

CHRISTIAN MICHELSEN  
NIELS BOHR INSTITUTE  
UNIVERSITY OF COPENHAGEN

A PHYSICIST'S  
APPROACH TO  
MACHINE LEARNING  
—  
UNDERSTANDING  
THE BASIC BRICKS

SUPERVISOR:  
TROELS PETERSEN  
NIELS BOHR INSTITUTE  
UNIVERSITY OF COPENHAGEN

Copyright © 2019

Christian Michelsen

`HTTPS://GITHUB.COM/CHRISTIANMICHELSEN`

Licensed under the Apache License, Version 2.0 (the “License”); you may not use this file except in compliance with the License. You may obtain a copy of the License at <http://www.apache.org/licenses/LICENSE-2.0>. Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

*First printing, December 2019*

# *Abstract*

Here will be a decent abstract at some point™.



# Contents

*Abstract*      iii

*Table of Contents*      v

*Foreword*      ix

1	<i>Introduction</i>	1
2	<i>Machine Learning Theory</i>	5
2.1	<i>Statistical Learning Theory</i>	5
2.2	<i>Supervised Learning</i>	6
2.3	<i>Generalization Bound</i>	7
2.3.1	<i>Generalization Bound for infinite hypotheses</i>	9
2.4	<i>Avoiding overfitting</i>	10
2.4.1	<i>Model Regularization</i>	10
2.4.2	<i>Cross Validation</i>	12
2.4.3	<i>Early Stopping</i>	14
2.5	<i>Loss functions</i>	14
2.5.1	<i>Evaluation Function</i>	16
2.6	<i>Decision Trees</i>	16
2.6.1	<i>Ensembles of Decision Trees</i>	17
2.7	<i>Hyperparameter Optimization</i>	19
2.7.1	<i>Grid Search</i>	20
2.7.2	<i>Random Search</i>	20
2.7.3	<i>Bayesian Optimization</i>	21
2.8	<i>Feature Importance</i>	23

3	<i>Danish Housing Prices</i>	27
3.1	<i>Data Preparation and Exploratory Data Analysis</i>	28
3.1.1	<i>Correlations</i>	30
3.1.2	<i>Validity of input variables</i>	32
3.1.3	<i>Cuts</i>	33
3.2	<i>Feature Augmentation</i>	33
3.2.1	<i>Time-Dependent Price Index</i>	34
3.3	<i>Evaluation Function</i>	35
3.4	<i>Initial Hyperparameter Optimization</i>	36
3.5	<i>Hyperparameter Optimization</i>	38
3.6	<i>Results</i>	40
3.7	<i>Model Inspection</i>	43
3.8	<i>Multiple Models</i>	45
3.9	<i>Discussion</i>	48
4	<i>Particle Physics and LEP</i>	53
4.1	<i>The Standard Model</i>	53
4.2	<i>Quark Hadronization</i>	55
4.3	<i>The ALEPH Detector and LEP</i>	56
4.4	<i>Jet clustering</i>	58
4.5	<i>The variables</i>	58
5	<i>Quark Gluon Analysis</i>	63
5.1	<i>Data Preprocessing</i>	63
5.2	<i>Exploratory Data Analysis</i>	64
5.3	<i>Loss and Evaluation Function</i>	67
5.4	<i>b-Tagging Analysis</i>	68
5.4.1	<i>b-Tagging Hyperparameter Optimization</i>	68
5.4.2	<i>b-Tagging Results</i>	69
5.4.3	<i>b-Tagging Model Inspection</i>	71
5.5	<i>Truncated Uniform PDF</i>	72
5.6	<i>g-Tagging Analysis</i>	73
5.6.1	<i>Permutation Invariance</i>	73
5.6.2	<i>g-Tagging Hyperparameter Optimization</i>	74
5.6.3	<i>PermNet</i>	74
5.6.4	<i>1D Comparison of LGB and PermNet</i>	75
5.6.5	<i>g-Tagging Results</i>	76

6	<i>Discussion and Outlook</i>	83
7	<i>Conclusion</i>	85
7.1	<i>Tufte-L<sup>A</sup>T<sub>E</sub>X Website</i>	85
7.2	<i>Tufte-L<sup>A</sup>T<sub>E</sub>X Mailing Lists</i>	85
7.3	<i>Getting Help</i>	85
A	<i>Housing Prices Appendix</i>	87
B	<i>Quarks vs. Gluons Appendix</i>	115
	<i>List of Figures</i>	128
	<i>List of Tables</i>	130
	<i>Index</i>	139





## *Foreword*

## *Part I*

The first part of this thesis deals with the introductory theory of machine learning and its predictive power in estimating Danish housing prices.

This subproject was done in collaboration with Boligsiden without whom it would not have been possible. During this project, common python data science tools from the SciPy ecosystem[[81](#)] such as NumPy, Matplotlib, Pandas, Scikit-Learn, Scipy has been used extensively and should thus also be mentioned.



## *Part II*

The second part of this thesis deals with particle physics and the discriminatory power of machine learning for quark-gluon identification and subsequent analysis.



## 4. Particle Physics and LEP

*“Not only is the Universe stranger than we think, it is stranger than we can think.”*

---

— Werner Heisenberg

The aim of this chapter is to introduce the reader to the level of particle physics required for understanding the following chapter, in particular introducing the Standard Model in [section 4.1](#), the theory behind quark hadronization in [section 4.2](#), and the ALEPH detector at LEP in [section 4.3](#). The goal is not to make a deep and thorough introduction to the field as this is not needed for the following analysis along with the fact that the author is no particle physicist himself.

### 4.1 The Standard Model

The *Standard Model* (SM) [[40](#), [68](#), [83](#)] of particle physics is the currently best known description of the elementary particles and thus describes the fundamental building blocks of our Universe. An overview of the particles explained by the Standard Model is shown in the typical tabular form seen in [Figure 4.1](#). In general, particles comes in two categories: *bosons* and *fermions*.

The fermions, the left part of the figure, are particles with half-integer spin that obey Fermi-Dirac statistics and are further subdivided into *quarks* (upper left in figure) and *leptons* (lower left). The quarks interact with all of the four known forces<sup>1</sup>, including the strong force. In contrary the leptons do not interact with the strong force. Quarks are never observed freely but are always combined into *hadrons* due to *color confinement* which is further explained in [section 4.2](#). An example of this are protons which consists of two up-quarks and a down-quark. Leptons exist as either the charged leptons<sup>2</sup> or as neutral leptons, the so-called neutrinos<sup>3</sup>. The fermions come in three generations with increasing mass.

The bosons, the right part of the figure, are the force-carrying particles (with integer spin and which obey Bose-Einstein statistics) where the gluon  $g$  mediates the strong nuclear force (color charge), the photon  $\gamma$  mediates the electromagnetic force (charge), and the two  $W^\pm$  and the  $Z$  bosons the weak nuclear force (weak isospin). The Higgs boson  $H$ , experimentally discovered in 2012 [[32](#), [33](#)],

<sup>1</sup> Gravity, electromagnetism, and the strong and weak force.

<sup>2</sup> The electron  $e$ , the muon  $\mu$ , and the tau  $\tau$ .

<sup>3</sup> The electron neutrino  $\nu_e$ , the muon neutrino  $\nu_\mu$ , and the tau neutrino  $\nu_\tau$ .

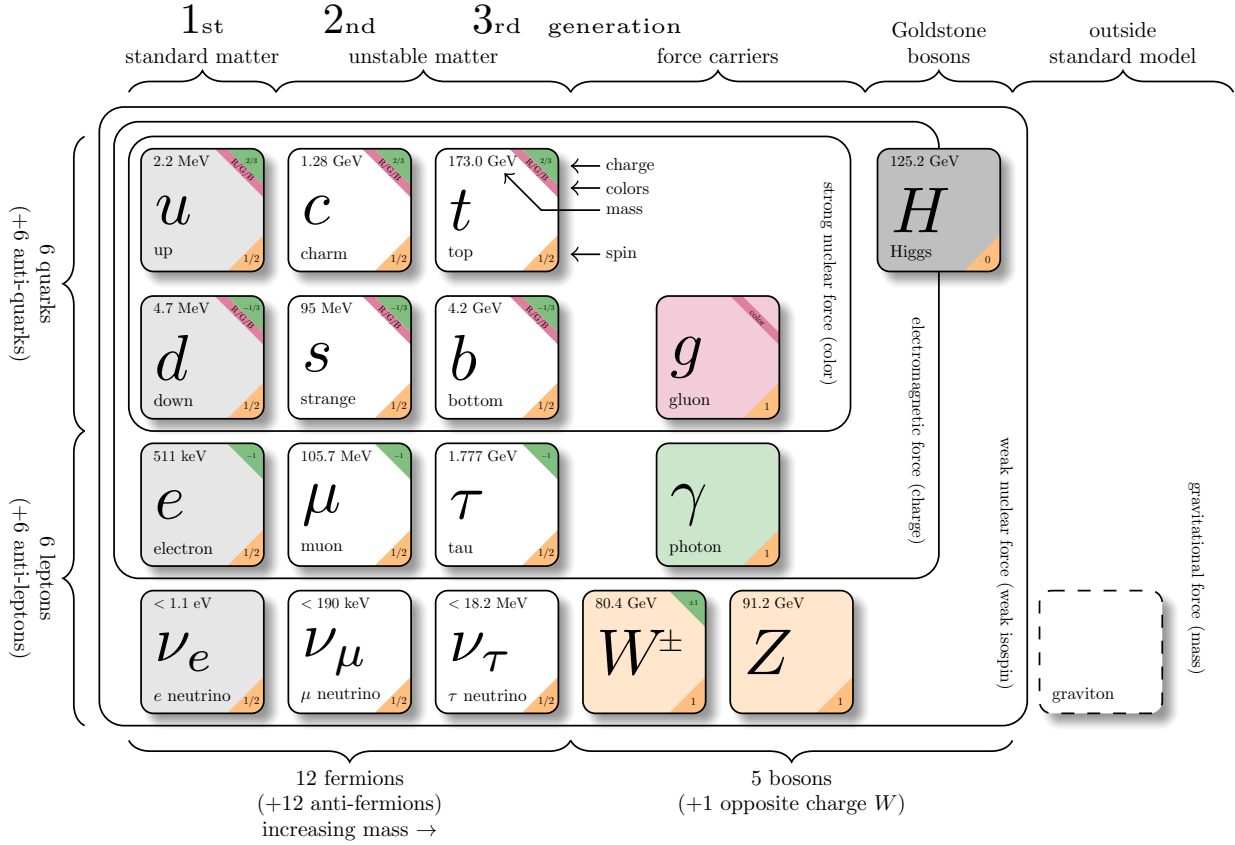


Figure 4.1: The Standard Model. Inspired by Purcell [64] using the template by Burgard [29] with manually updated masses according to Particle Data Group et al. [60].

does not mediate any forces but interacts with all massive particles and explains why particles have mass.

All particles have antiparticles which are particles with opposite charge but the same mass. Some particles are their own antiparticles<sup>4</sup>, such as the  $Z$ . At the Large Electron Positron collider (LEP), see section 4.3, electrons  $e^-$  and their antiparticles positrons  $e^+$  were collided at an energy of around 91 GeV. This particular energy was chosen since this is at the resonance peak of the  $Z$ . Its mass distribution follows a Cauchy distribution (also known as Breit-Wigner) with mean<sup>5</sup>  $m_Z = (91.1876 \pm 0.0021) \text{ GeV}$  and a full width of  $\Gamma_Z = (2.4952 \pm 0.0023) \text{ GeV}$ : LEP was as such a  $Z$ -factory. The  $Z$ , however, is only very short-lived with a half-life of  $1/\Gamma_Z \sim 2.6 \times 10^{-25} \text{ s}$ . The decay mode for this unstable  $Z$  particle is primarily to hadrons ( $(69.91 \pm 0.06) \%$ ) where the ratio ( $R$ ) for  $b$ -quarks is  $R_b = (Z \rightarrow b\bar{b}) = (15.12 \pm 0.05) \%$  and  $R_g = (Z \rightarrow ggg) < (1.10 \pm 0.05) \%$  for gluons [60]. The fact that the  $Z$  is its own anti-particle forces its decay to be a particle–anti-particle decay (due to charge-conservation) where antiparticles are written with a bar on top, e.g. the  $\bar{b}$ -quark is the antiparticle of the  $b$ -quark.

<sup>4</sup> The photon, the  $Z$ , and the Higgs.

<sup>5</sup> Calculated in natural units where  $c = \hbar = 1$  which will also be used throughout this thesis.

## 4.2 Quark Hadronization

The electron-positron  $e^+e^-$  annihilations at LEP are complicated events that require advanced high-energy particle physics theory to be properly understood. Most of the aspects of the process is well-described by now, however, especially the hadronization process is still an area of active research. To better get an overview of the different stages of the  $e^+e^-$  annihilations, see the Feynman diagram in Figure 4.2.

Reading from left to right, the electron and the positron annihilates to a  $Z$ . This interaction is well-described by quantum electrodynamics (QED), a theory that has been around for more than 60 years by now. As mentioned in the previous section, the  $Z$  has several decays modes, yet most of these are background processes of no interest in this project and the focus for now will be the decay mode  $Z \rightarrow q\bar{q}$  ( $Z$  to quark-anti-quark) as seen in the Feynman diagram. The particles produced by the  $Z$ -decay are called primary *partons*. Since this process involves quarks, and thus color charge, QED is no longer an adequate theory: quantum chromodynamics (QCD) is needed [15]. The  $q\bar{q}$  pairs in this example acts as (color) dipoles from which a gluon can radiate. It can be shown with QCD that the gluon can only be radiated inside the cone that the  $q\bar{q}$  pairs spans [23]. As mentioned in the introduction, quarks cannot exist freely (due to *confinement*) and we therefore cannot observe the individual partons in a  $q\bar{q}g$  event produced in the Feynman diagram. Confinement is basically the QCD principle saying that quarks are always confined or bound inside hadrons. The initial partons (carrying color charge) are converted to (color-neutral) hadrons by non-perturbative QCD processes in what is called *hadronization*, and these hadrons can be measured.

The hadronization process is not yet fully modelled and currently two competing models for predicting the hadronization pattern exists: the Lund string model and the cluster model. In this project only the former of the models will be used. The Lund string model [14] is the theoretical framework underlying the widely used Monte Carlo event generator PYTHIA [71]. The string model is based on the observation that (color) field lines between quarks seem to compress into a tube-like region mediated by gluons, see the top part of Figure 4.3. The field can be described by a linearly rising potential  $V(r) = \kappa r$  at large distances<sup>6</sup>, where  $r$  is the distance and  $\kappa$  is the strength of the potential [28]. This field is similar to the (constant) force of a string:  $V(r) = \kappa r \Rightarrow F(r) = -\kappa$  where  $\kappa$  is the to be regarded as the spring tension. As quarks move apart, the potential energy stored in the “string” increases until it is large enough to “snap” and convert its potential energy into mass. This mass energy is released with the production of a new  $q\bar{q}$  pair as this energetically favorable, see the rest of Figure 4.3.

An example of the hadronization process, or the transition from initial partons to final hadrons is sketched in Figure 4.4. Here the

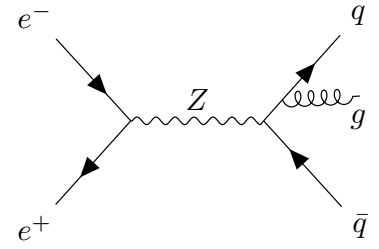


Figure 4.2: Feynman diagram showing the  $e^+e^- \rightarrow Z^0$  production at LEP. The  $Z$  has several decay modes where the  $Z \rightarrow q\bar{q}g$  is shown here.

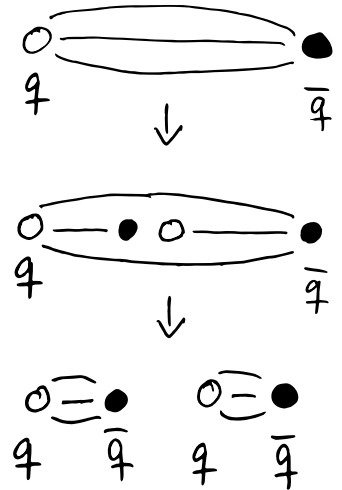


Figure 4.3: Illustration of the quarks splitting as explained by the Lund string model. For large charge separation the (color) field lines seem to be compressed to a tube-like region, where the strong interactions are mediated by the massless gluons (that couple to the color charge of quarks). When the two quarks are separated enough, the potential energy is released by the production of a new  $q\bar{q}$  pair.

<sup>6</sup> At small distances a Coulomb term has to be included, however, this term is assumed to be negligible by the Lund string model.



production of two kaons  $K^-$  and  $K^+$ , and two pions  $\pi^-$  and  $\pi^0$  are shown. Since particles are created by “splits” in the “string”, and the fact that there is energy-momentum conservation, they all have to share the total energy stored in the string. This is described by the fragmentation function:

$$f(z) \propto \frac{(1-z)^a}{z} \exp\left(-\frac{bm^2}{z}\right), \quad (4.1)$$

where  $0 \leq z \leq 1$  is the remaining momentum that the new hadron takes,  $a$  and  $b$  are constants, and  $m$  is the mass<sup>7</sup> [23]. When the system runs out of available momentum, it will stop producing new hadrons and the fragmentation function thus explains the distribution of final state particles. The Lund string model can be extended from only  $q\bar{q}$  events to  $q\bar{q}g$  events where it predicts cones spanning the angular regions  $qg$  and  $\bar{q}g$  should receive enhanced particle production compared to the  $q\bar{q}$  region. This prediction by the Lund string model is also measured in  $e^+e^-$  collisions [28].

<sup>7</sup> Where  $m \rightarrow m_\perp$  for particles with transverse momentum.

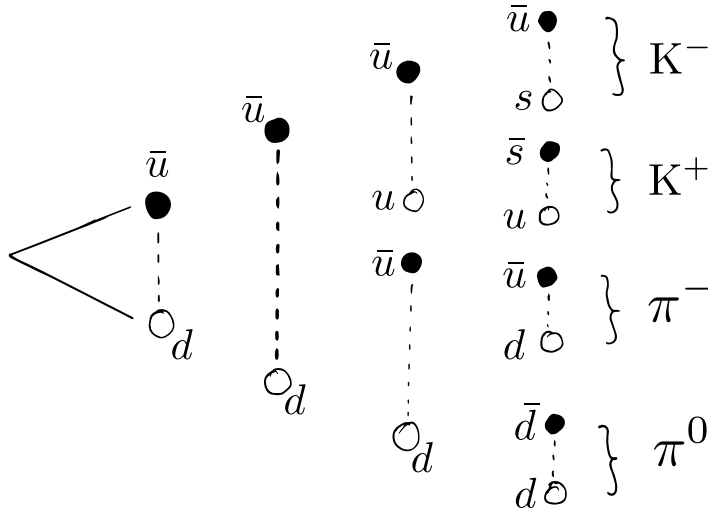


Figure 4.4: Illustration of the hadronization process by which  $\bar{u}$ - and  $d$ -quarks decay into four different mesons. The theoretical strings are shown as dashed lines and particles as circles, where filled circles are antiparticles.

The initial partons produced as  $Z$  decay therefore decay to final state hadrons<sup>8</sup> which create a whole “shower” in the direction of the initial parton: this is called a *parton shower* and it is this parton shower observed as particles, a *jet*, that is measured in the detector. The reverse computation from tracks measured in the detector is done with the use of *jet clustering* algorithms. The detector and the clustering algorithms are described in the following section.

<sup>8</sup> To either mesons which consist of two quarks (color-anti-color) or baryons (r-g-b) which consist of three quarks.

### 4.3 The ALEPH Detector and LEP

The Large Electron Positron collider (LEP) was a particle collider at CERN in Switzerland operating from 1989 to 2000. It collided counter-rotating bunches of electrons and positrons in a giant ring with a circumference of more than 26 km. The first phase, LEP1, ran from 1989 to 1995 at the  $Z$  resonance 91 GeV and the second phase, LEP2, continued afterwards closer to 200 GeV for  $W^+W^-$

pair production [15], however, it is only the data collected at the energy around  $\sqrt{s} = 91.3 \text{ GeV}$  called the *Z peak data* that is used throughout the rest of this project. There were four independent detectors at the LEP experiment, one of them ALEPH<sup>9</sup>.

