituents and the associated jet:
51 #rows -> jets ; columns -> particles
52 delta_eta=jets['jet_eta'].reshape(-1,1)-jets['part_eta']
53
54 #rows -> jets ; columns -> particles
55 R=np.ma.sqrt(delta_eta**2+delta_phi_new**2)
56
57 print(f"number_of_particles_loaded:\n
58 {len(R.compressed())}")
59
60 print(f"number_of_particles_with_R=0:\n
61 {len(np.where(R.compressed()==0.0)[0])}")
62
63 print(f"percentage_of_particles_with_R=0:\n
64 {np.round(len(np.where(R.compressed()==0.0)[0])*100/len(R.compressed()),3)}%")
65
66 print(f"number_of_particles_with_R>={R_max}:\n
67 {len(np.where(R.compressed()>=R_max)[0])}")
68
69 print(f"percentage_of_particles_with_R>={R_max}:\n
70 {np.round(len(np.where(R.compressed()>=R_max)[0])*100/len(R.compressed()),3)}%")
71
72 print(f"\n")
73
74 jets['part_pt']=np.ma.masked_where(R<0,jets['part_pt'])
75 jets['part_pt']=np.ma.masked_where(R>=R_max,jets['part_pt'])
76
77 R=np.ma.masked_where(R<0,R)
78 R=np.ma.masked_where(R>=R_max,R)
79
80 #-----
81
82 # **SETTING pT BINS:**
83
84 bins_pt=np.arange(min_pt,max_pt+width_pt,width_pt)
85 pt_middles = (bins_pt[1:]+bins_pt[:-1])/2.
86 num_pt_bins=len(pt_middles)
87
88 print(f"num_pt_bins={num_pt_bins}")
89 print(f"\n")
90
91 #-----
92
93 # **SETTING R BINS:**
94
95 bins_r=np.arange(0,R_max,width)

```

```

96 R_middles = (bins_r[1:]+bins_r[:-1])/2.
97 num_R_bins=len(R_middles)
98
99 print(f"num_R_bins={num_R_bins}")
100 print(f"\n",flush=True)
101
102 #-----
103
104 mask_binpt=[]
105
106 for k in range(num_pt_bins):
107
108     mask_binpt.append([index for index,value in enumerate(
109         jets['jet_pt']) if ((value>bins_pt[k])
110             and (value<=bins_pt[k+1]))])
111
112 num_jets=[len(sublist) for sublist in mask_binpt]
113
114 #-----
115
116 # **CREATING THE HISTOGRAMS AND CALCULATING THEIR COVARIANCE MATRICES:**
117
118 hist=np.zeros([num_pt_bins,num_R_bins])
119 cov_matrix=np.zeros([num_pt_bins,num_R_bins,num_R_bins])
120
121 for idx in range(num_pt_bins):
122
123     print(f"hist_cov_of_bin_{idx}/{num_pt_bins-1} in progress..")
124
125     weights=jets['part_pt'][mask_binpt[idx]]
126
127     hist[idx],cov_matrix[idx]=hist_cov(R[mask_binpt[idx]],bins_r,weights)
128
129 print(f"\n")
130
131 #-----
132
133 # **SAMPLING:**
134
135 idx=2
136
137 weights=jets['part_pt'][mask_binpt[idx]]
138 weights=weights/(num_jets[idx]*np.sum(weights,axis=1).reshape(-1,1))
139
140 n_samples = 500000
141 data = np.random.choice(R[mask_binpt[idx]].compressed(),
142     size=n_samples,p=weights.compressed())
143

```



```

144 #-----
145
146 # **TRANSFORMING SAMPLED R, PLOTTING IT AND FITTING:**
147
148 # Transforming:
149 data_trans,best_lambda = boxcox(data)
150
151 if not os.path.exists("./Plots"):
152     os.makedirs("./Plots")
153
154 fig,ax=plt.subplots(1,2,figsize=(12,5))
155
156 ax[0].set_title(f"{num_jets[idx]}_jets_with_pT_in_range_"
157                f"{bins_pt[idx]}-{bins_pt[idx+1]}_GeV")
158 ax[0].hist(data_trans,bins=len(bins_r),ec='black')
159 ax[0].set_xlabel("R-transformed samples")
160 ax[0].set_ylabel("counts")
161
162 # Fitting a Gaussian:
163 x=np.linspace(min(data_trans.flatten()),max(data_trans.flatten()),10000)
164 mu,sigma = norm.fit(data_trans)
165 pdf = norm.pdf(x,loc=mu,scale=sigma)
166
167 ax[1].set_title(f"{num_jets[idx]}_jets_with_pT_in_range_"
168                f"{bins_pt[idx]}-{bins_pt[idx+1]}_GeV")
169 ax[1].plot(x,pdf,'r-')
170 ax[1].hist(data_trans,bins=len(bins_r),ec='black',density=True)
171 ax[1].set_xlabel("R-transformed samples")
172 ax[1].set_ylabel("weights=counts/[sum(counts)*width]")
173
174 fig.savefig("./Plots/trans_boxcox.jpg")
175
176 #-----
177
178 # **SAMPLING FROM THE GAUSSIAN:**
179
180 prob=pdf/np.sum(pdf)
181 samples = np.random.choice(x,size=n_samples,p=prob)
182
183 plt.figure(figsize=(6,5))
184 plt.hist(samples,bins=len(bins_r),density=True,ec='black')
185 label="Fitted Gaussian"
186 label+=f"\n"
187 label+=fr"$\mu$={mu:.2f}, $\sigma$={sigma:.2f}"
188 plt.plot(x,pdf,'r-',label=label)
189 plt.legend()
190 plt.savefig(f"./Plots/samples_boxcox.jpg")
191

```

