

1 Spis treści

2	1 Fizyka	2
3	1.1 Plazma kwarkowo-gluonowa	2
4	1.2 Dżety	2
5	1.3 Identyfikacja dżetów b	3
6	1.4 Eksperyment ALICE	4
7	2 Uczenie maszynowe	5
8	2.1 Wzmacniane drzewa decyzyjne	5
9	2.2 Sieci neuronowe	6
10	2.3 Dyskusja użycia dwóch algorytmów	11
11	3 Dane	12
12	4 Analiza	14

1 Fizyka

1.1 Plazma kwarkowo-gluonowa

Plazma kwarkowo-gluonowa (ang. *quark-gluon plasma* – *QGP*) to stan materii, który istniał w pierwszych ułamkach sekund po Wielkim Wybuchu. Przewiduje się, że materia w takim stanie obecna jest także w jądrach gwiazd neutronowych [1].

Obecnie aby uzyskać dostęp do materii w stanie plazmy kwarkowo-gluonowej potrzebne są wysokoenergetyczne zderzenia cząstek. Powszechnie mówi się o niej w kontekście zderzeń ciężkich jonów, chociaż istnieją także prace doszukujące się obecności *QGP* w mniejszych systemach np. w zderzeniach proton-proton [2], [3].

QGP jest stanem o ekstremalnej gęstości i temperaturze. Jego cechą charakterystyczną jest obecność wolnych kwarków i gluonów (zbiorczo nazywanych partonami). W każdym innym stanie materii są one zawsze związane i tworzą hadrony. Zjawisko to, zwane uwięzieniem koloru (ang. *color confinement*) uniemożliwia obserwację cząstek obdarzonych ładunkiem kolorowym (a takimi są kwarki i gluony) w stanie wolnym. Wolne kwarki i gluony powstające w zderzeniach muszą zatem przejść przez proces hadronizacji, w którym rekombinują one ze spontanicznie wytwarzanymi nowymi partonami, tworząc hadrony. W wyniku tego procesu, z każdego partonu obecnego w początkowym etapie zderzenia może powstać wiele cząstek poruszających się podobnym kierunkiem, tworząc stożek z wierzchołkiem blisko punktu interakcji wiązek. Taki stożek skolimowanych cząstek nazywany jest dżetem cząstek.

1.2 Dżety

Przedstawiona powyżej definicja dżetu nie jest precyzyjna z punktu widzenia pracy eksperymentalnej. W detektorze obserwuje się tylko cząstki w stanie końcowym, nie jest znana natomiast ich historia (tj. parton, z którego powstały), nie jest zatem możliwe przyporządkowanie cząstki według jej pochodzenia. W związku z tym, konieczne jest użycie algorytmu klasteryzującego, dostającego na wejściu tylko obserwowalne eksperymentalnie cząstki. To jakie dżety zostaną zaobserwowane w danym zdarzeniu zależy od użytego algorytmu. Oznacza to, że precyzyjną definicję dżetu stanowi algorytm klasteryzujący wraz z zestawem parametrów. Obecnie najpowszechniej stosowanym algorytmem jest algorytm *anti-kt* [4].

Eksperymentalne ograniczenia związane z obserwacją tylko końcowego stanu oddziaływań nie występują w analizie danych z symulacji Monte Carlo (MC), gdzie ma się dostęp do pełnej informacji na temat historii każdej cząstki. Można pomyśleć, że daje to możliwość lepszej klasteryzacji dżetów. Rekonstrukcja przy użyciu informacji o partonach - matkach sprawia jednak, że definicja dżetu wykorzystana w badaniach danych symulacyjnych jest inna od tej wykorzystywanej w eksperymencie. Tracona jest przez to cecha odpowiedniości między obiektami nazywanymi dżetami w symulacji i w eksperymencie, która to cecha jest niewątpliwie jedną z podstawowych wymagań stawianych przed dobrą symulacją.

Dżety stanowią ważny element w badaniach plazmy kwarkowo-gluonowej. Dają one pośredni wgląd we właściwości *QGP* na podstawie jej wpływu na oddziałujące z nią partony. Przykładową obserwabłą mierzoną w zderzeniach ciężkich jonów jest czynnik modyfikacji jądrowej (ang. *nuclear modification factor*), który jest miarą strat energii przez parton przechodzący przez medium.