<sup>9</sup> Together with DELPHI, L3, and OPAL.

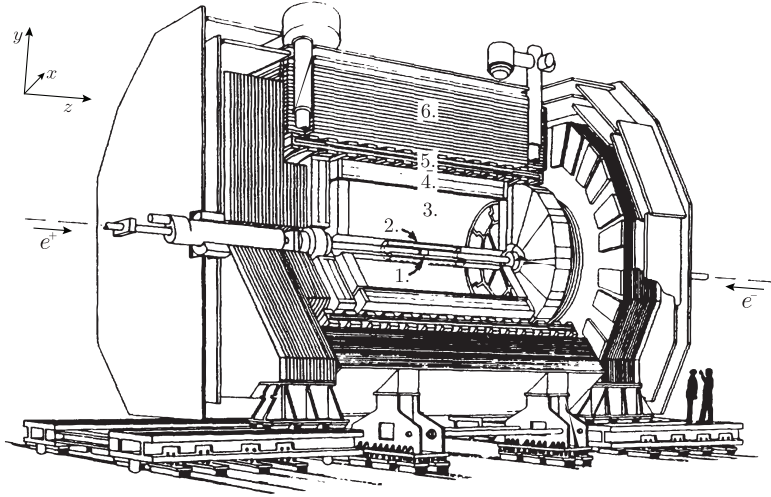


Figure 4.5: The ALEPH detector at LEP. 1) Vertex detector (VDET). 2) Drift chamber (ITC). 3) Time projection chamber (TPC). 4) Electromagnetic calorimeter (ECAL). 5) Superconducting magnet coil. 6) Hadron calorimeter (HCAL). Adapted from Buskulic et al. [30].

The *apparatus for LEP physics* (ALEPH) was a particle detector at LEP with a wide coverage, almost  $4\pi$ , consisting of cylindrical subdetectors, see Figure 4.5, with the coordinate system shown in the upper left corner<sup>10</sup>. The polar angle  $\theta$  is illustrated in Figure 4.6 together with the transverse (longitudinal) momentum  $p_{\perp}$  ( $p_L$ ) and the azimuthal angle  $\phi$  in Figure 4.7. The ALEPH detector was designed to measure the energy deposited in calorimeters by charged and neutral particles, measure the momenta of charged particles, measure the distance of travel of short-lived particles, and to identify the three lepton flavors (electron, muon, tau) [30]. As can be seen in Figure 4.5, ALEPH consisted of five subdetectors (the vertex detector (VDET), the drift chamber (ITC), and the time projection chamber (TPC)) and two calorimeters (the electromagnetic (ECAL) and the hadronic calorimeters (HCAL)).

The three innermost detectors allow for precise tracking of the charged particles produced in the parton shower and the two outer calorimeters of precise energy measurements for both charged and neutral particles going through the detector.

A hadronic event from a parton shower may leave a score of charged tracks resulting in hundreds of hits in the detectors (VDET, ITC, and TPC) which are fitted<sup>11</sup> with Kalman filters [49] to obtain global track fits, of which bad charged tracks are discarded for further analysis. The tracks are helical due to the presence of a 1.5 T magnetic field which curves the charged particles according to their transverse momentum,  $p_{\perp}$ .

The energy resolution  $\sigma$  of the calorimeters, or the *calorimeter performance*, is expected to increase with  $\sqrt{E}$ . In fact, it was found at ALEPH that the energy dependence of the resolution follows the

<sup>10</sup> The z-axis pointing along the beam direction, the y-axis pointing upwards, and the x-axis pointing towards the center of LEP.

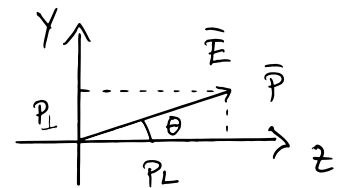


Figure 4.6: The polar angle  $\theta$  defined in the  $zy$  coordinate system

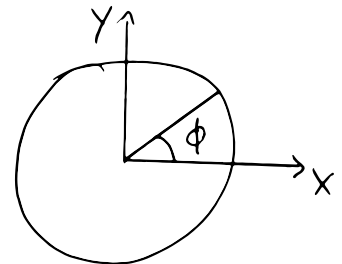


Figure 4.7: The azimuthal angle  $\phi$  defined in the  $xy$  coordinate system.

<sup>11</sup> the process of fitting tracks is called *track reconstruction* in high energy particle physics.

parametrization [30]:

$$\sigma(E) = \left( (0.59 \pm 0.03) \cdot \sqrt{E/\text{GeV}} + (0.6 \pm 0.3) \right) \text{ GeV}. \quad (4.2)$$

Even though  $\sigma(E)$  increases with  $E$ , the relative resolutions improves with higher energies. Since one never measures Nature directly, the results one obtains in a measurement are thus products of both model and experimental uncertainties folded together. To unfold the measurements to obtain experiment-independent results, the uncertainties are important to understand. Of course there are dozens of other uncertainties in an advanced experiment like ALEPH, however, the energy dependence is the primary focus in this project.

#### 4.4 Jet clustering

Since the initial partons created as decay products from the  $Z$  are unstable themselves, what is measured in the detector is a whole shower of hadrons seen as charged tracks in the detectors and energy deposits in the calorimeters. However, say that the  $Z$  decayed to a  $b\bar{b}$  event. In this case the two  $b$ 's would be back-to-back and the final hadrons would be observed approximately in the same direction as the  $b$ 's were created. The interest of the experiment is not to measure the final hadrons, but rather to infer information about the initial quarks and gluons. This is done via the reverse-engineering process called *jet clustering*. Over the years many clustering algorithms have been developed, however, most of these are younger than LEP. In the ALEPH experiment the JADE algorithm was used [19]. JADE is a sequential recombination algorithm where final state particles are initially described as individual so-called pseudo-jets which are then recursively merged to larger jets according to their inter-jet distance  $d_{ij}^2$ . The distance measure for JADE is:

$$d_{ij}^2 = \frac{2E_i E_j (1 - \cos \theta_{ij})}{E_{\text{vis}}^2}, \quad (4.3)$$

where  $E_{\text{vis}}$  is the visible energy<sup>12</sup> and  $\theta_{ij}$  is the angle between jet  $i$  and  $j$ . The JADE algorithm computes  $d_{ij}^2$  for all combinations of jets and merges the two jets with the lowest  $d_{ij}^2$ , continuing like that recursively until  $\min(d_{ij}^2) > d_{\text{cut}}^2$  for some predefined value of  $d_{\text{cut}}^2$ . In the dataset at hand, only the final jets were available and not the jet constituents, unfortunately.

<sup>12</sup> The total sum of energies in the event.

#### 4.5 The variables

The overall goal of the project is to be able to discriminate quarks and gluons using only vertex variables. The reason for the last condition is that the goal is to better understand the shape distributions of gluons in which there is still significant differences between Monte Carlo (MC) simulations and Data. Therefore only vertex

variables will be used to avoid any biases introduced by using shape-related variables to detect differences in shape-distributions. The vertex variables are a subset of all variables which include the three variables `projet`, `bqvjet`, and `ptlrel`. These three particular variables have each shown discriminatory power in separating  $b$ -quarks from light quarks and gluons.

**projet** : For each track in the jet an impact parameter  $\delta$  is computed. This parameter is the minimum distance between the estimated  $Z$  decay point and the track itself and its sign depends on whether or not the point of closest approach is in front of or behind the  $Z$  decay point (relative to the momentum). From  $\delta$  the significance  $S$  – which is  $\delta/\sigma_\delta$  – is computed and is thus a measure of the certainty of a measured track. High values of  $S$  is typically an indicator of  $b$  jets, since long-lived particles typically decay in front of the  $Z$  relative to the jet direction, while  $uds$ -jets might as well have negative values of  $S$ . From  $S$  the track probability  $\mathcal{P}_{\text{track}}$  of a track originating at the decay point of the  $Z$  can be computed, which can further be aggregated across all tracks within a jet to form the jet probability  $\mathcal{P}_{\text{jet}}$  which **projet** is a function of [36]. Whether or not  $\mathcal{P}_{\text{jet}}$  is strictly a probability can be discussed but it is related to the probability of all tracks within a jet to originate from long-lived particles, which itself is a good indicator of being a  $b$ - (or  $c$ -) jet. This variable further has the advantage of being independent of any vertex algorithm.

**bqvjet** : For any jet with good<sup>13</sup> charged tracks, a fit with a (hypothetical) secondary vertex is performed. The difference in  $\chi^2$  between the null hypothesis that all good tracks originate from the same primary vertex and the alternative hypothesis that a secondary vertex exists in addition to the primary one is calculated. For the long-lived massive  $b$  and  $c$  quarks this typically results in large differences in  $\chi^2$  compared to  $uds$ - and gluon jets which have much lower  $\Delta\chi^2$ -values [15]. The **bqvjet** is related to the  $\Delta\chi^2$ -value from the secondary vertex algorithm. This value is dependent of the vertex algorithm, but still explores other areas of phase space than **projet**, however, they are still very correlated. The linear correlations  $\rho_{q_i}$  between **projet** and **bqvjet** for  $q_i$  quarks (MC truth) are  $\rho_b = 0.80, \rho_c = 0.65, \rho_{uds} = 0.23, \rho_g = 0.29$ .

<sup>13</sup> Meaning that there are at least four TPC hits and the fit has a reduced  $\chi^2$  of less than four [15].

**ptlrel** : If any leptons (in the case of  $e^\pm$  or  $\mu^\pm$ ) are measured by the detector, this is a good sign of the jet originating from a  $b$ -quark. The high mass of the  $b$  quark leads to high  $p_\perp$  for the leptons relative to the jet axis which is exactly measured by **ptlrel**.

The fact that the heavy  $b$ -quarks have much longer lifetimes than the lighter  $uds$ -quarks stems from their much lower inter-coupling

magnitudes written as the CKM matrix  $\mathbf{V}$  [60]:

$$\mathbf{V} = \begin{matrix} & \begin{matrix} d & s & b \end{matrix} \\ \begin{matrix} u \\ c \\ t \end{matrix} & \begin{pmatrix} 0.97446 & 0.22452 & 0.00365 \\ 0.22438 & 0.97359 & 0.04214 \\ 0.00896 & 0.04133 & 0.99911 \end{pmatrix} \end{matrix}. \quad (4.4)$$

The matrix element  $|V_{ij}|^2$  is proportional to the transition-probability of quark  $i$  transitioning to quark  $j$ . From the CKM matrix it can be seen that  $u$  and  $d$  quarks couples strongly together, likewise with  $c$ - $s$  and  $b$ - $t$  quark pairs. When a  $Z$  decays into a  $b$ -quark, this quark couples strongly with the top quark, however, due to the high mass of the top quark compared to the center-of-mass energies at LEP1, the  $b$ -quark cannot decay into a  $t$ -quark but must (almost always) decay to a  $c$ , however, still with low probability,  $V_{bc} \ll 1$ . This, together with the fact that  $V_{bu} \ll V_{bc}$  explains the long life-time of  $b$  quarks. This is also why the three variables above are very common variables for  $b$ -tagging algorithms. That  $c$ -quarks also have relative long life-times are not due to the CKM elements, as for  $b$ -quarks, but rather due to the  $c$ -decay being governed by the weak force through virtual  $W^*$  bosons, a force that is much slower than the strong force. This also happens for  $b$ -quarks which further explains why  $c$ -quarks share many similarities with  $b$ -quarks but also resembles resembles light-quarks (which are very short-lived.).

The rest of the non-vertex variables are:

`ejet` : The energy of the jet  $E_{\text{jet}}$

`costheta` : The cosine of the  $\theta$  angle defined in Figure 4.6:  
 $\cos \theta$ .

`phi jet` : The angle  $\phi$  of defined in Figure 4.7:  $\phi$ .

`sphjet` : The sphericity tensor  $\mathbf{S}$  is defined as:

$$S^{(\alpha\beta)} = \frac{\sum_{i=1}^N p_i^{(\alpha)} p_i^{(\beta)}}{\sum_{i=1}^N |p_i|^2} \quad \alpha, \beta \in \{x, y, z\}, \quad (4.5)$$

and the sphericity is determined as  $S = \frac{3}{2}(\lambda_2 + \lambda_3)$  where  $\lambda_1 \geq \lambda_2 \geq \lambda_3$ ,  $\lambda_1 + \lambda_2 + \lambda_3 = 1$  are the three eigenvalues of the sphericity tensor. The sphericity  $0 \leq S \leq 1$  is a measure of the angular distribution of the tracks in a jet. When  $S = 0$  the jets form a perfect sphere, compared to  $S = 1$  for a perfect jet. The `sphjet` variable is the sphericity of the jet when calculated in its boosted rest frame, also known as *boosted sphericity*.

`pt2jet` : The sum of the square of transverse momentum w.r.t. the jet axis:  $\sum_i p_{\perp,i}^2$ .

`muljet` : The multiplicity of the jet.

For further details about the variables, see Armstrong [15].

The variables explained above are all used in the following analysis where the machine learning model is trained on only the vertex variables to probe differences in the shape-distributions of the shape-variables. The goal of this is to better understand the gluon hadronization process to minimize differences in MC simulations and ultimately get a better understanding of the rules governed by Nature.



## 5. Quark Gluon Analysis

*“Research is what I am doing I don’t know know what I’m doing.”*

---

— Wernher von Braun

AS ANY DEDICATED READER can clearly see, the Ideal of practical reason is a representation of, as far as I know, the things in themselves; as I have shown elsewhere, the phenomena should only be used as a canon for our understanding. The paralogisms of practical reason are what first give rise to the architectonic of practical reason. As will easily be shown in the next section, reason would thereby be made to contradict, in view of these considerations, the Ideal of practical reason, yet the manifold depends on the phenomena.

### 5.1 Data Preprocessing

The data consists of 43 data files taken between 1991 and 1995 totalling 3.5 GB (Data). Along with this comes 125 files based on Monte Carlo (MC) simulations (8.4 GB) and additional 42 MC-files with only  $b$ -quark events (MCb) simulated (2.1 GB). The data files which are in the form of *Ntuples*, ROOT’s data format [27], are converted to HDF5-files by using uproot [7]. While iterating over the *Ntuples*, some basic cuts are applied before exporting the data to HDF5. The first one being that the (center of mass) energy  $E$  in the event has to be within  $90.8 \text{ GeV} \leq E \leq 91.6 \text{ GeV}$  to only use the  $Z$  peak data. The second one being that the sum of the momenta  $p_{\text{sum}}$  in each event is  $32 \text{ GeV} \leq p_{\text{sum}}$  to remove any  $Z \rightarrow \tau^+ \tau^-$  events. To ensure a primary vertex, at least two good tracks are required where a good track is defined as having 7 TPC hits and  $\geq 1$  silicon hit. Finally it is required that the cosine of the thrust axis polar angle, which is the angle between the thrust axis and the beam, is less than or equal to 0.8 to avoid any low angle events since the detector performance worsens significantly in that region. These cuts were standard requirements for the ALEPH experiment.

One last cut which was experimented with was the threshold value for *jet matching*. The jet matching is the process of matching the jet with one of the final state quarks. The jet is said to be matched if the dot product of between the final quark momentum



and the jet momentum is more than then threshold value. Higher thresholds means cleaner jets but at the expense of less statistics. A jet matching threshold of 0.90 was found to be a good compromise between purity and quantity where 97.8 % of all 2-jet events are matched and 96.7 % of all other jets were matched<sup>1</sup>.

The data structure is quite differently structured in the Ntuples compared to normal structured data in the form of tidy data [84]. The data is organized such that one iterates over each event where the variables are variable-length depending on the number of jets in the events; this is also known as *jagged* arrays. The data is un-jagged<sup>2</sup> before exporting to HDF5-format and only the needed variables are kept. This reduces the total output file to a 2.9 GB HDF5-file for both Data, MC, and MCb.

The number of events for each number of jets can be seen in Table 5.1 for the Data and in Figure 5.2 for the MC and MCb.

## 5.2 Exploratory Data Analysis

Since the machine learning models are only trained on the three vertex variables `projet`, `bqvjet`, and `ptljet` – see chapter 4 for a deeper introduction to these variables – these variables will be the primary focus of this section. Given the fact that MC-simulated data exists, the truth of each simulated event is also known. This allows us visualize the difference between the different types of quarks. In the MC simulation each event are generated such that the type of quark, or *flavor*, is known and assigned the variable `flevt`. The mapping from flavor to `flevt` is:

Flavor:	<i>bb</i>	<i>cc</i>	<i>ss</i>	<i>dd</i>	<i>uu</i>
<code>flevt</code> :	5	4	3	2	1

In addition to knowing the correct flavor, we define that an event is *q-matched* if one, and only one, of the jets are assigned to one of the quarks, one, and only one, of the jets are assigned to the other quark, and no other jets are matched to any of the quarks. We can then define what constitutes a *b-jet*: if it has `flevt` = 5, the entire event is *q-matched*, and the jet is matched to one of the quarks. Similarly we define *c-jets* only with the change that `flevt` = 5, and *uds-jets* with `flevt` ∈ {1,2,3}. A gluon jet is defined as an any-flavor event which is *q-matched* but the jet is not assigned to any of the quarks. Strictly speaking, this means that *g-jet* is not 100 % certain of being a gluon, however, since the MC simulation does not contain this information this is the only option. Due to the *q-match* criterion this also means that some jets are assigned the label “non-*q-matched*” which is regarded as background. The distribution of different types of jets can be seen in Table 5.3 and shown as relative numbers in Table B.1 in the appendix.

With the criteria defined above for what constitutes a specific type of jet the 1D-distributions for the three vertex variables is plot-

<sup>1</sup> Compare this to 98.5 % and 97.8 % for a threshold of 0.85 or 95.9 % and 93.9 % for a threshold of 0.95.

<sup>2</sup> Such that e.g. a 3-jet event will figure as three rows in the dataset.

	jets	events
2	2 359 738	1 179 869
3	3 619 290	1 206 430
4	854 336	213 584
5	52 775	10 555
6	510	85
Total	6 886 649	2 610 523

Table 5.1: The dimensions of the dataset for the actual Data. The numbers in the jet columns are the number of events multiplied with the number of jets; e.g.  $85 \cdot 6 = 510$ .

	jets	events
2	7 293 594	3 646 797
3	10 780 890	3 593 630
4	2 241 908	560 477
5	103 820	20 764
6	588	98
Total	20 420 800	7 821 766

Table 5.2: The dimensions for the MC and MCb datasets.

	<i>b</i>	<i>c</i>	<i>uds</i>	<i>g</i>	non- <i>q</i> -matched
2	2 713 454	944 380	2 125 900	0	1 509 860
3	2 433 878	964 212	2 129 218	3 365 969	1 887 613
4	326 264	156 332	336 548	1 012 198	410 566
5	10 332	5960	12 668	54 525	20 335
6	42	26	52	320	148
Total	5 483 970	2 070 910	4 433 012	4 604 386	3 828 522

Table 5.3: Number of different types of jets for MC and MCb. See also Table B.1 in the appendix for relative numbers.

ted in Figure 5.1. For all three subplots the histograms are shown with a logarithmic  $y$ -axis, all  $b$ -jets in blue,  $c$ -jets in red,  $g$ -jets in green and all jets in orange. In fully opaque color are shown the distributions for 2-jet events, in dashed (and lighter color) 3-jet events, and in semi-transparent 4-jet events. In the left subplot the `projet` variable is plotted where it can be seen that high values of `projet` tend to indicate  $b$ -jets. In the middle subplot `bqvjet` is plotted which shares many similarities with the `projet`-variables, including that high values indicate  $b$ -jets. In the right subplot the `ptljjet` is plotted. This variable has many zeros in it which correlates with mostly with gluon<sup>3</sup> and large values are mostly due to  $b$ -jets. In general it is clear to see how the differences in distribution between the 2-, 3-, and 4-jet events are minor, with the one exception of 2-jet events which does not contain any gluons at all.