```

192 #-----
193
194 # **TRANFORMING THE FIT BACK TO THE ORIGINAL SPACE**
195
196 samples_back = inv_boxcox(samples,best_lambda)
197
198 hist_fit=[]
199 middles_fit=[]
200
201 # Making the fit as smooth as possible
202 dif=0.001
203 a=1/n_samples*np.ones(shape=n_samples)
204 for i in range(0,70):
205     h, be = np.histogram(samples_back,
206                           bins=len(bins_r),weights=a,range=(0+dif*i,R_max+dif*i))
207     hist_fit.extend(h)
208     bm = (be[1:]+be[:-1])/2.
209     middles_fit.extend(bm)
210
211 sorting=np.argsort(np.argsort(middles_fit))
212 hist_fit = [val for _, val in sorted(zip(sorting,hist_fit))]
213 middles_fit=[val for _, val in sorted(zip(sorting,middles_fit))]
214
215 #-----
216
217 # **PLOT OF THE FIT IN THE ORIGINAL SPACE**
218
219 plt.figure(figsize=(6,5))
220 plt.title(f"{num_jets[idx]}_jets_with_{T_in_range}"
221           f"{bins_pt[idx]}-{bins_pt[idx+1]}_GeV")
222 plt.scatter(R_middles,hist[idx],color='b',marker='x')
223 plt.plot(middles_fit,hist_fit,'r-')
224 plt.xlabel("R-original")
225 plt.ylabel("original_weights")
226 plt.savefig("./Plots/fit_boxcox.jpg")
227
228 #-----
229
230 # **FITTING WITH THE NEWLY DISCOVERED FUNCTION AND PREPARING FOR THE PLOTS**
231
232 def model(x,lmbda,loc,scale):
233     return norm.pdf(boxcox(x,lmbda=lmbda),loc=loc,scale=scale)
234
235 rs=np.linspace(min(R_middles),max(R_middles),1000)
236 p0=[0.05,-2.5,1.2]
237
238 num_params=len(signature(model).parameters)-1
239

```

```

240 params=np.zeros([num_pt_bins,num_params])
241 params_err=np.zeros([num_pt_bins,num_params])
242 y_mean=np.zeros([num_pt_bins,len(rs)])
243 y_std=np.zeros([num_pt_bins,len(rs)])
244
245 for idx in range(num_pt_bins):
246
247     print(f"fitting of bin {idx}/{num_pt_bins-1} in progress..")
248
249     params[idx],cov = curve_fit(
250         model,xdata=R_middles,ydata=hist[idx,:],
251         sigma=cov_matrix[idx],p0=p0)
252
253     params_err[idx] = np.sqrt(np.diag(cov))
254
255     n_samples=1000
256     params_samples=np.random.multivariate_normal(
257         params[idx,:],cov,size=n_samples)
258     y_samples=np.array([model(rs,*sample) for sample in params_samples])
259     y_mean[idx,:]=np.mean(y_samples,axis=0)
260     y_std[idx,:]=np.std(y_samples,axis=0)
261
262 print(f"\n")
263
264 #-----
265
266 # **PLOT THE FIT OF SOME pT BIN:**
267
268 def digit(x):
269
270     str="{:.10e}".format(x)
271
272     base,exp=str.split('e')
273     exp=int(exp)
274
275     return -exp
276
277 if not os.path.exists("./Plots"):
278     os.makedirs("./Plots")
279
280 idx=15
281
282 fig,ax=plt.subplots(figsize=(10,6))
283
284 label=r"$\frac{1}{\sqrt{2\pi}\sigma^2}$"
285 label+="_"
286 label+=r"$e^{\frac{(\frac{R^{\lambda}-1}{\lambda}-\mu)^2}{2\sigma^2}}$"
287 label+=f"\n"

```

```

288 err_rounded=round(params_err[idx,0],digit(params_err[idx,0]))
289 mean_rounded=round(params[idx,0],digit(err_rounded))
290 label+=fr"$\lambda$={mean_rounded}$\pm${err_rounded};"
291 label+=f"\n"
292 err_rounded=round(params_err[idx,1],digit(params_err[idx,1]))
293 mean_rounded=round(params[idx,1],digit(err_rounded))
294 label+=fr"$\mu$={mean_rounded}$\pm${err_rounded};"
295 label+=f"\n"
296 err_rounded=round(params_err[idx,2],digit(params_err[idx,2]))
297 mean_rounded=round(params[idx,2],digit(err_rounded))
298 label+=fr"$\sigma$={mean_rounded}$\pm${err_rounded};"
299
300 ax.set_title(f"{num_jets[idx]}_jets_with_PT_in_range_
301             f"{bins_pt[idx]}-{bins_pt[idx+1]}_GeV")
302 ax.errorbar(R_middles,hist[idx],np.sqrt(np.diag(cov_matrix[idx])),
303             ls='none',capsize=2)
304 ax.scatter(R_middles,hist[idx],color='b',s=4)
305 ax.plot(rs,y_mean[idx],color='red',label=label)
306 ax.fill_between(rs,y_mean[idx]-y_std[idx],y_mean[idx]+y_std[idx],
307                color='red',alpha=0.3)
308 ax.set_xlabel("R")
309 ax.set_ylabel("P_{}_probability")
310 ax.legend()
311 fig.savefig("./Plots/norm.pdf(boxcox).jpg")
312
313 #-----
314
315 # **PLOTING AND FITTING OF THE PARAMETERS:**
316
317 def lin(x,a,b):
318     x=x-np.min(pt_middles)
319     x=x/(np.max(pt_middles)-np.min(pt_middles))
320     return a+b*x
321
322 pts=np.linspace(min(pt_middles),max(pt_middles),1000)
323
324 p=[r"$\lambda$",r"$\mu$",r"$\sigma$"]
325
326 fig,ax=plt.subplots(3,1,figsize=(6,15))
327
328 for n,s in enumerate(p):
329
330     pars,cov=curve_fit(
331         lin,xdata=pt_middles,ydata=params[:,n],sigma=params_err[:,n])
332
333     pars_err=np.sqrt(np.diag(cov))
334
335     n_samples=1000

```