Oprócz globalnego wpływu medium na dżety, analizuje się także różnice między dżetami pochodzącymi z gluonów oraz kwarków o różnych zapachach (ang. *flavours*). Modele teoretyczne przewidują między innymi większe straty energii w wyniku interakcji z *QGP* dla dżetów gluonowych niż kwarkowych [5] oraz zależność strat energii od masy partonu [6] – w tym przypadku

58 precyzyjne pomiary rozróżniające typy dżetów pozwalają lepiej zrozumieć mechanizm odpo-
59 wiadający za straty energii przez partony. Zagadnienie rozpoznania z jakiego rodzaju partonu
60 powstał dany dżet, nazywane jest identyfikacją lub tagowaniem dżetu.

61 1.3 Identyfikacja dżetów b

62 Poza badaniami właściwości QGP , szczególne znaczenie ma identyfikacja dżetów pochodzących
63 z ciężkich kwarków: b i c . Są one ważnym elementem w poszukiwaniu łamania symetrii CP w
64 rozpadach hadronów B i D oraz innych sygnatur tzw. *Nowej Fizyki* wykraczającej poza ramy
65 Modelu Standardowego. Kwarki *piękne* pojawiają się także często w kanałach rozpadu cząstek
66 takich jak bozon Higgsa i kwark t .

67 Identyfikacja dżetów b jest sporym wyzwaniem ze względu na zdecydowanie częściej wystę-
68 pujące dżety lekkie, tj. powstałe z hadronizacji kwarków u, d, s oraz gluonów. Rozpoznawanie
69 dżetów b bazuje na charakterystycznych właściwościach hadronów zawierających kwark piękny:
70 relatywnie długim czasie życia oraz (w mniejszym stopniu) na ich pół-leptonowych rozpadach
71 o względnej częstości rozpadu w tym kanale (ang. *branching ratio*) na poziomie 10%.

72 przegląd algorytmów używanych w identyfikacji b-jetów: TBD

73 Używane w eksperymentach na LHC: ATLAS [7], [8], CMS i ALICE algorytmy można
74 podzielić na trzy kategorie: wykorzystujące wtórne wierzchołki, informację o odległości najbliż-
75 szego zbliżenia (parametrach zderzenia) (ang. *Distance of Closest Approach – DCA, Impact*
76 *Parameter – IP*) cząstek tworzących dżet oraz identyfikujące produkty półleptonowych rozpa-
77 dów pięknych lub powabnych hadronów.

78 Najprostszym algorytmem jest dyskryminacja na podstawie istotności statystycznej odległo-
79 ści wtórnego wierzchołka od pierwotnego L (wynik pomiaru podzielony przez jego niepewność).
80 Jest to metoda wykorzystywana w każdym z trzech wymienionych eksperymentów (por. ALI-
81 CE, ATLAS: algorytm SV0, CMS: algorytm SSV). Może być ona rozszerzona poprzez użycie
82 dodatkowych zmiennych opisujących wtórny wierzchołek jak na przykład jego masa lub uła-
83 mek niesionej przez wtórny wierzchołek całkowitej energii dżetu (por. ATLAS: SV1) lub użycie
84 "pseudowierzchołków" (kombinacji dwóch cząstek o dużych DCA) w celu poprawienia wydajno-
85 ści detekcji o przypadki, w których wtórny wierzchołek nie został zrekonstruowany (por. CMS:
86 CSV).

87 Algorytmy wykorzystujące informację o poszczególnych cząstkach mogą zasadniczo bazować
88 albo na sumie logarytmów prawdopodobieństw pochodzenia każdej cząstki z pierwotnego wierz-
89 chołka (por. ATLAS: IP3D, CMS: JP) lub tym samym prawdopodobieństwie ale dla wybranej,
90 np. drugiej lub trzeciej cząstki na liście posortowanej według malejącego IP (por. ALICE i
91 CMS: TC). Bardziej złożonym podejściem, w którym cząstki nie są traktowane jako niezależne,
92 jest użycie rekurencyjnych sieci neuronowych (ATLAS: RNNIP).

93 Do wykorzystania półleptonowego kanału rozpadu ciężkich hadronów do identyfikacji dżetów
94 b w eksperymentach ATLAS (SMT) i CMS użyto wzmacnianych drzew decyzyjnych trenowa-
95 nych na kilku ręcznie zdefiniowanych w tym celu zmiennych takich jak pęd leptonu transwer-
96 salny względem osi dżetu.