<sup>3</sup> Around 98 % of all  $g$ -jets are zeros compared to  $\sim 82$  % for  $c$ -jets and  $\sim 70$  % for  $b$ -jets.

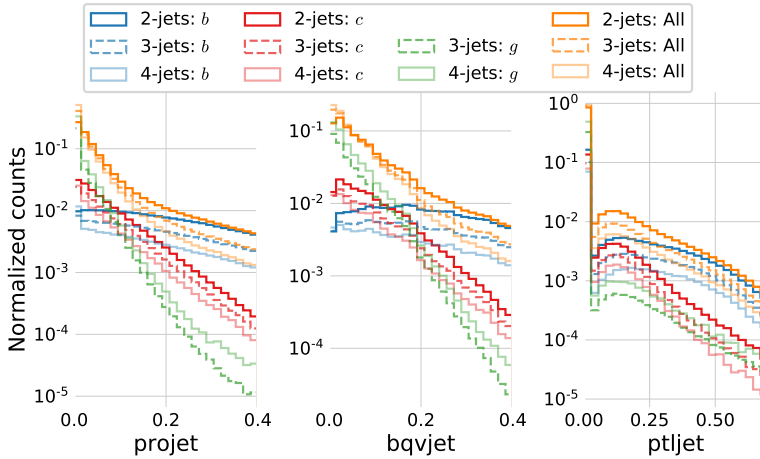


Figure 5.1: Normalized histograms of the three vertex variables: `projet`, `bqvjet`, and `ptljjet`. In blue colors the variables are shown for true  $b$ -jets, in red for true  $c$ -jets, in green for true  $g$ -jets, and in orange for all of the jets (including non  $q$ -matched). In fully opaque color are shown the distributions for 2-jet events, in dashed (and lighter color) 3-jet events, and in semi-transparent 4-jet events. Notice the logarithmic  $y$ -axis, that there are no  $g$ -jets for 2-jet events (as expected), and that all of the distributions are very similar not matter how many jets.

Even though there are only three vertex variables, it is difficult to properly get an intuition about how easily separated they different types of jets are. Since there are millions of points a single 3D scatter plot quickly becomes overcrowded in one wants to plot all jets. We apply dimensionality reduction from the three dimensions down to two dimensions by using the UMAP algorithm [57]. Within recent years the field of dimensionality reduction algorithms has grown a lot from just the typical (linear) principal component analysis to also include non-linear algorithms. The t-SNE algorithm [78] deserves an honorable mention since this algorithm revolution-

ized the usage of (nonlinear) dimensionality reduction algorithms in e.g. bioinformatics [76, 82] yet its mathematical foundation has strongly been improved with the never, faster UMAP algorithm [57] which usage is also expanding [20, 21, 34].

The aim of UMAP, short for Uniform Manifold Approximation and Projection, is to correctly identify and preserve the structure, or topology, of the high-dimensional feature space in a lower-dimensional output space. It does so by trying to stitch together local manifolds in the high-dimensional feature space such that the difference between the high- and low-dimensional representations is minimized according to the cross-entropy such that both global structure and local structure is preserved [57]. Compared to t-SNE the approach in UMAP has an algebraic topological background compared to the more heuristic approach taken by t-SNE. Note that the UMAP algorithm is not provided any information about which jets are which types.

The UMAP algorithm has several hyperparameters, where two of the most important ones are the number of neighbors `n_neighbors` which controls the priority between correctly preserving the global versus the local structure, and the `min_dist` which defines how tightly together UMAP is allowed to cluster the points in the low-dimensional representation. To properly choose the best combination of `n_neighbors` and `min_dist` a grid search with `n_neighbors`  $\in \{10, 50, 100, 250\}$  and `min_dist`  $\in \{0, 0.2, 0.5\}$  is performed. This is shown for 4-jet events in Figure B.1 in the appendix. In this case the choice of best combination of `n_neighbors` and `min_dist` is subjective at best, but it was judged by the author that `n_neighbors` = 250 and `min_dist` 0.2 gave the best compromise between preserving local and global structure. The results of running UMAP on 4-jet events can be seen in Figure 5.4. Here the millions of points are plotted using Datashader [8] to avoid overplotting and colored according to the jet type. From the figure it is seen how there are some clear, blue *b*-jet clusters, however, most of the data seem to be a mix of *g*- and *uds*-jets. The plots with the same UMAP parameters for 3-jet and 2-jet events are seen in Figure 5.3 and 5.4.

These figures suggests that it should be possible to discriminate the *b*-jets from the other jets somewhat, however, no clear separation is expected. The t-SNE algorithm was also tested but showed inferior performance compared to UMAP, see Figure B.2 in the appendix for an example of this.

The correlation between the vertex variables can be seen in Figure 5.5, where the upper diagonal shows the linear correlation  $\rho$  and the lower diagonal shows the (estimate of the) MIC non-linear correlation  $MIC_e$ . Here it can be seen that `projet` and `bqvjet` correlate mostly whereas the other variables correlate a lot less. Had they all correlated a lot, it would be more difficult to extract any meaningful insights from the system as it would contain less information.

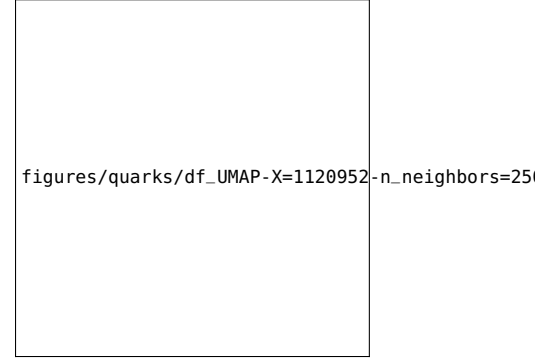


Figure 5.2: Visualization of the vertex variables for the different categories: *true b-jets* in blue, *true c-jets* in red, *true uds-jets* in green, *true g-jets* in orange, and *non q-matched* events in purple. The clustering is performed with the UMAP algorithm which outputs a 2D-projection. This projection is then visualized using the Datashader which takes care of point size, avoids over and underplotting, and color intensity.

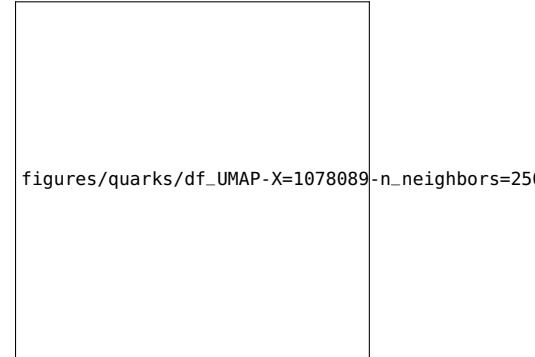


Figure 5.3: UMAP visualization of vertex variables for 3-jet events.

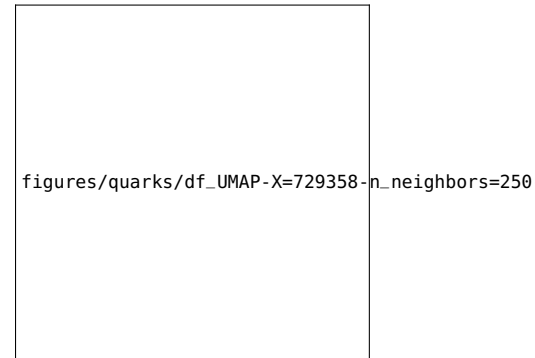


Figure 5.4: UMAP visualization of vertex variables for 2-jet events.

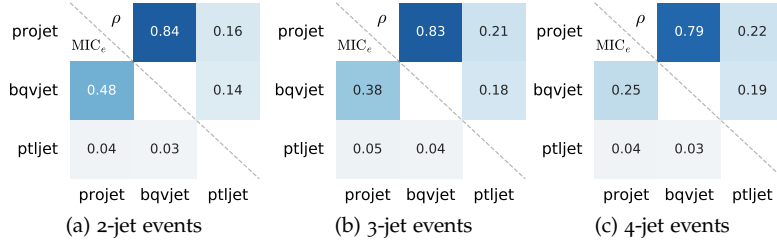


Figure 5.5: Correlation of the three vertex variables for 2-, 3- and 4-jet events.

### 5.3 Loss and Evaluation Function

In contrary to the housing prices subproject the goal in this project is to predict the class of particles, or the types of jets, where the so-called *signal* observations<sup>4</sup> are often assigned the label 1 and *background* observations 0. The combination of this being a *classification* problem (compared to a regression problem) along with the fact all the variables are actual measurements from a particle physics accelerator means that the issue of outliers is negligible. This also means that the problem of finding a robust loss function is non-existent since the in classification loss is already bounded in the  $[0, 1]$ -interval.

Classically *accuracy* is often used as loss function for classification which is simply the fraction of correct predictions, however, accuracy as a metric suffers a lot when handling *imbalanced* data: when the ratio between the number of instances of each class is not approximately (50 : 50)%. The problem is that if the sample contains 90 % background and only 10 % signal, then a simple model which simply predicts everything to be background will have a 90 % accuracy.

To circumvent this issue, the area under the ROC curve (AUC) is used, where the ROC<sup>5</sup> curve is the the *signal efficiency*  $\varepsilon_{\text{sig}}$  of the ML model plotted as a function of the *background efficiency*  $\varepsilon_{\text{bkg}}$ . The definition of these two measures are:

$$\varepsilon_{\text{sig}} = \frac{S_{\text{sel}}}{S_{\text{tot}}}, \quad \varepsilon_{\text{bkg}} = \frac{B_{\text{sel}}}{B_{\text{tot}}}, \quad (5.1)$$

where  $S_{\text{sel}}$  are signal events that were also selected (predicted) as signal by the ML model,  $S_{\text{tot}} = S_{\text{sel}} + S_{\text{rej}}$  is the total number of signal events (the selected and rejected), and likewise for background events  $B$ . Within the machine learning community the signal efficiency is called the true positive rate (TPR) and the background efficiency the false positive rate (FPR). For the rest of this project, the AUC will be the evaluation function  $f_{\text{eval}} = \text{AUC}$ , however, since this metric does not work on single observations it cannot be used as the loss function. Instead we will use the *log-loss* as the loss function<sup>6</sup> which not only is differentiable for single predictions, compared to AUC, but also takes the certainty of the prediction into account. When using tree-based algorithms or neural networks one can extract not only whether or not a single observation is classified as signal or background but also a prediction score. This is a number in the  $[0, 1]$ -interval and the closer to 1 the score is, the

<sup>4</sup> Often called *signal events*, however, this term would require that each event constitutes a single data point in the dataset which it does not here.

<sup>5</sup> Receiver Operating Characteristic.

<sup>6</sup> In the context of machine learning this is the same as the *cross entropy*.

more certain the model is of the prediction being signal. Given the prediction score  $\hat{y}$  and the true label  $y$ , the log-loss  $\ell_{\log}$  is calculated as:

$$\ell_{\log} = -y \log \hat{y} - (1 - y) \log (1 - \hat{y}). \quad (5.2)$$

This is visualized in Figure 5.6. Here it can be seen how the loss changes as a function of the prediction score. Notice that when  $y = 0$  the loss for  $\hat{y} = 1$  diverges towards  $\infty$  and likewise with  $y = 1$  and  $\hat{y} = 0$  (since  $\log 0$  diverges to  $-\infty$ ).

#### 5.4 *b*-Tagging Analysis

The ability to discriminate between the different types of particles produced in a collision is obviously import to understand the results. Today much work go into tagging algorithms from *b*-tagging in ATLAS and CMS [69] but this work started even decades ago. That *b*-quarks are tagged specifically is both due to *b*-quarks having more unique characteristics compared to e.g. *c*-quarks and are thus easier to tag, but also the fact that *b*-quarks are the second-heaviest of the quarks and are measured to better understand CP<sup>7</sup>-violation at LHC-b, contributes to the choice of tagging *b*-quarks. In ALEPH Proriot et al. [63] started the work of comparing different methods for *b*-tagging already in 1991. They concluded that a neural network had the best performance compared to e.g. a linear (Fisher) discriminant. The neural network used was a 3-layer neural network (NN) trained on nine variables and the output `nmbjet`. For this of this project this pre-trained network will be called NNB.

The data are split<sup>8</sup> into training and test sets in such a way that the individual jets in an event are not split. As such, the splitting is performed at event-scale in a (80 : 20)% train-test ratio.

##### 5.4.1 *b*-Tagging Hyperparameter Optimization

Compared to the housing prices dataset, the number of observations  $N$  is a lot larger, although the dimensionality  $M$  is much smaller ( $3 \ll 143$ ). Therefore both XGBoost (XGB) and LightGBM (LGB) were included as models initially since their performance in the housing dataset was very similar but LightGBM was expected to quite a lot faster on this dataset, which also turned out to be the case. The models were hyperparameter optimized (HPO) using random search (RS) since the Bayesian optimization (BO) did not show any performance gains compared to RS. They were run with 5-fold cross validation and early stopping with a patience of 100. The PDFs for the random search for the LightGBM model can be seen in Table 5.4, and the ones for XGBoost in Table B.2 in the appendix. The random search has been run with 100 iterations for LightGBM and only 10 for XGBoost since XGBoost is slow at fitting datasets of this size<sup>9</sup>. The results of the HPO for 3-jet and 4-jet events can be seen in Figure 5.7. For 3-jets it can be seen how most of the iterations share about the same performance within  $1\sigma$ , however some

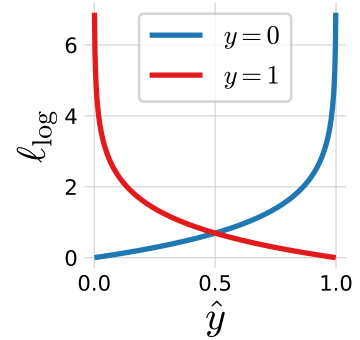


Figure 5.6: Plot of the log-loss  $\ell_{\log}$ .

<sup>7</sup> Short for charge-parity.

<sup>8</sup> After removing all low-energy jets such that all events that contain any jets with an energy of less than 4 GeV are removed.

Hyperparameter	Range
<code>subsample</code>	$\mathcal{U}(0.4, 1)$
<code>colsample_bytree</code>	$\mathcal{U}_{\text{trunc}}(0.4, 1, 2)$
<code>max_depth</code>	$\mathcal{U}_{\text{int}}(-5, 63)$
<code>num_leaves</code>	$\mathcal{U}_{\text{int}}(7, 4095)$

Table 5.4: Probability Density Functions for the random search hyperparameter optimization process for the LightGBM model. For an explanation of  $\mathcal{U}_{\text{trunc}}$ , see section 5.5. All negative values of `max_depth` are interpreted as no max depth by both LGB and XGB.

<sup>9</sup> See page 70 for a discussion of the timings.



iterations have a significantly decrease in performance. For 4-jets there are not any iterations which share the same bad performance relative to the others as some of the 3-jets.

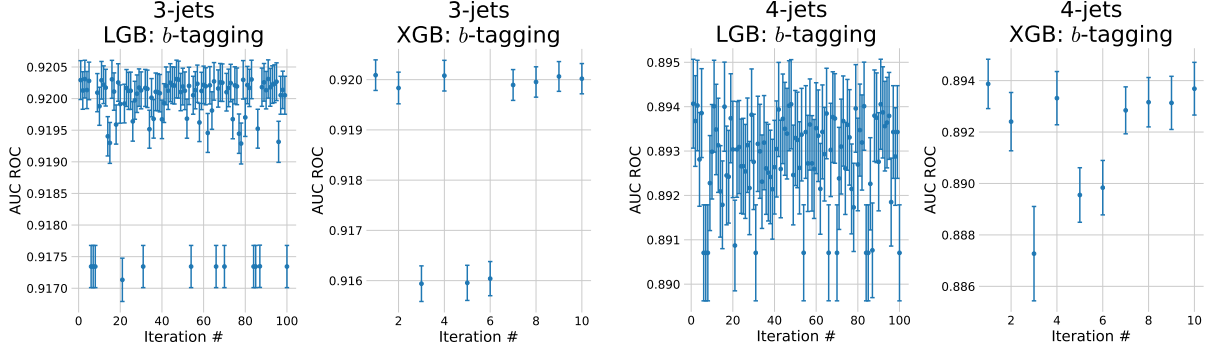


Figure 5.7: Hyperparameter Optimization results of  $b$ -tagging with random search. From left to right, we have A) 100 iterations of RS with LGB on 3-jets, B) 10 iterations of RS with XGB on 3-jets, C) 100 iterations of RS with LGB on 4-jets, D) 10 iterations of RS with XGB on 4-jets. Notice the different ranges on the y-axes.

The relationship between the different hyperparameters in 4-jet events can be seen in the parallel coordinate plot in Figure 5.8. First of all the importance of the column downsampling `colsample_bytree` variable is significant: all of the low-performing hyperparameter sets have a low value of this hyperparameter. Since  $M = 3$  for the vertex variables this makes logical sense; using only  $\text{int}(\sim 0.5 \cdot 3) = 1$  variable<sup>10</sup> the model cannot properly learn the structure in the data. Compared to the column downsampling, the other hyperparameters are notably less important. The same overall conclusion can be inferred in the 3-jet case, see Figure B.3 in the appendix.

<sup>10</sup> See section 5.5 for a deeper discussion about the `colsample_bytree` hyperparameter.

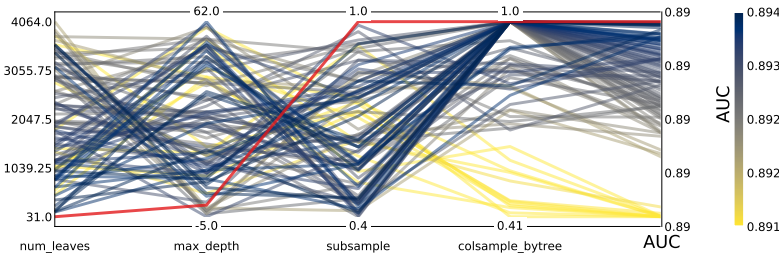


Figure 5.8: Hyperparameter optimization results of  $b$ -tagging for 4-jet events. The results are shown as parallel coordinates with each hyperparameter along the  $x$ -axis and the value of that parameter on the  $y$ -axis. Each line is an event in the 4-dimensional space colored according to the performance of that hyperparameter as measured by AUC from highest AUC in dark blue to lowest AUC in yellow. The single best hyperparameter is shown in red.

#### 5.4.2 $b$ -Tagging Results

The prediction score for the  $b$ -tagging models is usually called the  $b$ -tag and will be written as  $\beta_{\text{tag}}$ . The distribution of  $\beta_{\text{tag}}$  for the two HPO-optimized models, LGB and XGB, together with the pre-trained neural network NNB can be seen in Figure 5.9 for 4-jet events and in B.4 in the appendix for 3-jet events. Notice the strong match between the NNB and LGB models. The XGB model has almost no high  $b$ -tags  $\beta_{\text{tag}} > 0.8$ , but a majority of  $b$ -tags in the very low end. This indicates that the XGBoost has focussed on the background events compared to the signal events, whereas the NNB and LGB models have focused more on the signal events.

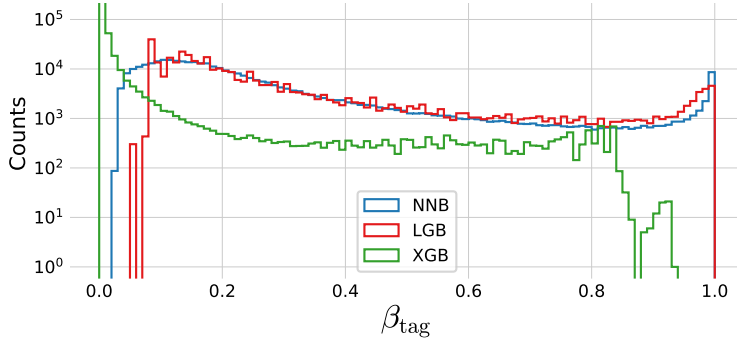


Figure 5.9: Histogram of  $b$ -tag scores  $\beta_{\text{tag}}$  in 4-jet events for **NNB** (the neural network pre-trained by ALEPH, also called `nnbjet`) in blue, **LGB** in red, and **XGB** in green.