```

336 pars_samples=np.random.multivariate_normal(pars,cov,size=n_samples)
337 param_samples=np.array([lin(pts,*sample) for sample in pars_samples])
338 param_mean=np.mean(param_samples,axis=0)
339 param_std=np.std(param_samples,axis=0)
340
341 label=r"$f(p_T)=a+b\cdot$"
342 label+=r"$\frac{p_T-\{\min(pt\_middles)\}}{\{\max(pt\_middles)-\min(pt\_middles)\}}$;"
343
344 label+=f"\n"
345 pars_err[0]=round(pars_err[0],digit(pars_err[0]))
346 pars[0]=round(pars[0],digit(pars_err[0]))
347 label+=f"$a=\{pars[0]\}\pm\{pars\_err[0]\}$;"
348 label+=f"\n"
349 pars_err[1]=round(pars_err[1],digit(pars_err[1]))
350 pars[1]=round(pars[1],digit(pars_err[1]))
351 label+=f"$b=\{pars[1]\}\pm\{pars\_err[1]\}$"
352
353 ax[n].errorbar(pt_middles,params[:,n],params_err[:,n],
354               ls='none',capsize=2)
355 ax[n].scatter(pt_middles,params[:,n],s=4)
356 ax[n].plot(pts,param_mean,color='red',label=label)
357 ax[n].fill_between(pts,param_mean-param_std,param_mean+param_std,
358                  color='red',alpha=0.3)
359 ax[n].set_xlabel(r"$jet\_p\_T$")
360 ax[n].set_ylabel(f"$s$")
361 ax[n].legend()
362
363 fig.savefig(f"./Plots/norm.pdf(boxcox)_params.jpg")

```

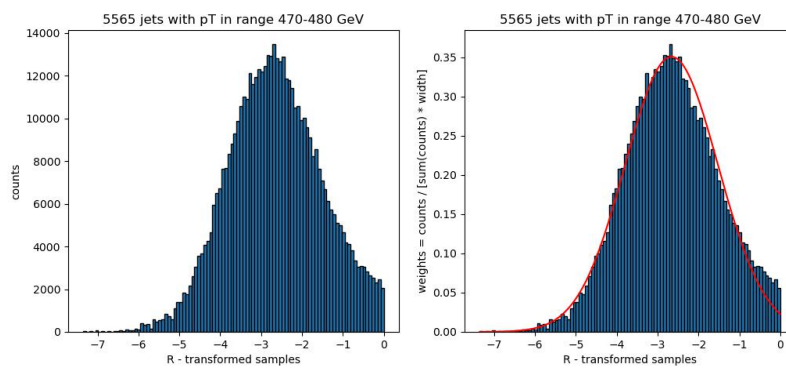


Figure 26: Fitting of transformed sampled data: trans_boxcox.jpg

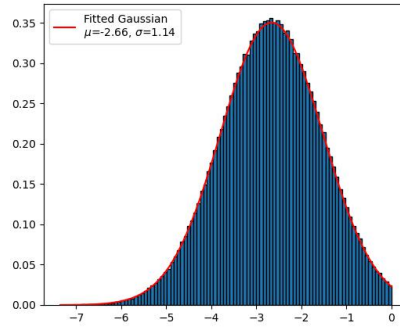


Figure 27: Sampling from the fitted distribution: samples_boxcox.jpg

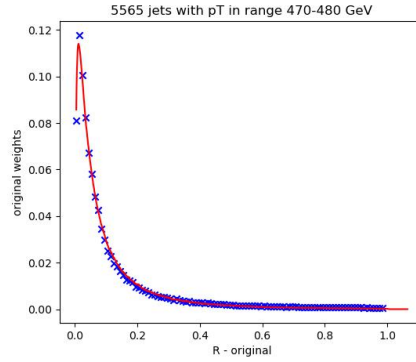


Figure 28: Fitted distribution transformed back to the original space: fit_boxcox.jpg

Still, even with the parametric transformation of the data samples to the space where the sample points are normally distributed, one key problem remains. Both the quantile transformation and the Box-Cox transformation work with the samples rather than with the actual data. The problem becomes the most apparent when one has to sample from the fitted distribution to then transfer these new samples back to the original space; the fit itself is not directly inversely transformed, but the samples of it. This does not allow one to find the expression of the final fit in the original space, even though the transformation is parametric. While analyzing the Box-Cox transformation, one function form that follows the data shape was discovered by accident and explored in the second part of the script.

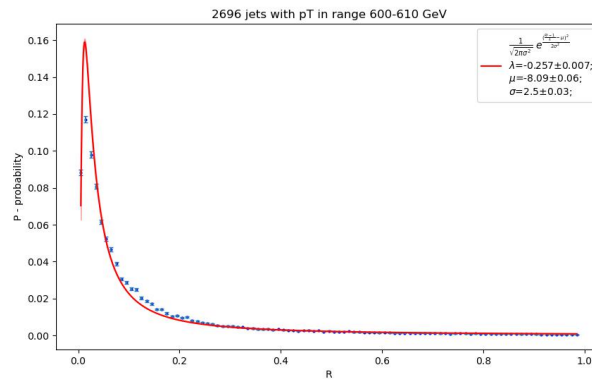


Figure 29: Probability fitted with the newly discovered function: norm.pdf(boxcox).jpg

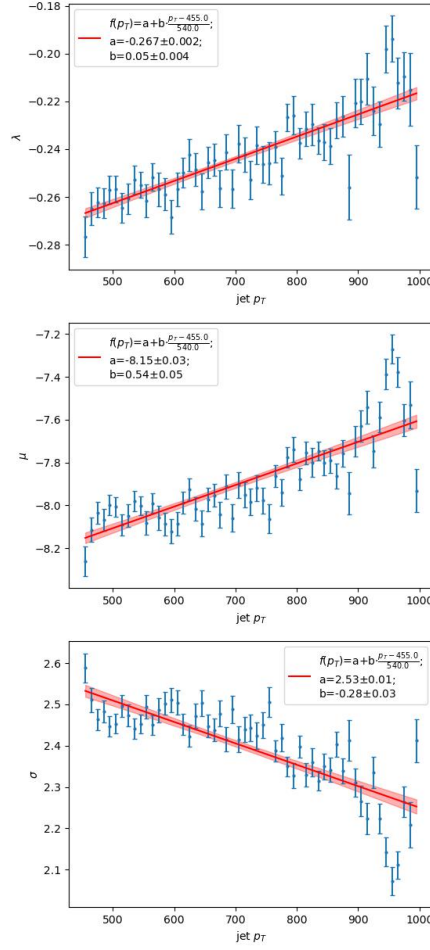


Figure 30: The parameters fitted across jet p_T bins: norm.pdf(boxcox).params.jpg

As can be seen from the figure above, the newly found function gave some hope but it tends to overshoot the peak.

9.4 Kernel density estimation

Kernel density estimation is a method for estimating the probability density function of a continuous variable by smoothing the data using a kernel function. It works only with the data samples, so again, sampling has to be done first to apply the method. The procedure places a smooth kernel function on each sample point and sums them up. The number of kernels then equals the number of samples. Because of this, kernel density estimation is not a fitting method in the traditional sense - it does not assume an underlying parametric model or learn parameters from data, but instead builds the estimate directly from the sample distribution, making it a purely data-driven, nonparametric approach. Although not a traditional fitting method, the underlying procedure was still explored in this subsection and in the following script.

kde.ipynb:

```
1 import os
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.neighbors import KernelDensity
5
6 # **THE SETUP:**
7
8 # The kernel function
9 def K(x):
10     return np.exp(-x**2/2)/np.sqrt(2*np.pi)
11
12 # Dummy dataset
13 dataset = np.array([1.33,1.33,1.33,0.3,0.97,0.97,0.96,1.1,0.1,1.4,0.4])
14
15 # x-value range for plotting KDEs
16 x_range = np.linspace(dataset.min()-0.3,dataset.max()+0.3,num=600)
17
18 # Bandwith values for experimentation
19 H = [0.3,0.1,0.03]
20 n_samples = dataset.size
21
22 # Line properties for different bandwith values
23 color_list = ['green','black','maroon']
24 alpha_list = [0.8,1,0.8]
25 width_list = [1.7,2.5,1.7]
26
27 #-----
28
29 # **CALCULATION AND PLOTTING:**
30
31 if not os.path.exists("./Plots"):
32     os.makedirs("./Plots")
33
34 plt.figure(figsize=(10,4))
35
36 # Iterate over bandwith values
37 for h, color, alpha, width in zip(H,color_list,alpha_list,width_list):
38
39     # **METHOD DONE MANUALLY:**
40
41     total_sum = 0
42     # Iterate over datapoints
43     for i, xi in enumerate(dataset):
44         total_sum += K((x_range-xi)/h)
45
46     # Normalization
47     y_range = total_sum/(h*n_samples)
```