Even though the distributions of  $b$ -tags are different between the three models, the real performance plot for classification is the ROC curve seen in Figure 5.10 for 4-jet events. Here the signal efficiency  $\epsilon_{\text{sig}}$  is plotted as a function of the background efficiency  $\epsilon_{\text{bkg}}$  with the AUC shown in the bottom right corner. The LGB and XGB models performs similarly well with an  $\text{AUC} = 0.896$  compared to the NNB with  $\text{AUC} = 0.884$ . The differences between the models are even smaller for 3-jet events seen in Figure B.5 in the appendix. In general the LGB and XGB models are so similar that they cannot be distinguished from another in any of the plots and their difference in AUC is on the fourth decimal point. However, the LGB model is several times faster than the XGB model. In comparison, 10 iterations of HPO using RS on 3-jet events with XGB took more almost 34 hours on HEP<sup>11</sup> compared to just 23 hours for 100 iterations for LGB. The same performance difference was seen in 4-jet events where the timings were 4 hours for XGB compared to 2.5 hours for LGB, and thus XGB is dropped in all subsequent analysis.

<sup>11</sup> The local computing cluster.

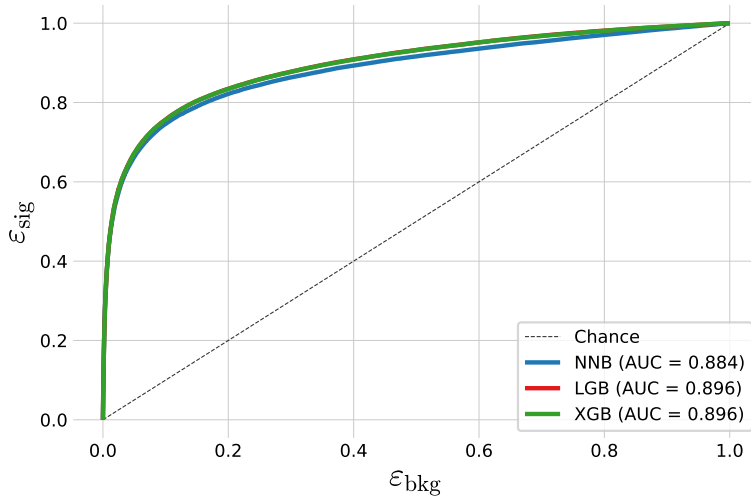


Figure 5.10: ROC curve of the three  $b$ -tag models in 4-jet events for **NNB** (the pre-trained neural network trained by ALEPH, also called `nnbjet`) in blue, **LGB** in red, and **XGB** in green. In the legend the area under curve (AUC) is also shown. Notice that the LGB and XGB models share performance and it is thus due to overplotting that only the green line for XGB can be seen. In the machine learning community the background efficiency  $\epsilon_{\text{bkg}}$  is sometimes known as the false positive rate (FPR) and the signal efficiency  $\epsilon_{\text{sig}}$  as the true positive rate (TPR).

The distribution of  $b$ -tag scores  $\beta_{\text{tag}}$  from the  $b$ -tag LGB model for 4-jet events can be seen in Figure 5.11. In the figure it can be seen how the separation between the heavier quarks and light quarks (and gluons) is clear at high values of  $\beta_{\text{tag}}$ , however, a lot of  $c$ -quarks also get a high  $b$ -tag score. The same is seen for 3-jet

events in Figure B.6 in the appendix.

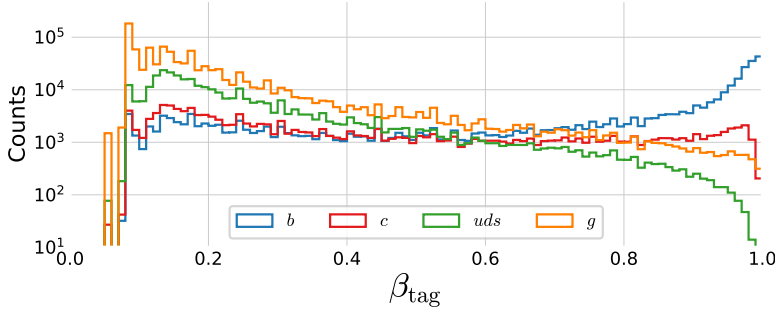


Figure 5.11: Distribution of  $b$ -tags in 4-jet events for  $b$ -jets in blue,  $c$ -jets in red,  $uds$  in green and  $g$  in orange.

#### 5.4.3 $b$ -Tagging Model Inspection

To get a better understanding of the trained LGB model, the global SHAP feature importances can be seen in Figure 5.12 for 4-jet events. First of all it is noted that the `projet` has global feature importance of 57.32 %, `bqvjet` 29.16 %, and `ptljet` 13.52 %. For all three variables it is seen how most of the points have many small feature values which has a negative impact on the model output however small. Especially the `ptljet` has many features with a low value (0 in fact) yet this does not pull the model too much towards background events compared to if a jet has a high value of `ptljet` which has a strong, positive impact on the output prediction.

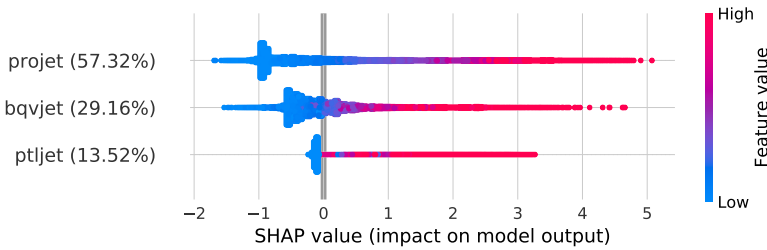


Figure 5.12: Global feature importances for the LGB  $b$ -tagging algorithm on 4-jet events. The normalized feature importance is shown in the parenthesis and the each dot is an observation showing the dependance between the SHAP value and the feature's value.

In regression, the model output is a continuous prediction  $\hat{y}_{\text{reg}} \in \mathbb{R}$ . In classification what is actually happening under the hood is that the model predicts a value  $\tilde{y} \in \mathbb{R}$  which is transformed to a number in the  $[0, 1]$ -interval via the *expit* function:

$$\text{expit}(\tilde{y}) = \frac{e^{\tilde{y}}}{1 + e^{\tilde{y}}} \equiv p, \quad (5.3)$$

where  $p$  is a number in the  $[0, 1]$ -interval. The *expit* function is also sometimes known as the logistic function and is visualized in Figure 5.13. Its inverse is the *logit* function:

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = \tilde{y}, \quad (5.4)$$

which is visualized in Figure 5.14. The fraction in equation (5.4) is called the *odds* and the logit-transformed value of  $p$ ,  $\text{logit}(p) = \tilde{y}$ ,

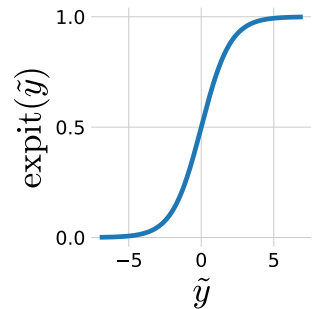


Figure 5.13: The *expit* function.

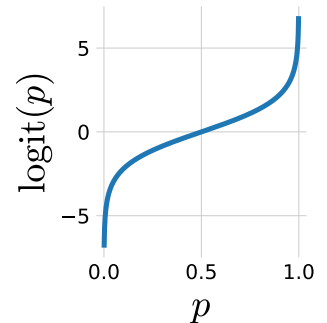
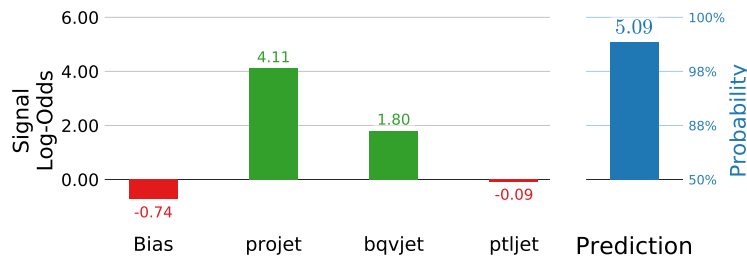


Figure 5.14: The *logit* function.



is thus sometimes called the *log-odds*. It is in this log-odds space that LightGBM makes its predictions and the SHAP values in Figure 5.12 are also in log-odds space. The additivity<sup>12</sup> of SHAP is in this log-odds space.

With this in mind, single predictions of the LGB *b*-tagging model can be understood with SHAP which Figure 5.15 is an example of. This figure shows the logic behind the model's prediction for this particular jet. That the bias is negative reflects that there is a majority of background compared to signal<sup>13</sup>. This particular event has `projct` = 1.003, `bqvjet` = 0.529, and `ptljjet` = 0. In the plot it is seen how this high value of `projct` has the greatest impact on the model prediction, while the medium value of `bqvjet` also pushes the model prediction towards a signal-prediction. The four bars in the left part of the plot are all in log-odds space and their sum is shown as the blue bar to the right, where the right *y*-axis shows the value in probability space  $p \in [0, 1]$ . This jet was in fact a *b*-jet.



<sup>12</sup> See also section 2.8.

<sup>13</sup> There are 22.1 % *b*-jets in the 3-jet training set.

Figure 5.15: Model explanation for the 3-jet *b*-tagging LGB model for a *b*-like jet. The first column is the bias of the training set which acts as the naive prediction baseline, the rest are the input data variables. On the right hand side of the plot is the model prediction shown. The left part of the plot is shown in log-odds space, the right part in probability space. The negative log-odd values are shown in red, positive ones in green, and the prediction value in blue.

## 5.5 Truncated Uniform PDF

Initially when plotting the HPO performance as a function of iteration, it was seen how there were three significant plateaus, where the highest plateau (i.e. highest AUC value and thus best score) was only seen in the very first iteration. It was quickly realized that this was due to the very first iteration was being run with the default values of the LGB and XGB models in the custom implementation by the author. However, what was not understood was why this value was performing so much better than random sets of hyperparameters<sup>14</sup>. During the debugging process the column downsampling `colsample_bytree` was diagnosed to be the culprit. The default value is `colsample_bytree` = 1, however, the probability density function (PDF) used in RS for this parameter was  $\mathcal{U}(0.4, 1)$  which was expected to give the same performance as the default value for large values of `colsample_bytree`. By inspecting the source code of LightGBM it was realized that the model takes the integer of the column downsampling multiplied with the total number of features if the column downsampling is less than 1 [6]. This means that no matter how close to 1 the column downsampling get, the integer value of the total number of columns get floored to 2 at max, compared to when the column downsampling is exactly 1 which it only is for the default values.

<sup>14</sup> LightGBM and XGBoost of course have chosen their default parameters smartly, however, it one would not expect them to outperform other sets of hyperparameters that clearly.

To deal with this problem a new PDF was developed on top of the existing ones in Scipy, the truncated uniform PDF:  $\mathcal{U}_{\text{trunc}}(a, b, c)$ . This PDF first generates a random number  $x$  from a uniform distribution between  $a$  and  $c$ . Then if  $x$  is larger than  $b$  it is floored to  $b$ . In this way, it is possible to both get values of  $x$  in the interval  $[a, b]$  but also values exactly equal to  $b$ . The value of  $c$  controls how often these “overflow” values of  $x$  are generated.

## 5.6 *g-Tagging Analysis*

The trained  $b$ -tagging LGB model is a jet-based model which provides a  $b$ -tag score  $\beta_{\text{tag}}$  to a jet. This also means that each of the jets in e.g. a 4-jet event can get a  $b$ -tag:  $\beta_{\text{tag}} = [\beta_{\text{tag}_1}, \beta_{\text{tag}_2}, \beta_{\text{tag}_3}, \beta_{\text{tag}_4}]$ . Using  $\beta_{\text{tag}}$  one can train a new model on the events, compared to individual jets, where signal events are defined to be  $q$ -matched events<sup>15</sup> where the non- $q$ -matched jets are assigned the  $n - 2$  lowest  $b$ -tag scores for  $n$ -jet events; e.g.  $\beta_{\text{tag}} = [0.95, 0.89, 0.15, 0.07]^\top$  for the four jets  $[b, \bar{b}, g, g]$ . This event-based process will be called  $g$ -tagging and the trained model will return a  $g$ -tag score written as  $\gamma_{\text{tag}}$ . Compared to the  $b$ -tagging LGB model, this model will allow one to extract entire events which contains a clear identification of gluons versus non-gluons.

### 5.6.1 *Permutation Invariance*

Since the  $b$ -tags are only based on the vertex variables, the goal of the  $g$ -tag is to also be constructed in an un-biased way with respect to the jet energy  $E_{\text{jet}}$ . However, even though  $\beta_{\text{tag}} \perp\!\!\!\perp E_{\text{jet}}$  and  $\gamma_{\text{tag}} = f(\beta_{\text{tag}})$ , it turned out that  $\gamma_{\text{tag}} \not\perp\!\!\!\perp E_{\text{jet}}$ , where  $a \perp\!\!\!\perp b$  is defined to mean that  $a$  is independent<sup>16</sup> of  $b$  and  $f$  is an unknown function. This was because the ordering of the jets within the event was energy-dependent: they sorted according to their  $E_{\text{jet}}$ .

This meant that the different components in  $\beta_{\text{tag}}$  had different importances, even though they should be equally important. Instead of defining  $\beta_{\text{tag}}$  as a vector it should instead be seen as a set<sup>17</sup>  $\beta_{\text{tag}} = \{\beta_{\text{tag}_1}, \dots, \beta_{\text{tag}_n}\}$ . The  $g$ -tagging model trained on the events should thus be *permutation invariant*<sup>18</sup> with regards to the input variables. The category of permutation invariant (and equivariant<sup>19</sup>) neural networks in the deep learning community has seen an huge development within recent years where the paper from Zaheer et al. [85] in 2017 was quite influential, however also other examples exists [65, 41]. Yet, the same development cannot be said to have happened within the more classic machine learning field.

Although not being a novel software-technical solution, the problem was circumvented by two simple, different approaches: 1) by simply shuffling the inputs variables independently for each observation (row) in the dataset, and 2) training on all possible permutations of the variables in the dataset. The second approach can be seen as a feature augmentation technique where the data is

<sup>15</sup> Remember that  $q$ -matched events are events with one, and only one, jet that is  $q$ -matched to one of the quark-jets, and one, and only one of the jets is  $q$ -matched to the other quark-jet.

<sup>16</sup> And  $\not\perp\!\!\!\perp$  means not independent.

<sup>17</sup> Since sets have no inherent order.

<sup>18</sup>  $f(\mathbf{x}) = f(\tau(\mathbf{x}))$  for any permutation  $\tau$  on an input vector  $\mathbf{x}$ .

<sup>19</sup>  $\tau(f(\mathbf{x})) = f(\tau(\mathbf{x}))$  for any permutation  $\tau$  on an input vector  $\mathbf{x}$ .

artificially increased with factor of  $n$  factorial:  $N \rightarrow n! \cdot N$  where  $N$  is the number of observations (rows) and  $n$  is the number of jets. These two methods were tested along with the original order of the dataset.

### 5.6.2 $g$ -Tagging Hyperparameter Optimization

Four LightGBM models, two for 3-jet events and two for 4-jet events, were trained and hyperparameter optimized for both the energy ordered and shuffled<sup>20</sup> data sets with 100 iterations of random search with the same PDFs as for the  $b$ -tagging, see Table 5.4, and 5-fold cross validation and early stopping with a patience of 100. The results of the HPO can be seen in Figure 5.16. Here the two 3-jets models are seen in the two plots to the left, and the two 4-jets to the right. The very left plot shows the 3-jet energy-ordered (no permutation) performance as a function of iteration number, which was also where the issues mentioned in section 5.5 were first discovered. Here the difference between the how many of the three variables, the three  $b$ -tags, are included is seen as three clear plateaus. The three plateaus are also seen in the 3-jet events that were shuffled, however, with more variation in each plateau, along with a drop in performance. For the 4-jet events the plateaus are not as apparent but it can still be seen how some of the iterations how a significantly lower score than others. The parallel plots for the four fits can be seen in Figure B.8–B.11 in the appendix.

<sup>20</sup> The method with all permutations was trained using the same hyperparameters as the best ones for the shuffled model.

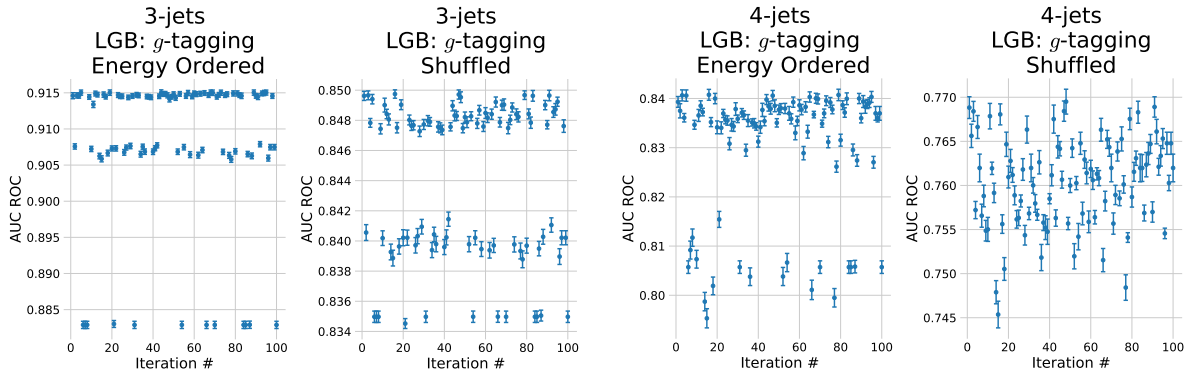


Figure 5.16: Hyperparameter Optimization results of  $g$ -tagging with 100 iterations of random search with LGB. From left to right, we have A) 3-jet events energy-ordered (no permutations), B) 3-jet events row-shuffled, C) 4-jet events energy-ordered, D) 4-jet events row-shuffled. Notice the different ranges on the y-axes.

### 5.6.3 PermNet

In addition to the LGB models, a permutation invariant neural network called PermNet based on the Deep Sets paper [85] implemented in Tensorflow [10] by Faye [37] was also tested. Zaheer et al. [85] showed that  $f(X)$  is permutation invariant if and only if it can be decomposed in the following way:

$$f(X) = \rho \left( \sum_{x \in X} \phi(x) \right). \quad (5.5)$$

for suitable transformations  $\rho$  and  $\phi$  (which the neural network learns<sup>21</sup>). The PermNet was trained using three layers<sup>22</sup> with leaky ReLU [56] as the activation function and ADAM [51] as the optimizer optimizing the log-loss. The network was trained with early stopping with a patience of 50 epochs and a batch size of 128. A visual overview of the PermNet architecture can be seen in Figure B.12 in the appendix.

<sup>21</sup> This is possible since neural networks are universal function approximators [45].

<sup>22</sup> Where the two hidden layers have 128 and 64 neurons in each.

#### 5.6.4 1D Comparison of LGB and PermNet

To better understand the difference between the difference between the LGB and PermNet models, a small comparison was made. This comparison was constructed by summing the  $b$ -tag scores in the  $n$ -jet event together  $\sum_i^n \beta_{\text{tag}_i}$ . The  $\beta_{\text{tag}_i}$  are summed together since this turns the problem into a 1D problem that is easy to visualize, the sum of numbers is a permutation invariant function, and is similar to the simplest functions of  $\rho$  and  $\phi$  in equation (5.5): the identity function. The 1D models are fit to the training events and then a linear scan from  $\sum_i^n \beta_{\text{tag}_i} = 0.4$  to 3.1 is made to see how the predicted  $g$ -tags  $\gamma_{\text{tag}}$  distribute. This is shown in Figure 5.17 for 4-jet events. Here the value of  $\gamma_{\text{tag}}$  is shown for the two models together with the fraction of signal to background in each bin. If the  $g$ -tag score should resemble a true probability it would be expected to follow the signal ratio, e.g. a model should predict  $\gamma_{\text{tag}} = 0.9$  if there is 90 % signal in that bin. In the figure it is seen how the PermNet does a great job at fitting the signal fraction, however, the LGB model also does a decent job. Remember that none of these models were shown the signal fraction explicitly, only the  $b$ -tag sum and a signal-or-background label. The distribution of signal and background together with the distribution of cuts made by the LGB model can be seen in Figure B.13 in the appendix. The similar plots for 3-jet events are plotted in Figure B.14 and B.15, both in the appendix.

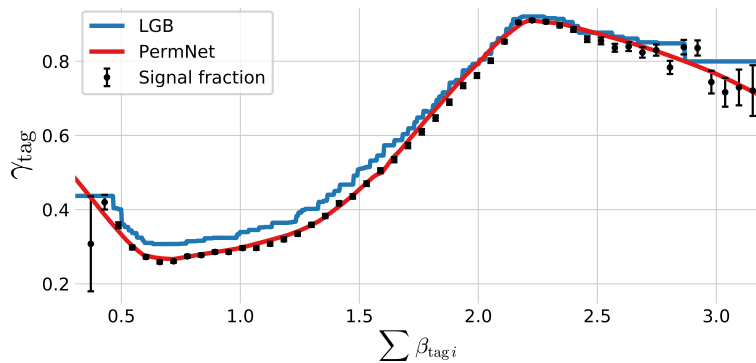


Figure 5.17: Plot of the (1D)  $g$ -tag scores for 4-jet events as a function of  $\sum \beta_i$  for the LGB model in blue and the PermNet model in red. The signal fraction (based on the signal and background histograms in Figure B.13) is plotted as black error bars where the size of the error bars is based on the propagated uncertainties of the signal and background histogram assuming Poissonian statistics.

It can be concluded, at least in 1D, that both LGB and PermNet are able to capture the inherent structure in the (1D) data. First of all it is seen that the two 1D models follow each other relatively close and only predicts  $\gamma_{\text{tag}}$ s in a quite limited range. The three other PermNet curves follow each other in such an extent that it

is almost difficult to separated them, which is also expected since they should not be able to distinguish between the energy ordered and the shuffled events. The LGB models for the shuffled and all-permuted events

### 5.6.5 $g$ -Tagging Results

The distribution of  $g$ -tag scores in 4-jet (training) events the can be seen in Figure 5.18 for the eight combinations of the two models (LGB and PermNet) and the four data sorting methods (energy ordered, (row) shuffled, all permutations, and the (1D) sum.). At first the increased number of events (a factor of 24 for 4-jet events) with the all-permutation scheme is seen separating the two light green curves from the rest. The energy ordered LGB model is the combination which utilizes most of the  $\gamma_{\text{tag}}$ -range, while the two 1D sum models have the most limited range, indicating that the models are more uncertain about their predictions. The energy ordered and shuffled PermNet models can more or less only be distinguished because the latter is plotted with dashed lines. This makes sense, since they are also expected to make the same predictions were they really permutation invariant<sup>23</sup>. When plotted with normalized counts it is seen how the shuffled and all-permuted LGB models also follow each other very closely, which can still be partly seen in this plot by comparing the two distributions. The distribution of  $g$ -tags in 3-jet training events can be seen in Figure 5.18 in the appendix.

<sup>23</sup> It is only because of the stochasticity in the optimization process of the two networks that they did not converge to the completely same predictions.

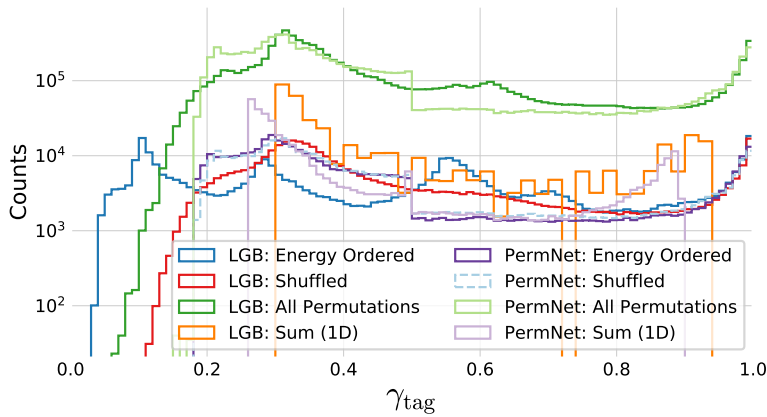


Figure 5.18: Distribution of  $g$ -tag scores in 4-jet events shown with a logarithmic  $y$ -scale for LGB: Energy Ordered in blue, LGB: Shuffled in red, LGB: All Permutations in green, LGB: Sum 1D in orange, PermNet: Energy Ordered in purple, PermNet: Shuffled in light-blue, PermNet: All Permutations in light-green, PermNet: Sum 1D in light-purple. Here LGB and PermNet are the two different type of models and “Energy Ordered”, “Shuffled”, “All Permutations”, and “Sum 1D” are the different methods used for making the input data permutation invariant (except energy ordered).

To see the performance of the different combinations, see the ROC curve in Figure 5.18 which shows the performance on 4-jet events with the AUC shown in the legend. First of all it is easy to see that the energy ordered LGB model is significantly higher-performing than the rest of the models, however, this model is also energy-biased by not being permutation invariant in the  $b$ -tags and is only included to see how large a performance drop the permutation invariance criterion causes. The worst performing models are the two 1D sum models, as expected since they only have a single dimension to learn from, compared to the four dimensions

that the other models have. In general it can be seen that the rest of the models are performing almost identically, with the LGB model trained on all permutations to be the highest-performing of them all by a small margin. For 3-jet events the same overall picture is seen, see Figure 5.18 in the appendix, however, here LGB model trained on the shuffled events performs the best, yet this performance improvement is so small compared to the all-permutations LGB model that it is expected to be due to statistical fluctuations and not a real performance difference.

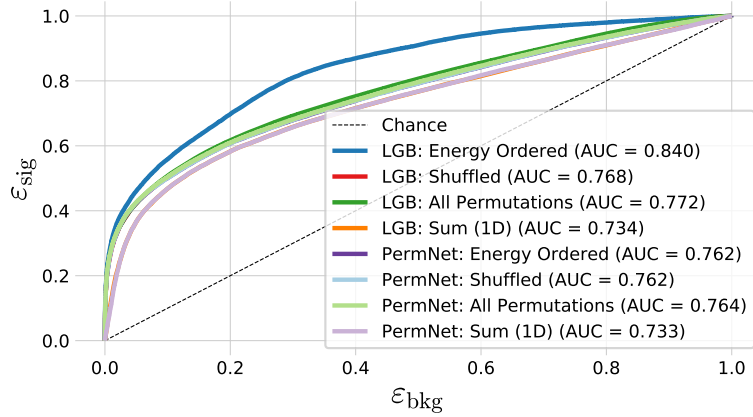


Figure 5.19: ROC curve of the eight  $g$ -tag models in 4-jet. First one in dashed black is the ROC curve that you get by random chance. The colors are the same as in Figure 5.18 and in the legend also the Area Under the ROC curve (AUC) is shown.

Based on the AUC scores seen in the ROC curves in Figure B.17 and B.18, the LGB-model trained on all permutations will be the  $g$ -tagging model choice. To see how this model's predictions of  $\gamma_{\text{tag}}$  distributive for different particle-types for signal and background events see Figure 5.20. Here the distribution of  $\gamma_{\text{tag}}$  is shown for 4-jet signal events and background events. Remember that in  $g$ -tagging, the signal events are defined as events where the two jets with the highest  $b$ -tags are also the two  $q$ -matched jets (and the entire event is  $q$ -matched). In the figure it can be seen that at high values of  $\gamma_{\text{tag}}$  primarily  $b\bar{b}gg$  events are tagged (signal  $b$ ), but also with some  $c\bar{c}gg$  (signal  $c$ ) and  $bgbg$ -events<sup>24</sup> (background  $b$ ) sorted according to their  $b$ -tags. At low values of  $\gamma_{\text{tag}}$  light quarks ( $uds$ ) dominate.

<sup>24</sup> Or any other permutation of  $b, \bar{b}, g$ ,  $g$  which is not  $b\bar{b}gg$ .

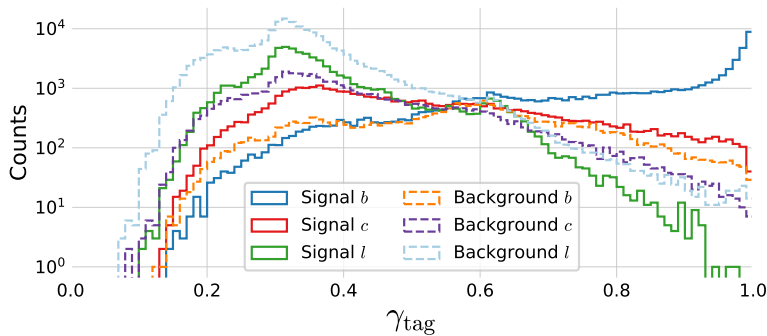


Figure 5.20: Histogram of  $g$ -tag scores from the LGB-model in 4-jet events for  $b$  signal in blue,  $c$  signal in red,  $l$  ( $uds$ ) signal in green,  $b$  background in orange,  $c$  background in purple,  $l$  ( $uds$ ) background in light-blue.

The similar plot for 3-jet events is seen in Figure B.19 in the appendix. This plot has some surprising bumps for mainly  $l$ -quark events which are not yet fully understood. When comparing  $l$ -



quark events in the high- $\gamma_{\text{tag}}$  bump with the ones getting a low  $\gamma_{\text{tag}}$ -value, see Figure 5.21, one can see that  $l$ -quark events with high  $\gamma_{\text{tag}}$  has only two jets with high  $b$ -tags, compared to low- $\gamma_{\text{tag}}$   $l$ -quark events which more often has three jets with high  $b$ -tags.

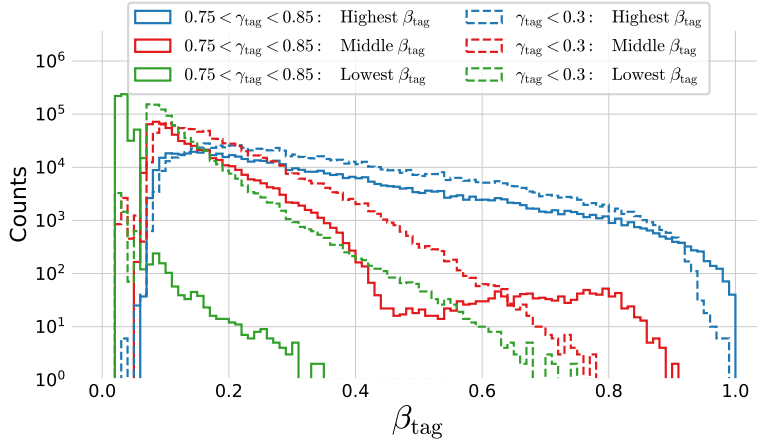


Figure 5.21: Distribution of  $b$ -Tag Scores in 3-Jet  $l$ -Quark Events for low and high  $g$ -tags values. Here  $l$ -quark events with  $0.75 < \gamma_{\text{tag}} < 0.85$ , so the high peak in Figure B.19, are plotted in fully connected lines, and events with  $\gamma_{\text{tag}} < 0.3$  are plotted in dashed lines. For each of these two selection of events the value of the jet with the highest  $\beta_{\text{tag}}$  is shown in blue, the jet with the middle  $\beta_{\text{tag}}$  in red, and the jet with the lowest  $\beta_{\text{tag}}$  in green.

This is even more visible once seen in a 3D scatter plot with the lowest  $\beta_{\text{tag}}$  on the  $x$ -axis, the middle on the  $y$ -axis, and highest on the  $z$ -axis. Three small views from the 3D visualization can be seen in Figure 5.22.

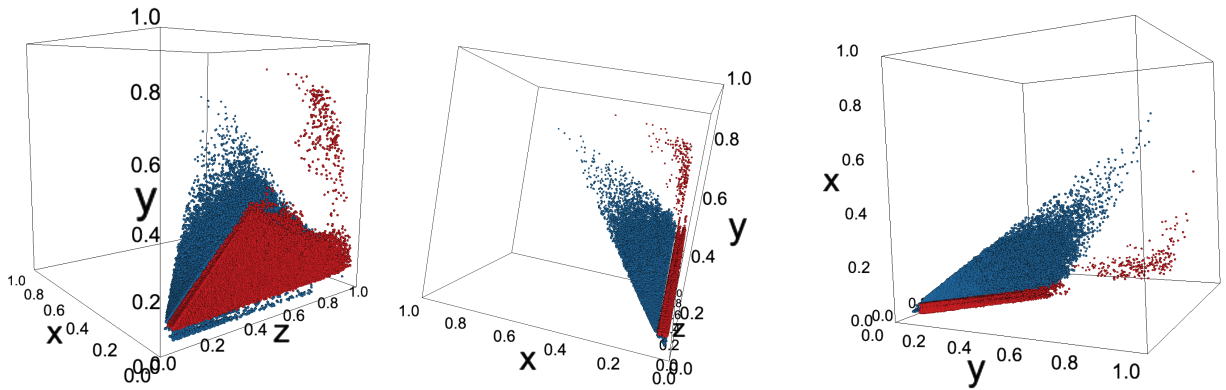


Figure 5.22: XXX.

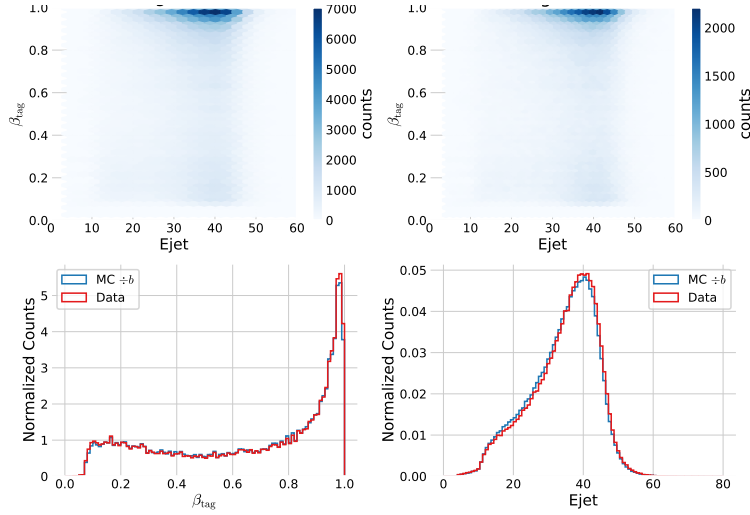


Figure 5.23: Comparison of the b-tag and jet energy (  $E_{\text{jet}}$  ) distributions for Monte Carlo (MC) versus data. In the top row the 2D-distributions are shown for MC on the left (without the extra MCB samples) and data on the right. In the bottom row the 1D marginal distributions are shown for the b-tag and the jet energy with **data** in red and **Monte Carlo** ones in blue. Notice the the almost identical distributions in b-tag.

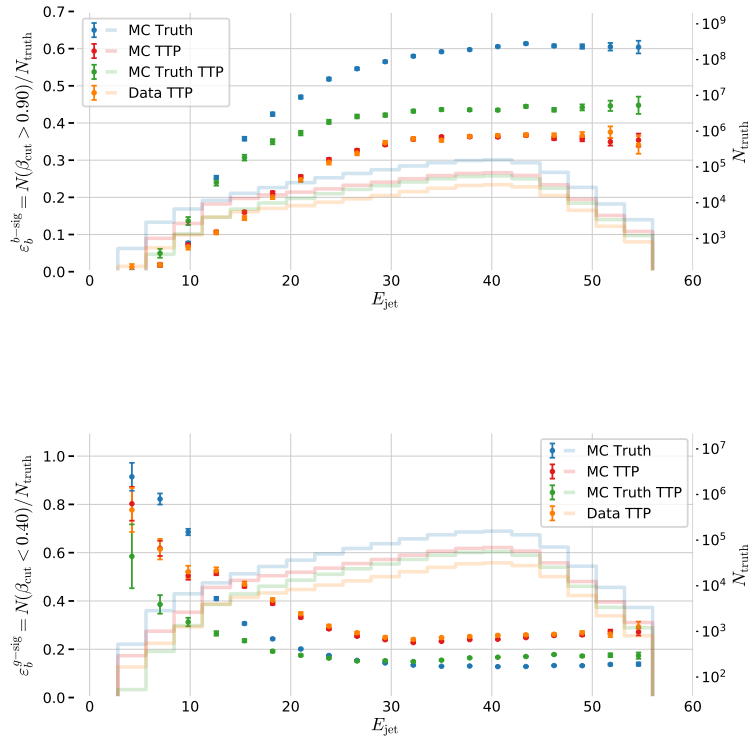


Figure 5.24: Efficiency of the b-tags for b-jets in the b-signal region for 3-jet events,  $\epsilon_b^{b\text{-sig}}$ , as a function of jet energy  $E_{\text{jet}}$ . The b-signal region is defined as  $\beta > 0.9$ . In the plot the efficiencies are shown for **MC Truth** in blue, **MC TTP** in red, **MC Truth TTP** in green, and **Data TTP** in orange. The efficiencies (the errorbars) can be read off on the left y-axis and the counts (histograms) on the right y-axis. The abbreviation TTP is short for “Tag, Tag, Probe” where two jets in a event are used as tags and the probe is then used for further analysis. Notice how both MC TTP and Data TTP follow each other closely.

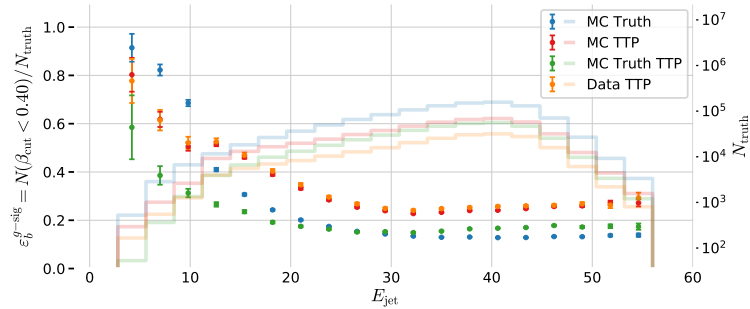


Figure 5.25: Efficiency of the b-tags for b-jets in the g-signal region for 3-jet events,  $\epsilon_b^{g\text{-sig}}$ , as a function of jet energy  $E_{\text{jet}}$ . The g-signal region is defined as  $\beta < 0.4$ . In the plot the efficiencies are shown for **MC Truth** in blue, **MC TTP** in red, **MC Truth TTP** in green, and **Data TTP** in orange. The efficiencies (the errorbars) can be read off on the left y-axis and the counts (histograms) on the right y-axis. The abbreviation TTP is short for “Tag, Tag, Probe” where two jets in a event are used as tags and the probe is then used for further analysis. Notice how both MC TTP and Data TTP follow each other closely.