```

48
49 # Plotting
50 plt.plot(x_range, y_range,
51          color=color, alpha=alpha, linewidth=width, label=f"{h}")
52
53 #-----
54
55 # **BUILT-IN METHOD**
56
57 kde = KernelDensity(
58     kernel='gaussian', bandwidth=h).fit(dataset[:, np.newaxis])
59
60 log_density = kde.score_samples(x_range[:, np.newaxis])
61
62 # Plotting
63 plt.plot(x_range, np.exp(log_density),
64          color='grey', alpha=0.8, linewidth=width, linestyle='--', label=f"{h}")
65
66 #-----
67
68 plt.hist(dataset, bins=50)
69 plt.xlabel("$x$", fontsize=22)
70 plt.ylabel("$f(x)$", fontsize=22, rotation="horizontal", labelpad=20)
71 plt.legend(fontsize=10, shadow=True, title="$h$", title_fontsize=16)
72 plt.savefig("./Plots/kde.jpg")

```

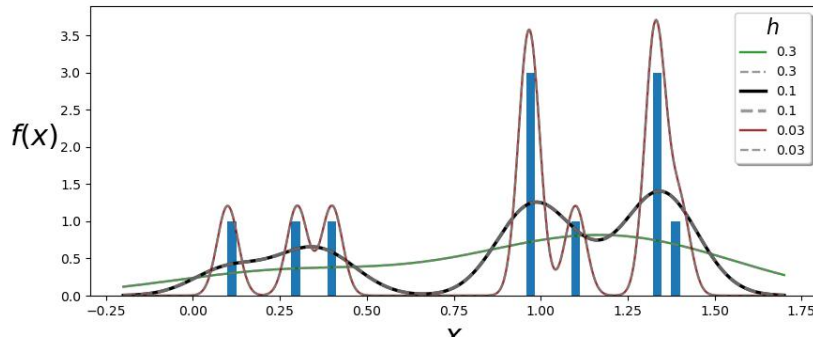


Figure 31: Showcase of Kernel Density Estimation method: kde.jpg

10 Results

10.1 Moments

The initial objective was to examine how the radius distributions behave across jet p_T bins for different data sets, using moment analysis. The procedure is as follows: A jet p_T bin is selected, and then the moment is calculated for each jet individually. A particle's contribution to the jet's moment is defined as its radius raised to the relevant power and multiplied by its contribution to the probability. As described in Section 8, this contribution is given by the transverse momentum of the particle divided

by the total transverse momentum of all particles in the same jet, further normalized by the number of jets considered, which in this case is one. The n-th moment of a jet is then calculated as follows.

$$M_J^{(n)} = \sum_{i(\text{particles})} R_{i(J)}^n \cdot \frac{p_{T i(J)}}{\sum_{k(\text{particles})} p_{T k(J)}} \quad (59)$$

Once the moment of each jet within a given jet p_T bin is calculated, the mean and standard error can be easily calculated. These summary statistics were then used to generate the plots shown below. **It is important to note that these are not the moments of the radius probability distributions corresponding to the jet p_T bins, but rather the average values of the jet-wise moments for all jets falling within a given p_T bin.**

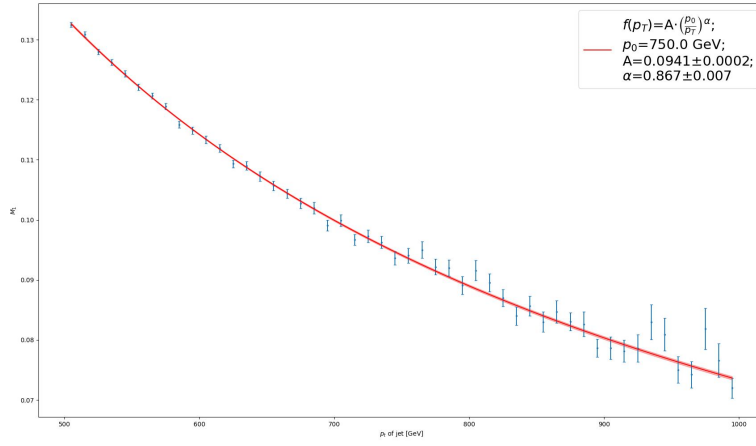


Figure 32: First moment behavior in WToQQ_070.root (~ 100000 jets considered): moments.root.jpg

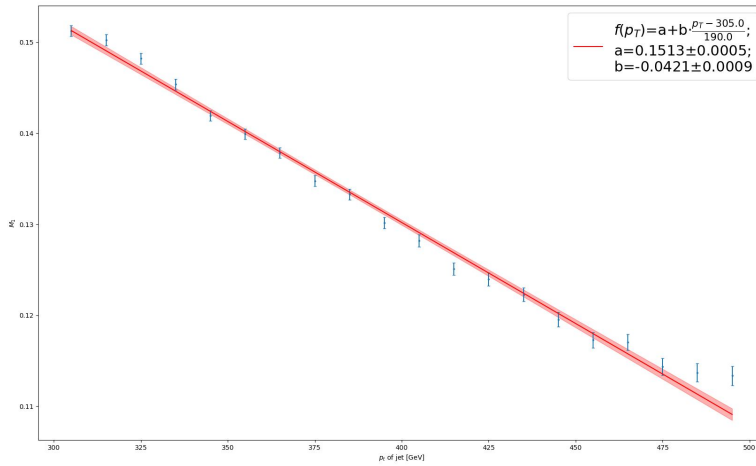


Figure 33: First moment behavior in RunG_batch0.h5 (~ 400000 jets considered): moments_lowpt.jpg

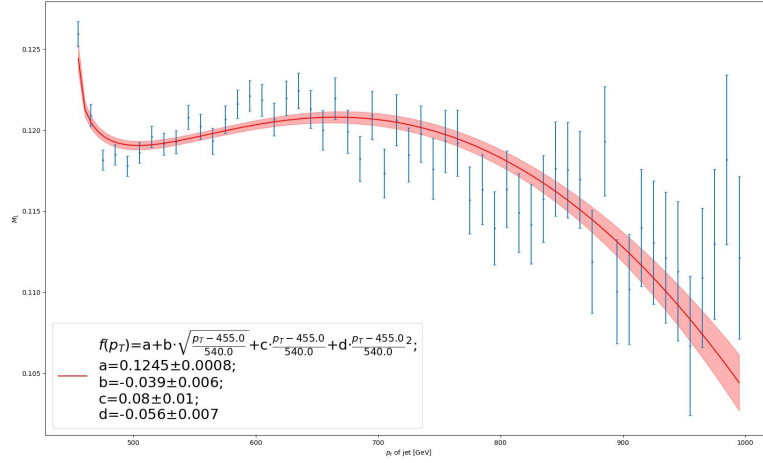


Figure 34: First moment behavior in QCDjj-pT_450-1200_train01.h5 (~ 400000 jets considered): moments_highpt.jpg

10.2 Radius probability distribution

The next step was to look at the radius probability distribution directly. For each data set, a single jet p_T bin was selected to illustrate the shape of the distribution.

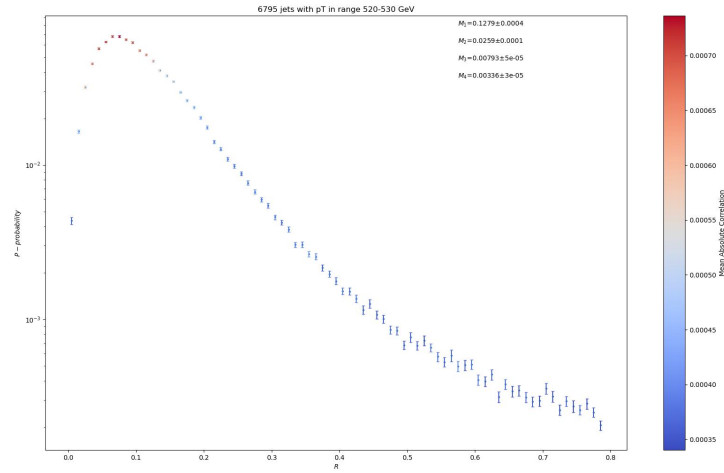


Figure 35: A jet p_T bin from WToQQ_070.root: probability_root.jpg

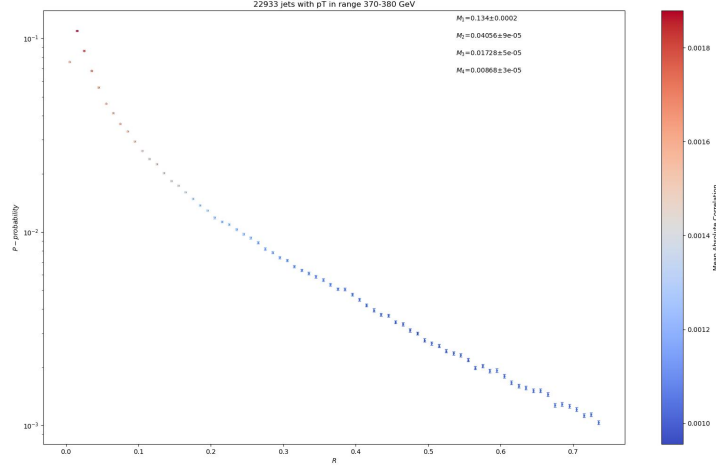


Figure 36: A jet p_T bin from RunG_batch0.h5: probability_lowpt.jpg

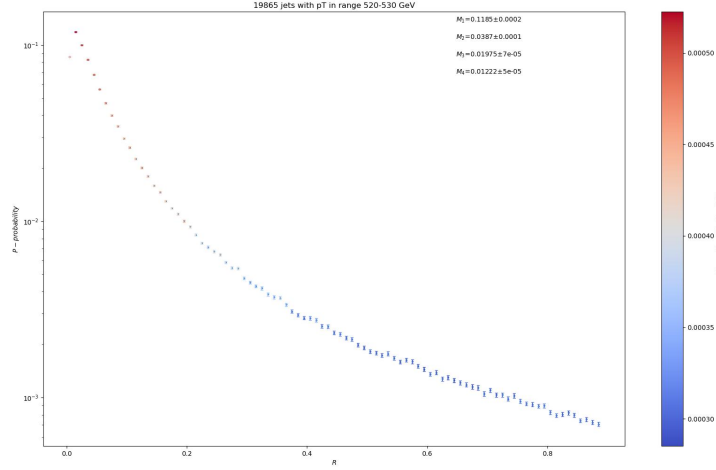


Figure 37: A jet p_T bin from QCDjj_pT_450_1200_train01.h5: probability_highpt.jpg

The moments of the given histograms are calculated as in the expression (60) below. Moment error is calculated with the help of the error propagation formula given in the expression (63) which, for this specific case, results in the expression (64). **It is important to note that these are the moments of the radius probability distribution corresponding to the jet p_T bin, not the average values of the jet-wise moments for all jets falling within the given p_T bin.**

$$\overline{M}_n = \sum_{k(bins)} (k\text{-th bin middle})^n \cdot \overline{P}(k\text{-th bin}) \quad (60)$$

$$x = F(a, b, c, \dots) \quad (61)$$

$$\bar{x} = F(\bar{a}, \bar{b}, \bar{c}, \dots) \quad (62)$$

$$x_{stdev}^2 = \sum_{k=a,b,c,\dots} \left(\frac{\partial F}{\partial k} \bigg|_{\substack{a=\bar{a} \\ b=\bar{b} \\ c=\bar{c}}} \cdot k_{stdev} \right)^2 \quad (63)$$

$$M_n^{(stdev)} = \sqrt{\sum_{k(bins)} (\text{k-th bin middle})^{2n} \cdot P^{(stdev)}(\text{k-th bin})^2} \quad (64)$$

$$M_n = \overline{M_n} \pm M_n^{(stdev)} \quad (65)$$

10.3 Radius probability density function