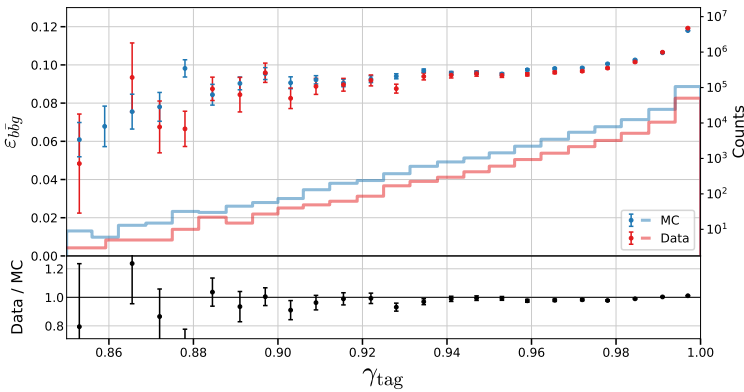
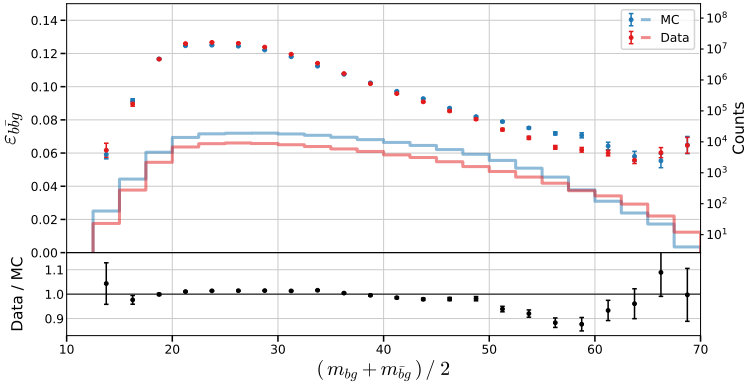
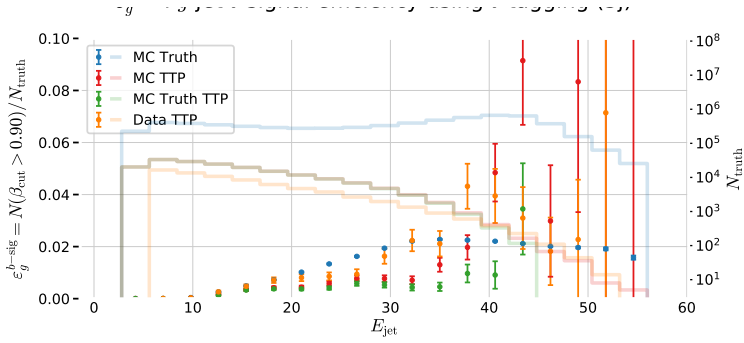
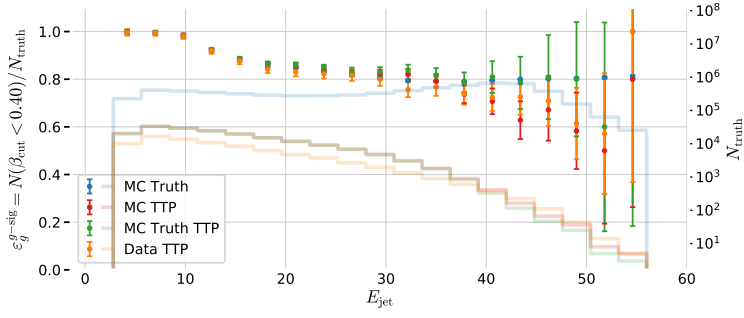


Figure 5.26: Efficiency of the b-tags for g-jets in the g-signal region for 3-jet events,  $\epsilon_g^{g\text{-sig}}$ , as a function of jet energy  $E_{\text{jet}}$ . The g-signal region is defined as  $\beta < 0.4$ . In the plot the efficiencies are shown for MC Truth in blue, MC TTP in red, MC Truth TTP in green, and Data TTP in orange. The efficiencies (the errorbars) can be read off on the left y-axis and the counts (histograms) on the right y-axis. The abbreviation TTP is short for “Tag, Tag, Probe” where two jets in a event are used as tags and the probe is then used for further analysis. Notice how both MC TTP and Data TTP follow each other closely.

Figure 5.27: Efficiency of the b-tags for g-jets in the b-signal region for 3-jet events,  $\epsilon_g^{b\text{-sig}}$ , as a function of jet energy  $E_{\text{jet}}$ . The b-signal region is defined as  $\beta > 0.9$ . In the plot the efficiencies are shown for MC Truth in blue, MC TTP in red, MC Truth TTP in green, and Data TTP in orange. The efficiencies (the errorbars) can be read off on the left y-axis and the counts (histograms) on the right y-axis. The abbreviation TTP is short for “Tag, Tag, Probe” where two jets in a event are used as tags and the probe is then used for further analysis.

Figure 5.28: Proxy efficiency of the g-tags for  $b\bar{b}g$  3-jet events as a function of the mean of the two invariant masses  $m_{b\bar{g}}$  and  $m_{\bar{b}g}$ . The proxy efficiency  $\epsilon_{b\bar{b}g}$  is measured by finding  $b\bar{b}g$ -events where  $\beta_b > 0.9$ ,  $\beta_{\bar{b}} > 0.9$ , and  $\beta_g < 0.4$ . and then calculating  $\epsilon_{b\bar{b}g} = \epsilon_b^{b\text{-sig}} \cdot \epsilon_{\bar{b}}^{b\text{-sig}} \cdot \epsilon_g^{g\text{-sig}}$ . In the top plot  $\epsilon_{b\bar{b}g}$  is shown for MC in blue and Data in red where the counts in each bin can be read on right y-axis. In the bottom plot the ratio between Data and MC is shown.

Figure 5.29: Proxy efficiency of the g-tags for  $b\bar{b}g$  3-jet events as a function of the event's g-tag. The proxy efficiency  $\epsilon_{b\bar{b}g}$  is measured by finding  $b\bar{b}g$ -events where  $\beta_b > 0.9$ ,  $\beta_{\bar{b}} > 0.9$ , and  $\beta_g < 0.4$ . and then calculating  $\epsilon_{b\bar{b}g} = \epsilon_b^{b\text{-sig}} \cdot \epsilon_{\bar{b}}^{b\text{-sig}} \cdot \epsilon_g^{g\text{-sig}}$ . In the top plot  $\epsilon_{b\bar{b}g}$  is shown for MC in blue and Data in red where the counts in each bin can be read on right y-axis. In the bottom plot the ratio between Data and MC is shown.

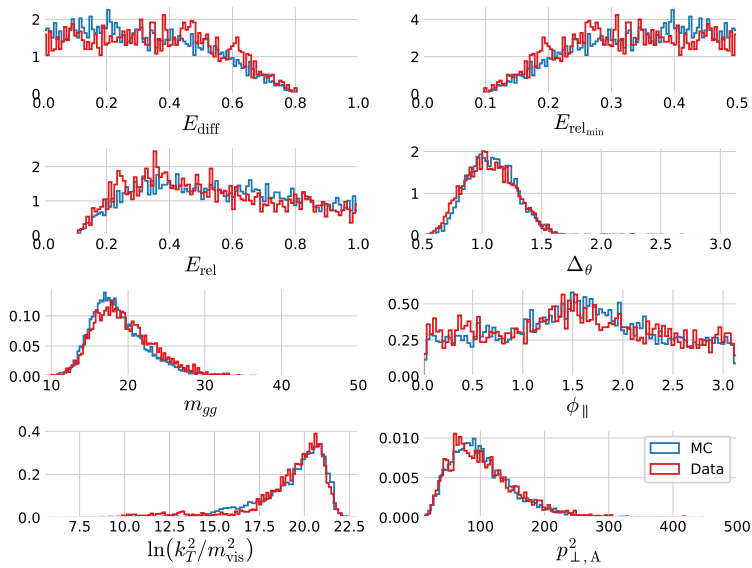
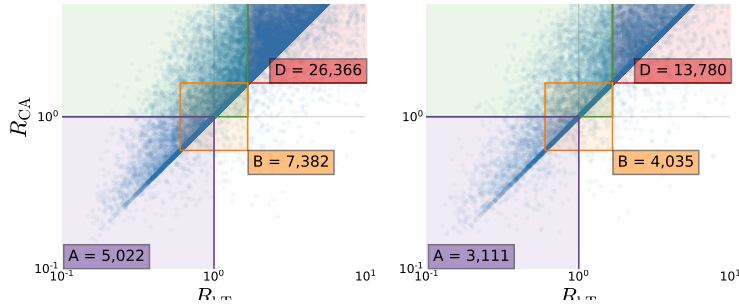
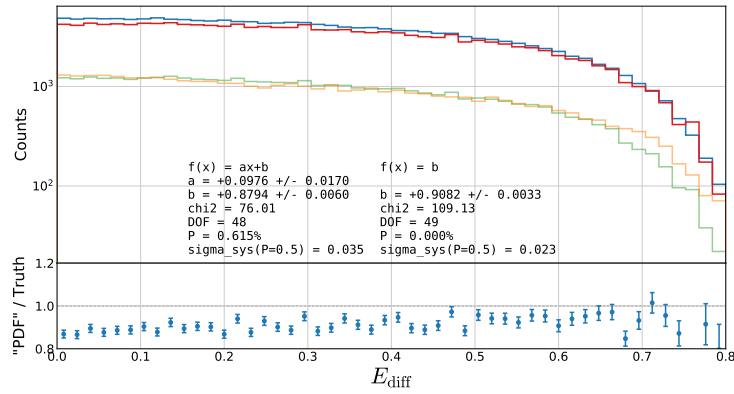
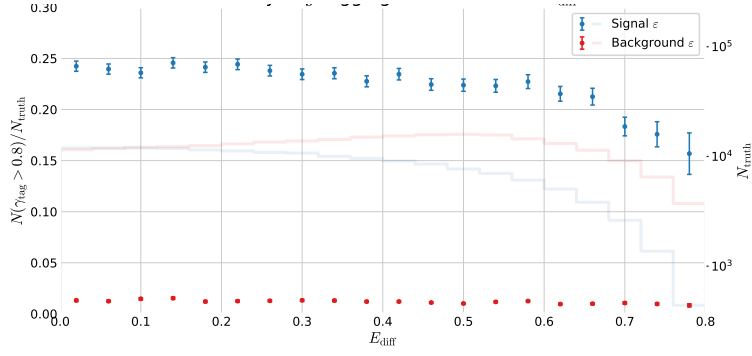


Figure 5.30: Efficiency of the g-tags for 4-jet events as a function of normalized gluon gluon jet energy difference in Monte Carlo. The efficiency is measured as the number of events with a g-tag higher than 0.8 ( $\gamma > 0.8$ ) out of the total number and the normalized gluon gluon jet energy difference  $A$  is  $A = \frac{E_{g\max} - E_{g\min}}{E_{g\max} + E_{g\min}}$  where  $E_{g\max}$  ( $E_{g\min}$ ) refers to the energy of the gluon with the highest (lowest) energy. The efficiency is plotted for **signal events** according to MC Truth in blue and **background events** according to MC Truth in red.

Figure 5.31: Closure plot between MC Truth and the corrected g-tagging model in 4-jet events for the normalized gluon gluon jet energy difference. The corrected g-tagging model is described in further detail in section XXX **TODO!**. In the top part of the plot the **MC Truth** is shown in blue, the **corrected g-tagging model** "PDF 2gg" in red, the **g-signal distribution** in semi-transparent green and the **g-sideband distribution** in semi-transparent orange. In the bottom part of the plot the ratio between MC Truth and the output of the corrected g-tagging model is shown. The normalized gluon gluon jet energy difference  $A$  is  $A = \frac{E_{g\max} - E_{g\min}}{E_{g\max} + E_{g\min}}$  where  $E_{g\max}$  ( $E_{g\min}$ ) refers to the energy of the gluon with the highest (lowest) energy.

Figure 5.33: R kt CA cut region A **XXX TODO!**



## *B. Quarks vs. Gluons Appendix*

	$b$	$c$	$uds$	$g$	non- $q$ -matched
2	37.2 %	12.9 %	29.1 %	0.0 %	20.7 %
3	22.6 %	8.9 %	19.7 %	31.2 %	17.5 %
4	14.6 %	7.0 %	15.0 %	45.1 %	18.3 %
5	10.0 %	5.7 %	12.2 %	52.5 %	19.6 %
6	7.1 %	4.4 %	8.8 %	54.4 %	25.2 %

Table B.1: Number of different types of jets for MC and MCb written in relative numbers such that each row sum to 100 %. See also Table 5.3.

figures/quarks/viz\_UMAP\_test\_0.5\_input2b\_njet=4\_algorithm=UMAP.pdf

Figure B.1: Grid search of the two parameters `n_neighbors` and `min_dist` for the UMAP algorithm run on 4-jet events. For an explanation of these, see [section 5.2](#).

figures/quarks/viz\_TSNE\_MULTI\_test\_0.5\_input2b\_njet=4\_algorithm=tsne\_multi.pdf

Figure B.2: Visualization of the t-SNE algorithm as a function of the `perplexity` parameters for 4-jet events.

Hyperparameter	Range
subsample	$\mathcal{U}(0.4, 1)$
colsample_bytree	$\mathcal{U}_{\text{trunc}}(0.4, 1, 2)$
max_depth	$\mathcal{U}_{\text{int}}(1, 20)$
min_child_weight	$\mathcal{U}_{\text{int}}(0, 10)$

Table B.2: Probability Density Functions for the random search hyperparameter optimization process for the XGBoost model.

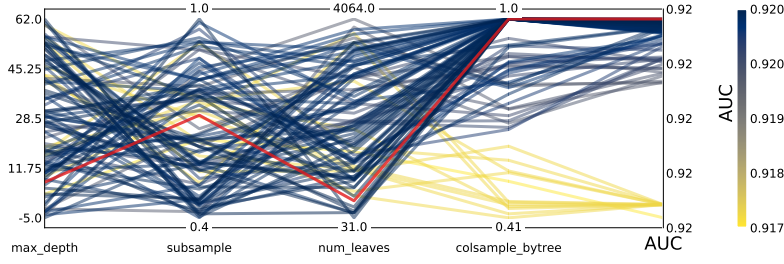


Figure B.3: Hyperparameter optimization results of  $b$ -tagging for 3-jet events. The results are shown as parallel coordinates with each hyperparameter along the  $x$ -axis and the value of that parameter on the  $y$ -axis. Each line is an event in the 4-dimensional space colored according to the performance of that hyperparameter as measured by AUC from highest AUC in dark blue to lowest AUC in yellow. The single best hyperparameter is shown in red.

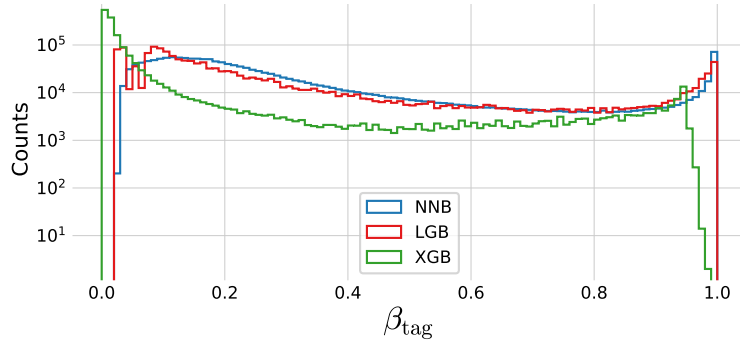


Figure B.4: Histogram of  $b$ -tag scores  $\beta_{\text{tag}}$  in 3-jet events for **NNB** (the neural network pre-trained by ALEPH, also called `nnbjct`) in blue, **LGB** in red, and **XGB** in green.

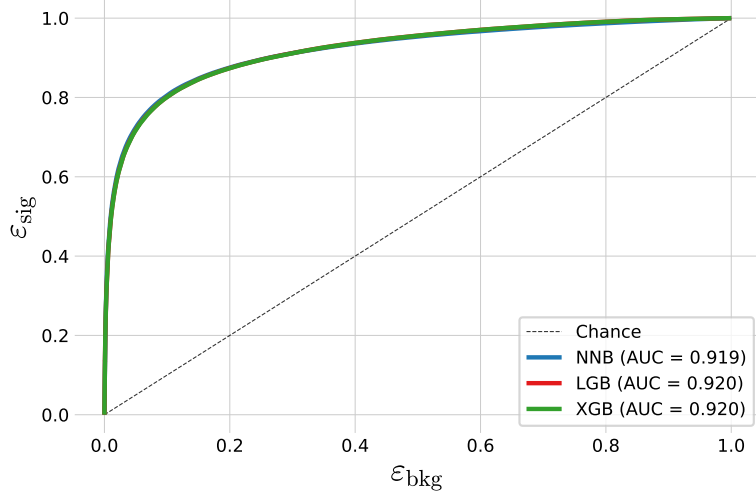


Figure B.5: ROC curve of the three  $b$ -tag models in 3-jet events for NNB (the pre-trained neural network trained by ALEPH, also called `nbnjet`) in blue, LGB in red, and XGB in green. In the legend the area under curve (AUC) is also shown. Notice that the LGB and XGB models share performance and it is thus due to overplotting that only the green line for XGB can be seen. In the machine learning community the background efficiency  $\epsilon_{\text{bkg}}$  is sometimes known as the false positive rate (FPR) and the signal efficiency  $\epsilon_{\text{sig}}$  as the true positive rate (TPR).

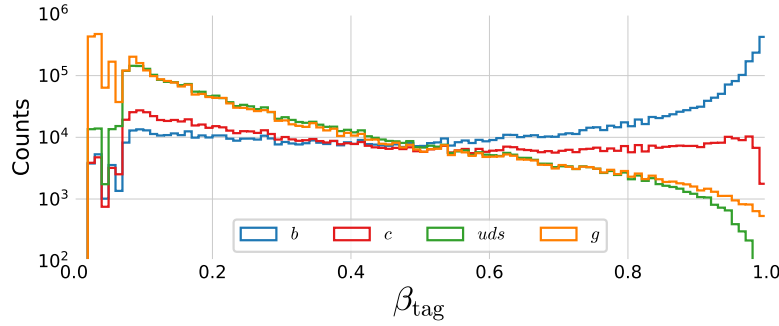


Figure B.6: Distribution of  $b$ -tags in 3-jet events for  $b$ -jets in blue,  $c$ -jets in red,  $uds$  in green and  $g$  in orange.

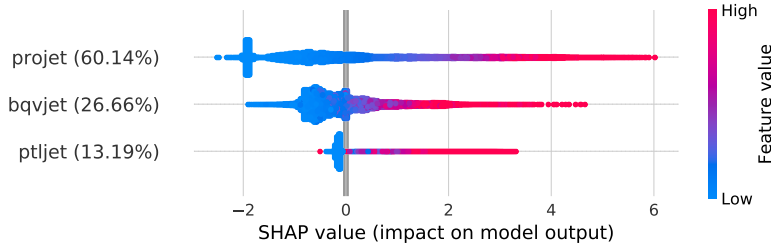


Figure B.7: Global feature importances for the LGB  $b$ -tagging algorithm on 3-jet events. The normalized feature importance is shown in the parenthesis and each dot is an observation showing the dependence between the SHAP value and the feature's value.

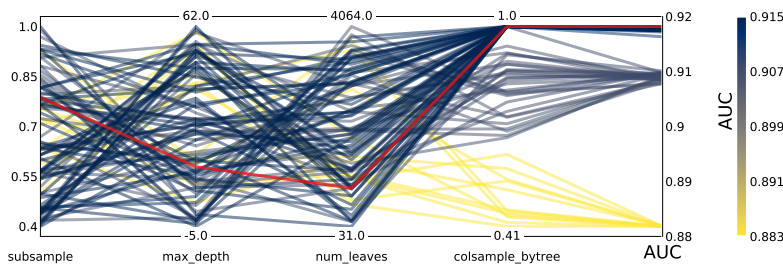


Figure B.8: Hyperparameter optimization results of  $g$ -tagging for 3-jet events for energy ordered jets.



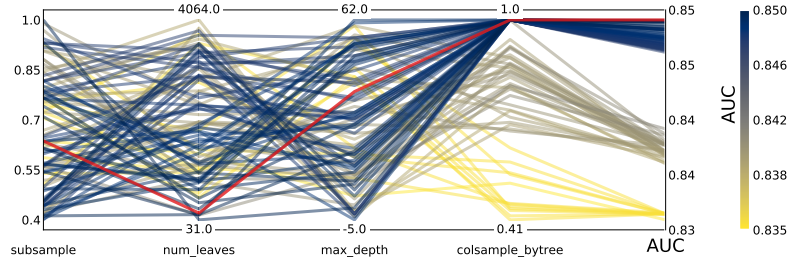


Figure B.9: Hyperparameter optimization results of  $g$ -tagging for 3-jet events for (row) shuffled jets.

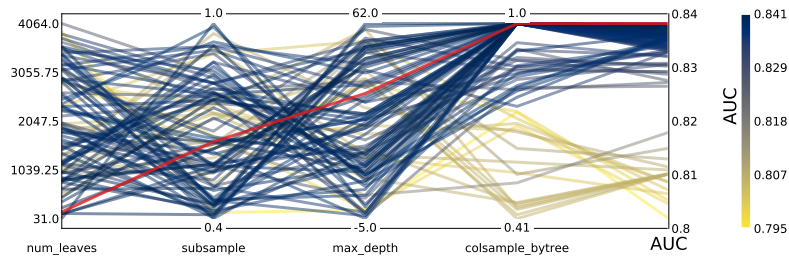


Figure B.10: Hyperparameter optimization results of  $g$ -tagging for 4-jet events for energy ordered jets.

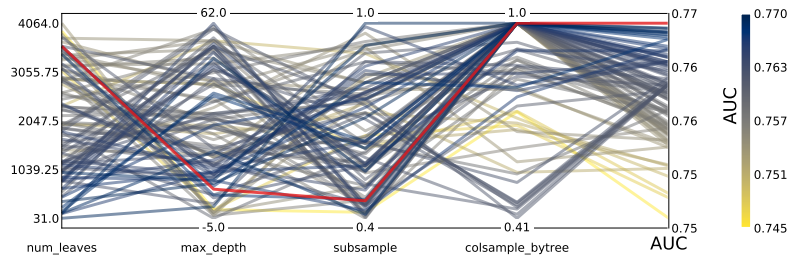


Figure B.11: Hyperparameter optimization results of  $g$ -tagging for 4-jet events for (row) shuffled jets.

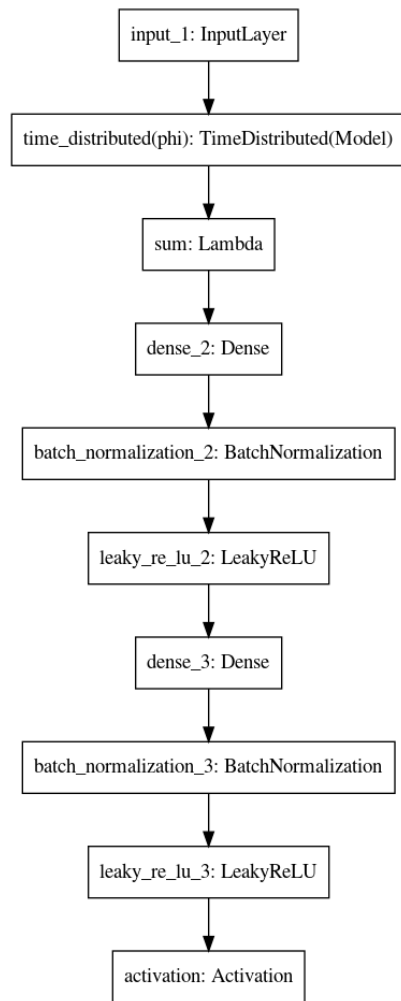


Figure B.12: Architecture of the PermNet neural network.

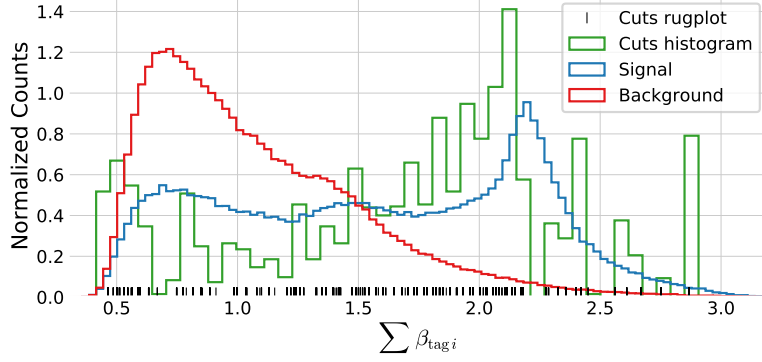


Figure B.13: Histogram of the distribution of **signal** in blue and **background** in red for the 1-dimensional sum of  $b$ -tags for 4-jet events. A histogram of the **cut values** from the LGB model trained on this data is shown in green together with a rug plot of the cut values in black. Notice how most of the cuts match up with the signal peak at around a  $\sum \beta_i \sim 2.1$ .

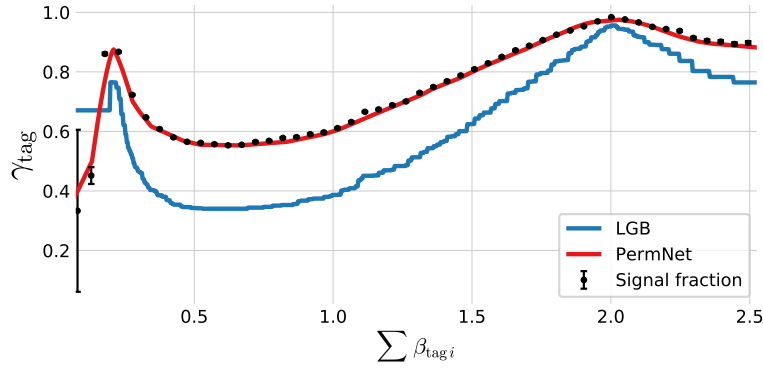


Figure B.14: Plot of the (1D)  $g$ -tag scores for 3-jet events as a function of  $\sum \beta_i$  for the **LGB** model in blue and the **PermNet** model in red. The signal fraction (based on the signal and background histograms in Figure B.15) is plotted as black error bars where the size of the error bars is based on the propagated uncertainties of the signal and background histogram assuming Poissonian statistics.

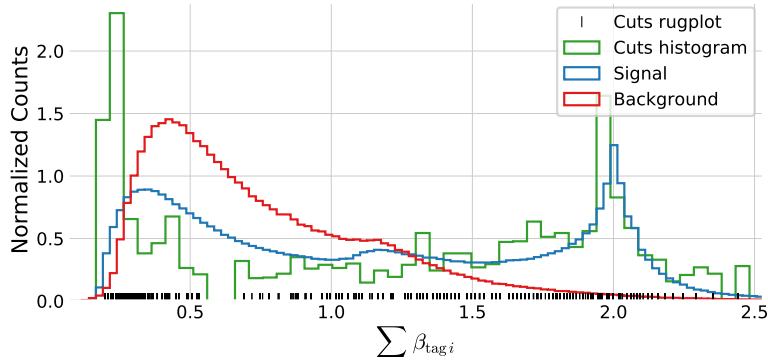


Figure B.15: Histogram of the distribution of **signal** in blue and **background** in red for the 1-dimensional sum of  $b$ -tags for 3-jet events. A histogram of the **cut values** from the LGB model trained on this data is shown in green together with a rug plot of the cut values in black. Notice how most of the cuts match up with the signal peak at around a  $\sum \beta_i \sim 2.1$ .

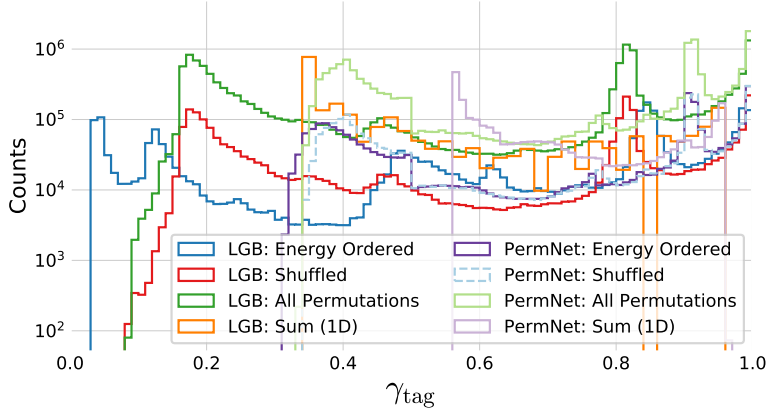


Figure B.16: Distribution of  $g$ -tag scores in 3-jet events shown with a logarithmic  $y$ -scale for LGB: Energy Ordered in blue, LGB: Shuffled in red, LGB: All Permutations in green, LGB: Sum 1D in orange, PermNet: Energy Ordered in purple, PermNet: Shuffled in light-blue, PermNet: All Permutations in light-green, PermNet: Sum 1D in light-purple. Here LGB and PermNet are the two different type of models and “Energy Ordered”, “Shuffled”, “All Permutations”, and “Sum 1D” are the different methods used for making the input data permutation invariant (except energy ordered).

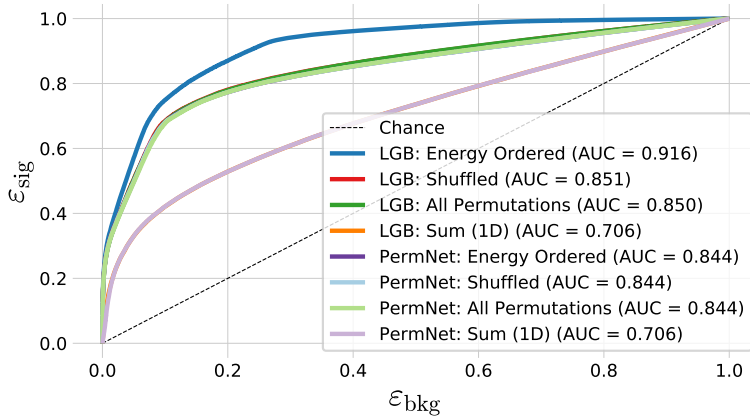
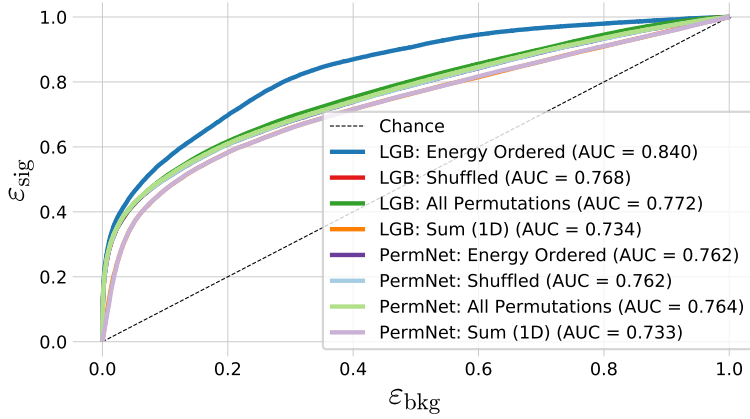


Figure B.17: ROC curve of the eight  $g$ -tag models in 4-jet events. First one in dashed black is the ROC curve that you get by random chance. The colors are the same as in Figure 5.18 and in the legend also the Area Under the ROC curve (AUC) is shown. Notice that the XGB model which uses the energy ordered data produced the best model, however, this model is not permutation invariant. Of the permutation invariant models (the rest), the XGB model trained on all permutations of the  $b$ -tags performs highest. The lowest performing models are the two models trained only on the 1-dimensional sum of  $b$ -tags, as expected, however, still with a better performance than expected by the author.

Figure B.18: ROC curve of the eight  $g$ -tag models in 3-jet. First one in dashed black is the ROC curve that you get by random chance. The colors are the same as in Figure 5.18 and in the legend also the Area Under the ROC curve (AUC) is shown.

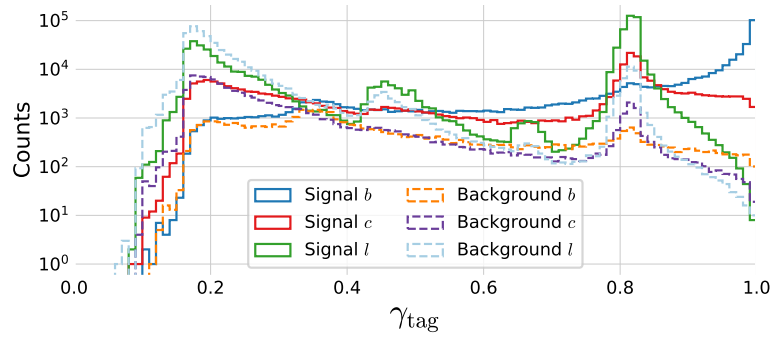


Figure B.19: Histogram of  $g$ -tag scores from the LGB-model in 3-jet events for  $b$  signal in blue,  $c$  signal in red,  $l$  ( $uds$ ) signal in green,  $b$  background in orange,  $c$  background in purple,  $l$  ( $uds$ ) background in light-blue.

# List of Figures

2.1	The learning problem.	6
2.2	Approximation-Estimation Tradeoff	10
2.3	Regularization Strength	11
2.4	Regularization Effect of $L_2$	12
2.5	Regularization Effect of $L_1$	12
2.6	$k$ -Fold Cross Validation	13
2.7	$k$ -Fold Cross Validation for Time Series Data	13
2.8	Objective Functions.	16
2.9	Objective Functions Zoom In.	16
2.10	Decision Tree Cuts In Feature Space	16
2.11	Decision Tree	17
2.12	Grid Search	20
2.13	Random Search	21
2.14	Bayesian Optimization	22
3.1	Danish Housing Price Index	27
3.2	Distributions for the housing price dataset	28
3.3	Distributions for the housing price dataset	29
3.4	Histogram of prices of houses and apartments sold in Denmark	30
3.5	Linear correlation between variables and price	31
3.6	Comparison of the Linear Correlation $\rho$ and the Non-Linear MIC.	31
3.7	Non-linear correlation between variables and price	32
3.8	Validity of input features	32
3.9	Validity Dendrogram	33
3.10	Prophet Forecast for apartments	35
3.11	Prophet Trends	35
3.12	XXX	37
3.13	Parallel Coordinate Plot of the Initial Hyperparameter Optimization for Apartments	38
3.14	Initial HPO Results for the Weight Half-life $T_{\frac{1}{2}}$	38
3.15	Initial HPO Results for the Loss Function	38
3.16	XXX	39
3.17	Hyperparameter optimization: random search results	40
3.18	Early Stopping results	40
3.19	Performance of XGB-model on apartment prices	41
3.20	Standard Deviation and MAD of the Static and Dynamic XGB Forecasts	42
3.21	Market Index based on the Static and Dynamic XGB Forecasts	43
3.22	SHAP Prediction Explanation for apartment	44

3.23 Feature importance of apartments prices using XGB	45
3.24 Feature importance of apartments prices using XGB XXX	45
3.25 Multiple Models XXX	46
3.26 SHAP plot villa TFIDF XXX	49
4.1 The Standard Model	54
4.2 Feynman diagram for the jet production at LEP	55
4.3 Quark splitting	55
4.4 Hadronization process	56
4.5 The ALEPH detector	57
4.6 Polar angle	57
4.7 Azimuthal angle	57
5.1 Histograms of the vertex variables	65
5.2 UMAP visualization of vertex variables for 4-jet events	66
5.3 UMAP visualization of vertex variables for 3-jet events	66
5.4 UMAP visualization of vertex variables for 2-jet events	66
5.5 Correlation of Vertex Variables	67
5.6 Plot of the log-loss $\ell_{\log}$	68
5.7 Hyperparameter Optimization of $b$ -tagging	69
5.8 Parallel Plot of HPO Results for 4-Jet $b$ -Tagging	69
5.9 $b$ -Tag Scores in 4-Jet Events	70
5.10 ROC curve for 4-jet $b$ -tagging	70
5.11 Distribution of $b$ -Tags in 4-Jet Events	71
5.12 Global Feature Importances for the LGB $b$ -Tagging Algorithm on 4-Jet Events	71
5.13 The expit Function	71
5.14 The logit Function	71
5.15 SHAP 3-Jet Model Explanation for $b$ -like Jet	72
5.16 Hyperparameter Optimization of $g$ -tagging	74
5.17 1D Sum Models Predictions and Signal Fraction for 4-jets events	75
5.18 $g$ -Tag Scores in 4-Jet Events	76
5.19 ROC Curve for $g$ -Tag in 4-Jet Events	77
5.20 Distribution of $g$ -Tag Scores in 4-Jet Events for Signal and Background	77
5.21 Distribution of $b$ -Tag Scores in 3-Jet $l$ -Quark Events for Low and High $g$ -Tag Values	78
5.22 XXX	78
5.23 Monte Carlo – Data bias for $b$ -tags and jet energy	79
5.24 $b$ -Tagging Efficiency $\epsilon_b^{b\text{-sig}}$ as a function of jet energy	79
5.25 $b$ -Tagging Efficiency $\epsilon_b^{g\text{-sig}}$ as a function of jet energy	79
5.26 $b$ -Tagging Efficiency $\epsilon_g^{g\text{-sig}}$ as a function of jet energy	80
5.27 $b$ -Tagging Efficiency $\epsilon_g^{b\text{-sig}}$ as a function of jet energy	80
5.28 $g$ -Tagging proxy efficiency for $b\bar{b}g$ -events as function of the mean invariant mass	80
5.29 $g$ -Tagging proxy efficiency for $b\bar{b}g$ -events as function of $g$ -tag	80
5.30 $g$ -Tagging efficiency for 4-jet events in MC as a function of normalized gluon gluon jet energy difference	81
5.31 Closure plot between MC Truth and the corrected $g$ -tagging model in 4-jet events for the normalized gluon gluon jet energy difference	81

5.32 R kt CA overview XXX <b>TODO!</b>	81
5.33 R kt CA cut region A XXX <b>TODO!</b>	81
A.1 Validity Heatmap	87
A.2 Distributions for the housing price dataset	88
A.3 Distributions for the housing price dataset	89
A.4 Distributions for the housing price dataset	90
A.5 Distributions for the housing price dataset	91
A.6 Distributions for the housing price dataset	92
A.7 Distributions for the housing price dataset	93
A.8 Distributions for the housing price dataset	94
A.9 Distributions for the housing price dataset	95
A.10 Distributions for the housing price dataset	96
A.11 Distributions for the housing price dataset	97
A.12 Distributions for the housing price dataset	98
A.13 Distributions for the housing price dataset	99
A.14 Distributions for the housing price dataset	100
A.15 Distributions for the housing price dataset	101
A.16 Linear Correlations	103
A.17 MIC non-linear correlation	104
A.18 Prophet Forecast for apartments	105
A.19 Prophet Trends	105
A.20 Overview of initial hyperparameter optimization of the housing model for houses	109
A.21 XXX	110
A.22 XXX	110
A.23 XXX	110
A.24 XXX	111
A.25 XXX	111
A.26 XXX	111
A.27 Performance of XGB-model on apartment prices	112
B.1 UMAP Parameter Grid Search	117
B.2 Visualization of the t-SNE algorithm	117
B.3 Parallel Plot of HPO results for 3-jet $b$ -Tagging	118
B.4 $b$ -tag scores in 3-jet events	118
B.5 ROC curve for 3-jet $b$ -tagging	119
B.6 Distribution of $b$ -Tags in 3-Jet Events	119
B.7 Global Feature Importances for the LGB $b$ -Tagging Algorithm on 3-Jet Events	119
B.8 Parallel Plot of HPO Results for 3-Jet $g$ -Tagging for Energy Ordered Jets	119
B.9 Parallel Plot of HPO Results for 3-Jet $g$ -Tagging for Shuffled Jets	120
B.10 Parallel Plot of HPO Results for 4-Jet $g$ -Tagging for Energy Ordered Jets	120
B.11 Parallel Plot of HPO Results for 4-Jet $g$ -Tagging for Shuffled Jets	120
B.12 PermNet Architecture	121
B.13 1D LGB Model Cuts for 4-jets events	122
B.14 1D Sum Models Predictions and Signal Fraction for 3-jets events	122



B.15 1D LGB Model Cuts for 3-jets events	122
B.16 $g$ -Tag Scores in 3-Jet Events	123
B.17 ROC curve for $g$ -tag in 4-jet events	123
B.18 ROC Curve for $g$ -Tag in 3-Jet Events	123
B.19 Distribution of $g$ -Tag Scores in 3-Jet Events for Signal and Background	124