After many attempts to fit the radius probability distribution for each chosen jet p_T bin with a single functional form, it was realized that it would be easier to fit the radius probability density since the troubling peak at low radius then disappears. Once the peak was gone, the exponential polynomial function was utilized to obtain the final model. The chosen functional form for the parametric fit was a linear combination of Legendre polynomials inside an exponential function.

$$f(x) = C \cdot \exp(a_1 L_1(x) + a_2 L_2(x) + a_4 L_4(x) + a_5 L_5(x) + a_6 L_6(x) + a_7 L_7(x) + a_8 L_8(x)), \quad (66)$$

$$x = 2R - 1$$

The term involving $L_3(x)$ was not included, as this specific model configuration had initially been selected during an earlier phase when parts of the procedure were not yet fully correct. Despite that, the model has since proven effective even under the corrected conditions and was therefore retained in this simplified form. Importantly, the implementation remains flexible, allowing for straightforward testing of alternative configurations if needed.

The modeling was done in three steps described below.

Fitting Procedure

- **Step 1: Initial independent fits**

- Fit the radius probability density function independently for each jet p_T bin, using only the variances of the radius bins.
- Each fit yields a set of top-level parameters, specific to that jet p_T bin.
- Then, for each top-level parameter, fit a smooth function to model its dependence on jet p_T . The parameters of these functions are referred to as bottom-level parameters, and they remain uncorrelated and are fit separately for each top-level parameter.

- **Step 2: Refined independent fits with full covariance**

- Repeat the fits for each jet p_T bin, now using the full covariance matrix of the radius bins.
- The top-level parameters are determined again independently for each bin.
- The initial guesses for these parameters are obtained by evaluating the smooth trends from Step 1, that is, using bottom-level parameters from the previous step.
- As in Step 1, the bottom-level parameters are not jointly optimized; they serve only to guide the initial guesses.

• **Step 3: Global fit with shared bottom-level parameters**

- Perform a global fit of all jet p_T bins simultaneously by combining their covariance matrices into a single block-diagonal matrix.
- In this step, the top-level parameters for each p_T bin are no longer treated as independent values; instead, they are constrained to follow a common functional form governed by shared bottom-level parameters.
- These bottom-level parameters, which define how the top-level parameters vary with jet p_T , are now jointly optimized across all jet p_T bins and may become correlated through global fit.
- The initial guess for the global fit is constructed by concatenating the bottom-level parameters obtained in Step 2.

The following plots present the fits of the top-level parameters from Step 2 for each data set.

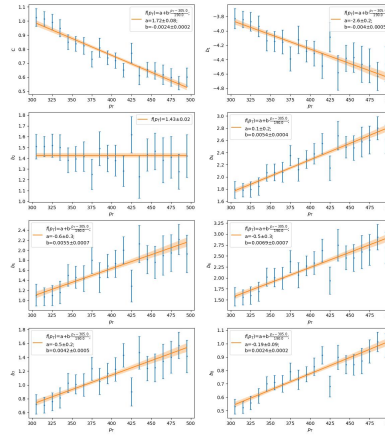


Figure 38: Parameter fits for RunG_batch0.h5: param_fits_lowpt.jpg

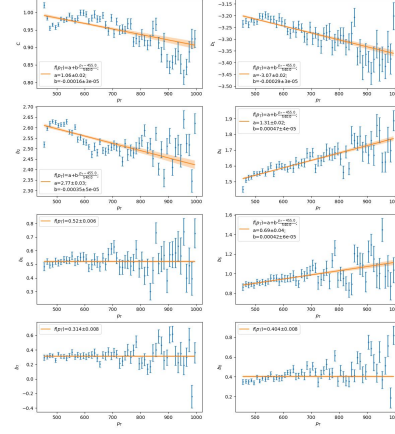


Figure 39: Parameter fits for QCDjj_pT.450.1200_train01.h5: param_fits_highpt.jpg

The bottom-level parameters presented in the legends of each data set are then used for the initial guess of the global fit.

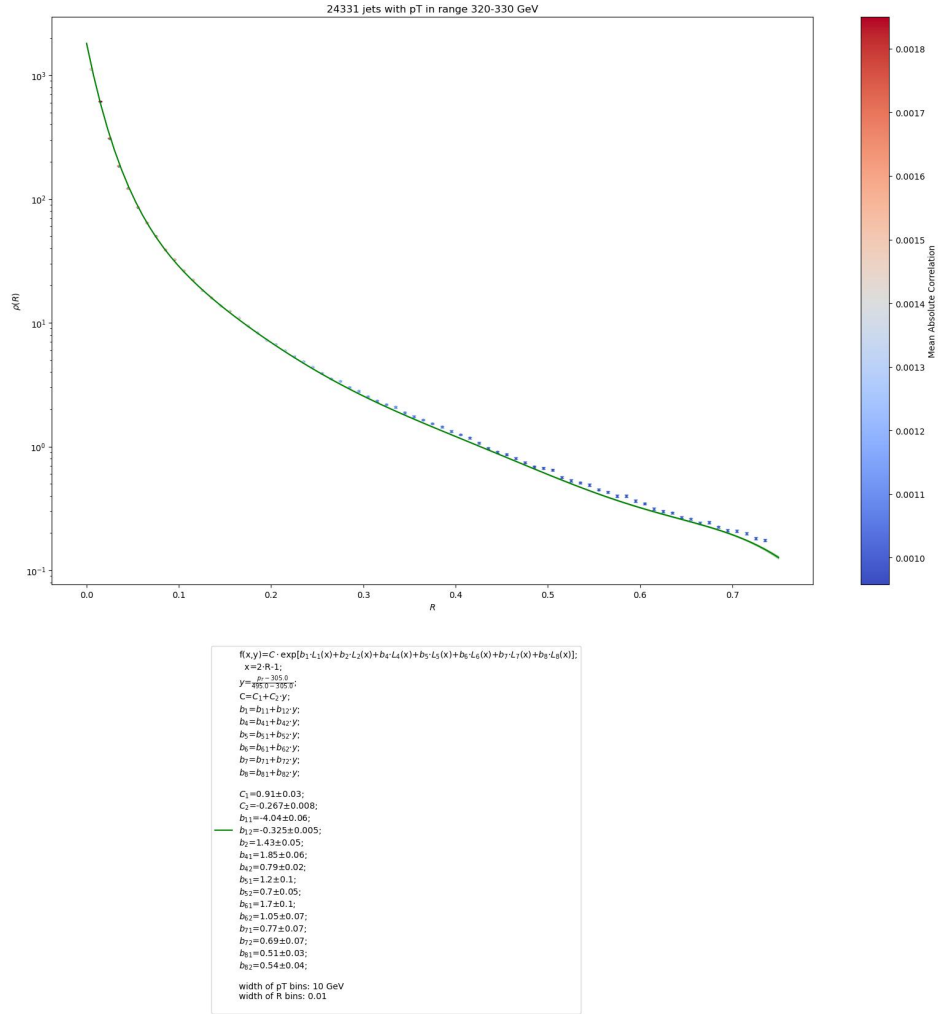


Figure 40: The final fit for RunG_batch0.h5: FINAL_FIT_lowpt.jpg

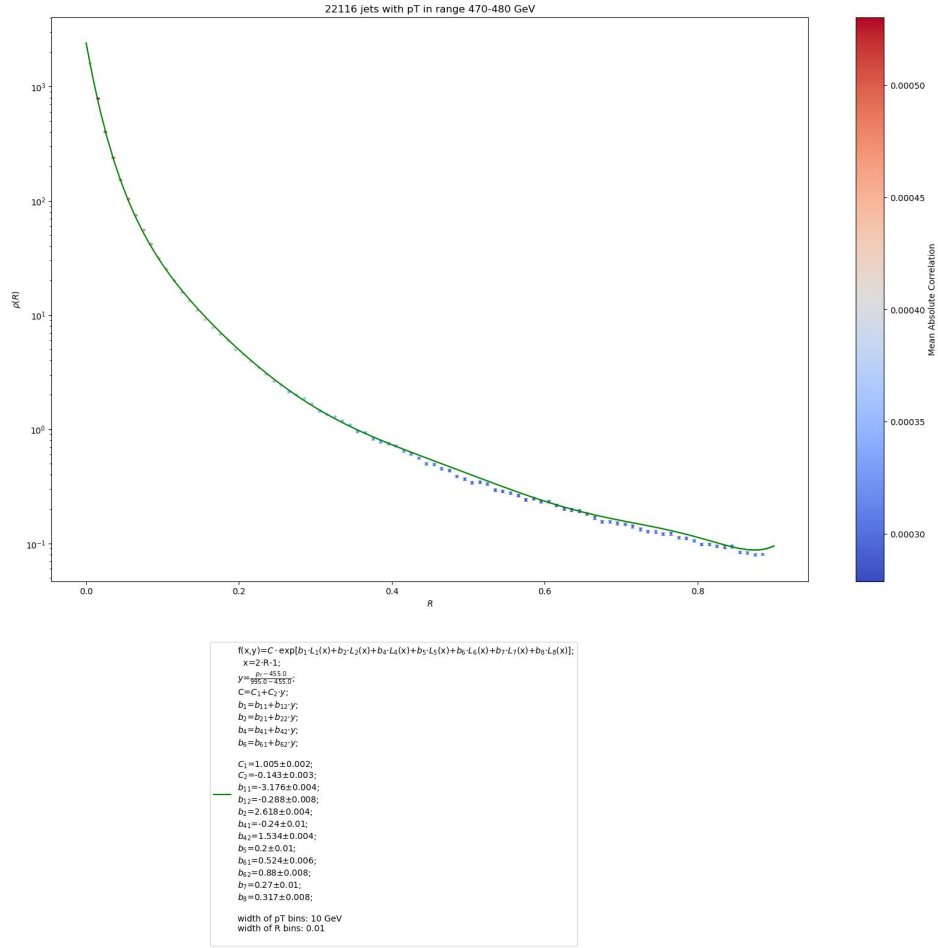


Figure 41: The final fit for QCDjj_pT_450_1200_train01.h5: FINAL_FIT_highpt.jpg

The only issue arises in the higher jet p_T range, where the model fails to capture the ideal tail behavior - instead of tapering off, the tail does not decline as expected.

10.4 Bayesian machine scientist

This last subsection allows for exploring alternative models. That said, as mentioned earlier, machine scientist may yield a function that is difficult to work with. In addition, because it does not incorporate data uncertainties, the resulting function may not ensure stable performance in the subsequent error-weighted fitting. As before, the data is structured so that the width of the jet p_T bins is 10 GeV, while the width of the radius bins is 0.01.

RunG_batch0.h5 file - probability distribution

Jet p_T range 300-500 GeV, ≈ 400000 jets loaded, $R_{max} = 0.75$

Machine scientist 2000 steps:

```

((((((sin((((R + _a0_) / _a9_) / ((R + _a0_) ** _a6_))) + (_a0_ * _a5_) ** pt_middles)
* (R ** _a11_) ** (pt_middles ** _a2_) + _a7_) / _a10_)

```

```

{"_a0_": -0.002757271313716955,

```



```

"a2_": -0.718451171271083,
"a5_": 85.60360674050297,
"a6_": 1.231599559670358,
"a7_": -5.551102369878217e-05,
"a9_": 2.273954508004391,
"a10_": 0.15121348929160183,
"a11_": 37.49953585484166,
"a1_": 1.0,
"a3_": 1.0,
"a4_": 1.0,
"a8_": 1.0,
"a12_": 1.0}

```

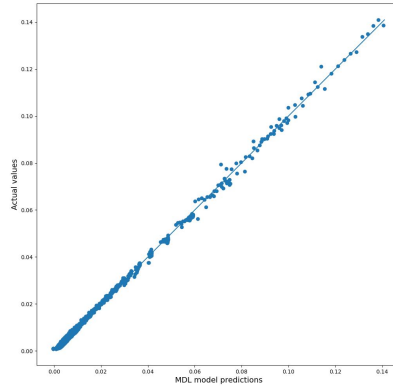


Figure 42: Quality of the probability model for RunG_batch0.h5 after 2000 steps: ms_prob_lowpt_2000.jpg

RunG_batch0.h5 file - probability density function

Jet p_T range 300-500GeV, ≈ 400000 jets loaded, $R_{max} = 0.75$

Machine scientist 2000 steps:

$$\frac{((((\exp((R * (_{a10_} ** pt_middles)))) + _{a7_}) ** _{a1_}) * (((pt_middles * _{a6_}) + _{a9_}) / (_{a1_} + pt_middles))) * R) + _{a2_} ** pt_middles}$$