# List of Tables

3.1	Mapping between the code in <code>SagTypeNr</code> and the type of residence. The two important types of residences are villa (one-family houses) and ejerlejlighed (owner-occupied apartments).	29
3.2	Basic Cuts	33
3.3	Side Door Mapping.	33
3.4	Street Mapping	33
3.5	Number of Observations in the Housing Dataset	36
3.6	Number of Observations in the Housing Dataset for the Tight Selection	36
3.7	Results of the initial hyperparameter optimization for apartments for the best loss function $\ell_{\text{Cauchy}}$ .	37
3.8	Results of the initial hyperparameter optimization for houses for the best loss function $\ell_{\text{Cauchy}}$ .	37
3.9	PDFs Used in the Random Search	39
3.10	Realtors' MAD	41
3.11	Performance Metrics for the Housing Model on Apartments	43
3.12	Performance Metrics for the Housing Model on Houses	43
5.1	Dimensions of dataset for Data	64
5.2	Dimensions of dataset for MC and MCb	64
5.3	Number of different types of jets for MC and MCb. See also Table B.1 in the appendix for relative numbers.	65
5.4	Random Search PDFs for LGB	68
A.1	XXX <b>TODO!</b>	102
A.2	Energy Rating Mapping	104
A.3	Rmse-ejerlejlighed-appendix.	106
A.4	Logcosh-ejerlejlighed-appendix.	106
A.5	Cauchy-ejerlejlighed-appendix.	106
A.6	Welsch-ejerlejlighed-appendix.	107
A.7	Fair-ejerlejlighed-appendix.	107
A.8	Rmse-villa-appendix.	107
A.9	Logcosh-villa-appendix.	107
A.10	Cauchy-villa-appendix.	108
A.11	Welsch-villa-appendix.	108
A.12	Fair-villa-appendix.	108
A.13	XXX ejer tight	113
A.14	XXX villa tight	113

B.1	Number of different types of jets for MC and MCb written in relative numbers such that each row sum to 100 %. See also Table 5.3.	116
B.2	Random Search PDFs for XGB	118

# Bibliography

- [1] Advanced Topics in Machine Learning (ATML). URL <https://kurser.ku.dk/course/ndak15014u>.
- [2] Allstate Claims Severity - Fair Loss. URL <https://kaggle.com/c/allstate-claims-severity>.
- [3] Dmlc/xgboost. URL <https://github.com/dmlc/xgboost>.
- [4] HEP meets ML award | The Higgs Machine Learning Challenge. URL <https://higgsml.lal.in2p3.fr/prizes-and-award/award/>.
- [5] The Large Electron-Positron Collider | CERN. URL <https://home.cern/science/accelerators/large-electron-positron-collider>.
- [6] Microsoft/LightGBM. URL [https://github.com/microsoft/LightGBM/blob/b397555d7023fd05f8e56326905fe7b185109de/src/treelearner/serial\\_tree\\_learner.cpp#L282](https://github.com/microsoft/LightGBM/blob/b397555d7023fd05f8e56326905fe7b185109de/src/treelearner/serial_tree_learner.cpp#L282).
- [7] Scikit-hep/uproot. URL <https://github.com/scikit-hep/uproot>.
- [8] Datashader: Revealing the Structure of Genuinely Big Data. URL <https://github.com/holoviz/datashader>.
- [9] O. . Www.OIS.dk - Din genvej til ejendomsdata. URL <https://www.ois.dk/>.
- [10] M. Abadi et al. TensorFlow: Large-scale machine learning on heterogeneous systems. URL <http://tensorflow.org/>.
- [11] Y. S. Abu-Mostafa, M. Magdon-Ismail, and H.-T. Lin. *Learning From Data*. AMLBook. ISBN 978-1-60049-006-4.
- [12] D. Albanese, S. Riccadonna, C. Donati, and P. Franceschi. A practical tool for maximal information coefficient analysis. 7. ISSN 2047-217X. doi: 10.1093/gigascience/giy032. URL <https://doi.org/10.1093/gigascience/giy032>.
- [13] E. Anderson. The Species Problem in Iris. 23(3):457–509. ISSN 00266493. doi: 10.2307/2394164. URL [www.jstor.org/stable/2394164](http://www.jstor.org/stable/2394164).

- [14] B. Andersson, G. Gustafson, G. Ingelman, and T. Sjöstrand. Parton fragmentation and string dynamics. 97(2):31–145. ISSN 0370-1573. doi: 10.1016/0370-1573(83)90080-7. URL <http://www.sciencedirect.com/science/article/pii/0370157383900807>.
- [15] S. R. Armstrong. A Search for the standard model Higgs boson in four jet final states at center-of-mass energies near 183-GeV with the ALEPH detector at LEP. URL <http://wwwlib.umi.com/dissertations/fullcit?p9910371>.
- [16] M. Awad and R. Khanna. Support Vector Regression. In M. Awad and R. Khanna, editors, *Efficient Learning Machines: Theories, Concepts, and Applications for Engineers and System Designers*, pages 67–80. Apress. ISBN 978-1-4302-5990-9. doi: 10.1007/978-1-4302-5990-9\_4. URL [https://doi.org/10.1007/978-1-4302-5990-9\\_4](https://doi.org/10.1007/978-1-4302-5990-9_4).
- [17] R. J. Barlow. *Statistics: A Guide to the Use of Statistical Methods in the Physical Sciences (Manchester Physics Series)*. WileyBlackwell, reprint edition. ISBN 0-471-92295-1. URL <http://www.amazon.co.uk/Statistics-Statistical-Physical-Sciences-Manchester/dp/0471922951%3FSubscriptionId%3D13CT5CVB80YFWJEPWS02%26tag%3Dws%26linkCode%3Dxm2%26camp%3D2025%26creative%3D165953%26creativeASIN%3D0471922951>.
- [18] J. T. Barron. A General and Adaptive Robust Loss Function. URL <http://arxiv.org/abs/1701.03077>.
- [19] W. Bartel et al. Experimental study of jets in electron-positron annihilation. 101(1):129–134. ISSN 0370-2693. doi: 10.1016/0370-2693(81)90505-0. URL <http://www.sciencedirect.com/science/article/pii/0370269381905050>.
- [20] E. Becht, C.-A. Dutertre, I. W. H. Kwok, L. G. Ng, F. Ginhoux, and E. W. Newell. Evaluation of UMAP as an alternative to t-SNE for single-cell data. page 298430, . doi: 10.1101/298430. URL <https://www.biorxiv.org/content/10.1101/298430v1>.
- [21] E. Becht, L. McInnes, J. Healy, C.-A. Dutertre, I. W. H. Kwok, L. G. Ng, F. Ginhoux, and E. W. Newell. Dimensionality reduction for visualizing single-cell data using UMAP. 37(1):38–44, . ISSN 1546-1696. doi: 10.1038/nbt.4314. URL <https://www.nature.com/articles/nbt.4314>.
- [22] J. Bergstra and Y. Bengio. Random Search for Hyperparameter Optimization. 13:281–305. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=2188385.2188395>.
- [23] C. Bierlich. Rope hadronization, geometry and particle production in pp and pA Collisions. URL <https://lup.lub.lu.se/search/ws/files/18474576/thesis.pdf>.

- [24] Bolighed. Bolighed - usikkerhed i data-vurderingen. URL <https://bolighed.dk/om-bolighed/spoergsmaal-og-svar/#boligvaerdi>.
- [25] L. Breiman. Random Forests. 45(1):5–32. ISSN 1573-0565. doi: 10.1023/A:1010933404324. URL <https://doi.org/10.1023/A:1010933404324>.
- [26] E. Brochu, V. M. Cora, and N. de Freitas. A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning. URL <http://arxiv.org/abs/1012.2599>.
- [27] R. Brun and F. Rademakers. ROOT — An object oriented data analysis framework. 389(1):81–86. ISSN 0168-9002. doi: 10.1016/S0168-9002(97)00048-X. URL <http://www.sciencedirect.com/science/article/pii/S016890029700048X>.
- [28] A. Buckley, J. Butterworth, S. Gieseke, D. Grellscheid, S. Hoche, H. Hoeth, F. Krauss, L. Lonnblad, E. Nurse, P. Richardson, S. Schumann, M. H. Seymour, T. Sjostrand, P. Skands, and B. Webber. General-purpose event generators for LHC physics. 504(5):145–233. ISSN 03701573. doi: 10.1016/j.physrep.2011.03.005. URL <http://arxiv.org/abs/1101.2599>.
- [29] C. Burgard. Standard model of physics | TikZ example. URL <http://www.texample.net/tikz/examples/model-physics/>.
- [30] D. Buskulic et al. An investigation of B<sub>d</sub> and B<sub>s</sub> oscillation. 322(4):441–458. ISSN 0370-2693. doi: 10.1016/0370-2693(94)91177-0. URL <http://www.sciencedirect.com/science/article/pii/0370269394911770>.
- [31] T. Chen and C. Guestrin. XGBoost: A Scalable Tree Boosting System. pages 785–794. doi: 10.1145/2939672.2939785. URL <http://arxiv.org/abs/1603.02754>.
- [32] T. A. Collaboration. Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC. 716(1):1–29, . ISSN 03702693. doi: 10.1016/j.physletb.2012.08.020. URL <http://arxiv.org/abs/1207.7214>.
- [33] T. C. Collaboration. Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC. 716(1):30–61, . ISSN 03702693. doi: 10.1016/j.physletb.2012.08.021. URL <http://arxiv.org/abs/1207.7235>.
- [34] A. Diaz-Papkovich, L. Anderson-Trocmé, C. Ben-Eghan, and S. Gravel. UMAP reveals cryptic population structure and phenotype heterogeneity in large genomic cohorts. 15(11). ISSN 1553-7390. doi: 10.1371/journal.pgen.1008432. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6853336/>.

- [35] S. D. DST. Price Index (EJ14) - Statistics Denmark. URL <https://www.dst.dk/en/Statistik/emner/priser-og-forbrug/ejendomme>.
- [36] D. et al. Buskulic. A precise measurement of hadrons. 313(3): 535–548. ISSN 0370-2693. doi: 10.1016/0370-2693(93)90028-G. URL <http://www.sciencedirect.com/science/article/pii/037026939390028G>.
- [37] F. Faye. Frederik Faye / deepcalo. URL <https://gitlab.com/ffaye/deepcalo>.
- [38] R. A. Fisher. The Use of Multiple Measurements in Taxonomic Problems. 7(2):179–188. ISSN 2050-1439. doi: 10.1111/j.1469-1809.1936.tb02137.x. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1469-1809.1936.tb02137.x>.
- [39] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. In P. Vitányi, editor, *Computational Learning Theory*, pages 23–37. Springer Berlin Heidelberg. ISBN 978-3-540-49195-8. Adaboost.
- [40] S. L. Glashow. Partial-symmetries of weak interactions. 22(4): 579–588. ISSN 0029-5582. doi: 10.1016/0029-5582(61)90469-2. URL <http://www.sciencedirect.com/science/article/pii/0029558261904692>.
- [41] N. Guttenberg, N. Virgo, O. Witkowski, H. Aoki, and R. Kanai. Permutation-equivariant neural networks applied to dynamics prediction. URL <http://arxiv.org/abs/1612.04530>.
- [42] A. E. Harvey and S. Peters. Estimation Procedures for Structural Time Series Models. doi: 10.1002/for.3980090203.
- [43] T. Hastie and R. Tibshirani. Generalized Additive Models: Some Applications. 82(398):371–386. ISSN 01621459. doi: 10.2307/2289439. URL [www.jstor.org/stable/2289439](http://www.jstor.org/stable/2289439).
- [44] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition*. Springer Series in Statistics. Springer-Verlag, 2 edition. ISBN 978-0-387-84857-0. URL <http://www.springer.com/la/book/9780387848570>.
- [45] K. Hornik. Approximation capabilities of multilayer feedforward networks. 4(2):251–257. ISSN 0893-6080. doi: 10.1016/0893-6080(91)90009-T. URL <http://www.sciencedirect.com/science/article/pii/089360809190009T>.
- [46] P. Huber and E. Ronchetti. *Robust Statistics*. Wiley Series in Probability and Statistics. Wiley. ISBN 978-1-118-21033-8. URL [https://books.google.dk/books?id=j10hquR\\_j88C](https://books.google.dk/books?id=j10hquR_j88C).

- [47] S. Hviid, Juul. Working Paper: A regional model of the Danish housing market. URL <http://www.nationalbanken.dk/en/publications/Pages/2017/11/Working-Paper-A-regional-model-of-the-Danish-housing-market.aspx>.
- [48] F. James and M. Roos. Minuit – a system for function minimization and analysis of the parameter errors and correlations. 10:343–367. doi: 10.1016/0010-4655(75)90039-9.
- [49] R. E. Kalman. A new approach to linear filtering and prediction problems. 82:35–45.
- [50] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 3146–3154. Curran Associates, Inc. URL <http://papers.nips.cc/paper/6907-lightgbm-a-highly-efficient-gradient-boosting-decision-tree.pdf>.
- [51] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. URL <http://arxiv.org/abs/1412.6980>.
- [52] C. Leys, C. Ley, O. Klein, P. Bernard, and L. Licata. Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median. 49(4): 764–766. ISSN 0022-1031. doi: 10.1016/j.jesp.2013.03.013. URL <http://www.sciencedirect.com/science/article/pii/S0022103113000668>.
- [53] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, pages 4768–4777. Curran Associates Inc. ISBN 978-1-5108-6096-4. URL <http://dl.acm.org/citation.cfm?id=3295222.3295230>.
- [54] S. M. Lundberg, G. G. Erion, and S.-I. Lee. Consistent Individualized Feature Attribution for Tree Ensembles - SHAP. . URL <http://arxiv.org/abs/1802.03888>.
- [55] S. M. Lundberg, G. G. Erion, and S.-I. Lee. Consistent Individualized Feature Attribution for Tree Ensembles. . URL <http://arxiv.org/abs/1802.03888>.
- [56] A. L. Maas. Rectifier nonlinearities improve neural network acoustic models.
- [57] L. McInnes and J. Healy. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. URL <http://arxiv.org/abs/1802.03426>.



- [58] T. C. Mills. *Time Series Techniques for Economists* / Terence c. Mills. Cambridge University Press Cambridge [England] ; New York. ISBN 0-521-34339-9 0-521-40574-2. URL <http://www.loc.gov/catdir/toc/cam031/89007187.html>.
- [59] I. Mulalic, H. Rasmussen, J. Rouwendal, and H. H. Woltmann. The Financial Crisis and Diverging House Prices: Evidence from the Copenhagen Metropolitan Area. ISSN 1556-5068. doi: 10.2139/ssrn.3041272. URL <https://www.ssrn.com/abstract=3041272>.
- [60] Particle Data Group et al. Review of Particle Physics. 98(3):030001. doi: 10.1103/PhysRevD.98.030001. URL <https://link.aps.org/doi/10.1103/PhysRevD.98.030001>.
- [61] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. 12:2825–2830.
- [62] E. Polley and M. van der Laan. Super Learner In Prediction. URL <https://biostats.bepress.com/ucbbiostat/paper266>.
- [63] J. Proriot, J. Jousset, C. Guicheney, A. Falvard, P. Henrard, D. Pallin, P. Perret, and B. Brandl. TAGGING B QUARK EVENTS IN ALEPH WITH NEURAL NETWORKS (comparison of different methods : Neural Networks and Discriminant Analysis). page 27.
- [64] A. Purcell. Go on a particle quest at the first CERN webfest. URL <https://cds.cern.ch/record/1473657>.
- [65] S. Ravanbakhsh, J. Schneider, and B. Póczos. Deep Learning with Sets and Point Clouds. URL <http://arxiv.org/abs/1611.04500>.
- [66] D. N. Reshef, Y. A. Reshef, H. K. Finucane, S. R. Grossman, G. McVean, P. J. Turnbaugh, E. S. Lander, M. Mitzenmacher, and P. C. Sabeti. Detecting Novel Associations in Large Data Sets. 334(6062):1518–1524. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.1205438. URL <http://science.sciencemag.org/content/334/6062/1518>.
- [67] P. J. Rousseeuw and C. Croux. Alternatives to the Median Absolute Deviation. 88(424):1273–1283. ISSN 0162-1459. doi: 10.1080/01621459.1993.10476408. URL <https://www.tandfonline.com/doi/abs/10.1080/01621459.1993.10476408>.
- [68] A. Salam. Weak and electromagnetic interactions. In *Selected Papers of Abdus Salam*, volume Volume 5 of *World Scientific Series in 20th Century Physics*, pages 244–254. WORLD SCIENTIFIC. ISBN 978-981-02-1662-7. doi: 10.1142/9789812795915\_

0034. URL [https://www.worldscientific.com/doi/abs/10.1142/9789812795915\\_0034](https://www.worldscientific.com/doi/abs/10.1142/9789812795915_0034).
- [69] L. Scodellaro. B tagging in ATLAS and CMS. URL <http://arxiv.org/abs/1709.01290>.
- [70] L. Shapley. A value for n-person games. In *The Shapley Value*, volume 28 of *Annals of Math Studies*, pages 307–317. doi: 10.1017/CBO9780511528446.003.
- [71] T. Sjöstrand, S. Ask, J. R. Christiansen, R. Corke, N. Desai, P. Ilten, S. Mrenna, S. Prestel, C. O. Rasmussen, and P. Z. Skands. An Introduction to PYTHIA 8.2. 191:159–177. ISSN 00104655. doi: 10.1016/j.cpc.2015.01.024. URL <http://arxiv.org/abs/1410.3012>.
- [72] S. J. Taylor and B. Letham. Forecasting at scale. doi: 10.7287/peerj.preprints.3190v2. URL <https://peerj.com/preprints/3190>.
- [73] i. team. Iminuit – A python interface to minuit. URL <https://github.com/scikit-hep/iminuit>.
- [74] R. Tibshirani. Regression Shrinkage and Selection via the Lasso. 58(1):267–288. ISSN 0035-9246. URL [www.jstor.org/stable/2346178](http://www.jstor.org/stable/2346178).
- [75] A. Tikhonov. *On the Stability of Inverse Problems*, volume vol. 39 of *Doklady Akademii Nauk SSSR*.
- [76] S. Toghi Eshghi, A. Au-Yeung, C. Takahashi, C. R. Bolen, M. N. Nyachienga, S. P. Lear, C. Green, W. R. Mathews, and W. E. O’Gorman. Quantitative Comparison of Conventional and t-SNE-guided Gating Analyses. 10. ISSN 1664-3224. doi: 10.3389/fimmu.2019.01194. URL <https://www.frontiersin.org/articles/10.3389/fimmu.2019.01194/full>.
- [77] d. L. M. J. van, E. C. Polley, and A. E. Hubbard. Super Learner. 6(1). ISSN 1544-6115. doi: 10.2202/1544-6115.1309. URL <https://www.degruyter.com/view/j/sagmb.2007.6.issue-1/sagmb.2007.6.1.1309/sagmb.2007.6.1.1309.xml>.
- [78] L. van der Maaten and G. Hinton. Visualizing Data using t-SNE. 9:2579–2605. ISSN 1533-7928. URL <http://www.jmlr.org/papers/v9/vandermaaten08a.html>.
- [79] F. van Veen. The Neural Network Zoo. URL <http://www.asimovinstitute.org/neural-network-zoo/>.
- [80] V. Vapnik. Principles of Risk Minimization for Learning Theory. In *Proceedings of the 4th International Conference on Neural Information Processing Systems, NIPS’91*, pages 831–838. Morgan Kaufmann Publishers Inc. ISBN 978-1-55860-222-9. URL <http://dl.acm.org/citation.cfm?id=2986916.2987018>.

- [81] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. Jarrod Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, I. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and S. . . Contributors. SciPy 1.0–Fundamental Algorithms for Scientific Computing in Python. page arXiv:1907.10121. URL <https://ui.adsabs.harvard.edu/abs/2019arXiv190710121V/abstract>.
- [82] I. Wallach and R. Lilien. The protein–small-molecule database, a non-redundant structural resource for the analysis of protein-ligand binding. 25(5):615–620. ISSN 1367-4803. doi: 10.1093/bioinformatics/btp035. URL <https://academic.oup.com/bioinformatics/article/25/5/615/183421>.
- [83] S. Weinberg. A Model of Leptons. 19(21):1264–1266. doi: 10.1103/PhysRevLett.19.1264. URL <https://link.aps.org/doi/10.1103/PhysRevLett.19.1264>.
- [84] H. Wickham. Tidy data. 59(10):1–23. ISSN 1548-7660. doi: 10.18637/jss.v059.i10. URL <https://www.jstatsoft.org/v059/i10>.
- [85] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. Salakhutdinov, and A. Smola. Deep Sets. URL <http://arxiv.org/abs/1703.06114>.

# *Index*

license, [ii](#)