```

{"_a1_": -1.180835600652233,
 "_a2_": 0.987189266832617,
 "_a6_": 0.0121791494053336,
 "_a7_": -0.9988267443494638,
 "_a9_": 2.259618173395361,
 "_a10_": 1.000357042489772,
 "_a0_": 1.0,

```

```

"a3_": 1.0,
"a4_": 1.0,
"a5_": 1.0,
"a8_": 1.0,
"a11_": 1.0,
"a12_": 1.0}

```

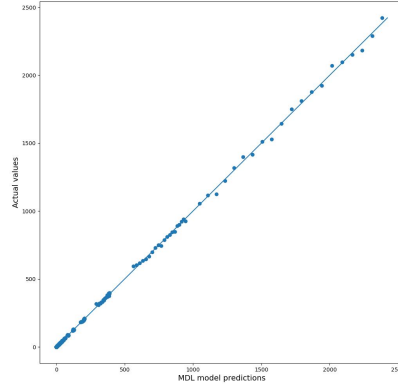


Figure 43: Quality of the probability density model for RunG_batch0.h5 after 2000 steps: ms_dens_lowpt_2000.jpg

QCDjj_pT_450_1200_train01.h5 file - probability distribution

Jet p_T range 450-1000GeV, ≈ 400000 jets loaded, $R_{max} = 0.9$

Machine scientist 3000 steps:

```

((((((pt_middles + _a8_) * ((pt_middles ** _a2_) * (R ** _a12_))) + ((R / _a7_) + ((R
** _a3_) * (_a5_ ** _a9_)))) * (_a9_ * _a2_)) / ((_a3_ + _a1_) + pt_middles)) * _a5_)

```

```

{"_a1_": 2237.0787838570122,
 "_a2_": -2.42575003999361,
 "_a3_": 0.19927590151212884,
 "_a5_": 0.004443476340136262,
 "_a7_": 16.219474949990126,
 "_a8_": -258.73252705457355,
 "_a9_": 0.8992038653292346,
 "_a12_": -0.7672415482817039,
 "_a0_": 1.0,
 "_a4_": 1.0,
 "_a6_": 1.0,
 "_a10_": 1.0,
 "_a11_": 1.0}

```

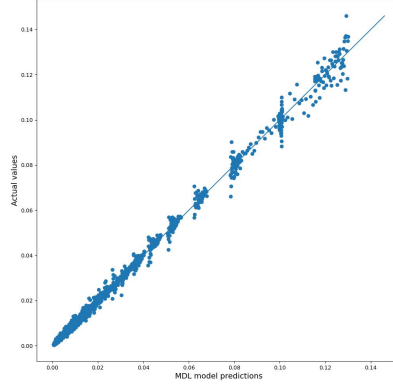


Figure 44: Quality of the probability model for QCDjj_pT_450_1200_train01.h5 after 3000 steps:
ms_prob_highpt_3000.jpg

QCDjj_pT_450_1200_train01.h5 file - probability density function

Jet p_T range 450-1000GeV, ≈ 400000 jets loaded, $R_{max} = 0.9$

Machine scientist 3000 steps:

```
((((((((-a4_ * _a4_) * _a7_) ** (pt_middles * _a12_)) + _a0_) + _a0_) ** (_a6_ * R)) +
_a2_) / (R / (_a11_ + ((-a4_ ** pt_middles) / R)))) * _a12_)
```

```
{ "_a0_": 4.517402226061548,
  "_a2_": -0.8545345597067965,
  "_a4_": 0.999760720675125,
  "_a6_": -20.7029615044342,
  "_a7_": 0.9933407738924719,
  "_a11_": -1.7808752481862429,
  "_a12_": -0.4823028433387591,
  "_a1_": 1.0,
  "_a3_": 1.0,
  "_a5_": 1.0,
  "_a8_": 1.0,
  "_a9_": 1.0,
  "_a10_": 1.0}
```

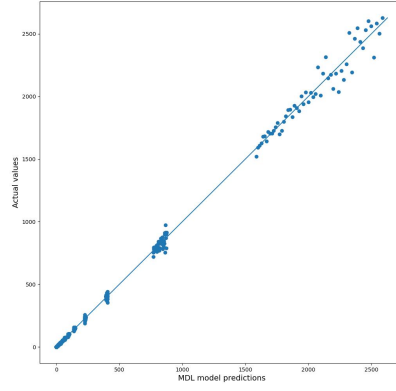


Figure 45: Quality of the probability density model for `QCDjj_pT_450_1200_train01.h5` after 3000 steps: `ms_dens_highpt.3000.jpg`

From the figures above it can be concluded, as expected, that the probability density function is easier to model than the probability distribution.

Bibliography

- [1] Qu, H., Li, C., & Qian, S. (2022). JetClass: A Large-Scale Dataset for Deep Learning in Jet Physics (1.0.0) [Data set]. Zenodo. <https://doi.org/10.5281/zenodo.6619768>
- [2] https://www.lhc-closer.es/taking_a_closer_look_at_lhc/0.momentum
- [3] Tomažič, T. Analiza generiranih podatkov, 2024.
- [4] K. Zoch, J. A. Raine, D. Sengupta, and T. Golling. RODEM Jet Datasets. Available on Zenodo: [10.5281/zenodo.12793616](https://doi.org/10.5281/zenodo.12793616). Aug. 2024. arXiv: 2408.11616 [hep-ph]
- [5] Amram, Oz, Anzalone, Luca, Birk, Joschka, Faroughy, Darius A., Hallin, Anna, Kasieczka, Gregor, Shih, David. (2024). Aspen Open Jets: a real-world ML-ready dataset for jet physics (Version 1.0) [Data set]. <http://doi.org/10.25592/uhhfdm.16505>
- [6] R. Guimer, I. Reichardt, A. Aguilar-Mogas, F. A. Massucci, M. Miranda, J. Pallars, M. Sales-Pardo, A Bayesian machine scientist to aid in the solution of challenging scientific problems. Sci. Adv. 6, eaav6971 (2020). <https://seeslab.info/publications/bayesian-machine-scientist-aid-solution-challenging-scientific-problems/>