97 Znaczącą poprawę zdolności predykcyjnej można uzyskać łącząc kilka różnych modeli. Algo-
98 rytmy łączące może pobierać na wejściu albo tylko predykcjach klasyfikatorów niższego poziomu
99 (CMS: cMVA2) lub dodatkowo także ich zmienne wejściowe (ATLAS: MV2, DL1). Do scalania
100 używane są zwykle wzmacniane drzewa decyzyjne lub sieci neuronowe.

101 Przykład innego podejścia zaprezentowała współpraca przy eksperymencie LHCb, gdzie
102 wykorzystano dwa zestawy wzmacnianych drzew decyzyjnych operujących na zmiennych zwią-
103 zanych z wtórnymi wierzchołkami. Pierwszy zapewnia separację dżetów lekkich od ciężkich a
104 drugi odróżnia dżety b od c . Do wyboru punktu pracy, zamiast jednowymiarowego rozkładu
105 predykcji używany jest dwuwymiarowy rozkład wag przypisany przez oba klasyfikatory.

¹⁰⁶ **1.4 Eksperyment ALICE**

¹⁰⁷ TBD

109 Uczenie maszynowe jest bardzo szerokim i obecnie dynamicznie się rozwijającym obszarem
110 wiedzy. Występuje w wielu odmianach łącząc w sobie w zależności od wariantu wiele dziedzin
111 takich jak matematyka (statystyka, algebra) informatyka (algorytmika, teoria informacji) a
112 także elementy robotyki i sterowania. Dziedzinami, w których jest najczęściej wykorzystywane
113 są min. widzenie maszynowe, przetwarzanie języka naturalnego, autonomiczne roboty i pojazdy,
114 systemy decyzyjno - eksperckie, optymalizacyjne oraz rekomendacyjne.

115 W tej pracy wykorzystywana jest gałąź uczenia maszynowego nazywana uczeniem nadzoro-
116 wanym lub "uczeniem z nauczycielem" (ang. *supervised learning*), gdzie uczenie występuje na
117 podstawie poprawnie oznaczonych przykładów. Terminami bliskoznacznymi dla tak rozumia-
118 nego uczenia maszynowego są uczenie statystyczne (ang. *statistical learning*) i rozpoznawanie
119 wzorców (ang. *pattern recognition*).

120 Problem identyfikacji dżetów jest klasycznym przykładem zagadnienia klasyfikacji, gdzie
121 poprawna odpowiedź jest jedną ze skończonej ilości opcji (klas) w przeciwieństwie do regresji,
122 gdzie szukana odpowiedź algorytmu ma charakter ciągły.

123 Występuje wiele algorytmów uczenia maszynowego takich jak regresja liniowa i logistyczna,
124 drzewa decyzyjne i ich wariacje, maszyny wektorów wspierających, sztuczne sieci neuronowe
125 oraz wiele innych. Uczenie polega na znalezieniu pewnej funkcji dopasowującej do przyjmowa-
126 nego na wejściu zestawu (wektora) cech (zmiennych, kolumn) pewną odpowiedź (predykcję),
127 która minimalizuje zadaną funkcję straty. Jej rolę w przypadku regresji często pełni błąd śred-
128 niokwadratowy a w przypadku klasyfikacji np. entropia krzyżowa (ang. *cross entropy*)¹. Różne
129 algorytmy szukają przy tym funkcji dopasowującej należącej do różnych klas funkcji: przy-
130 kładowo klasyczne drzewa decyzyjne przeszukują tylko przestrzeń funkcji dających się opisać
131 skończonym zbiorem reguł "jeśli – to" (ang. *if – else*).

132 W pracy wykorzystane zostały dwa rodzaje algorytmów: wzmacniane drzewa decyzyjne oraz
133 sieci neuronowe.

134 **2.1 Wzmacniane drzewa decyzyjne**

135 Wzmacniane drzewa decyzyjne są jednym z rozwinięć klasycznego algorytmu drzewa decyzyj-
136 nego. Pojedyncze drzewo decyzyjne dzieli przestrzeń cech uczących przy pomocy prostopadłych
137 cięć, na mniejsze/większe niż zadaną wartość w przypadku zmiennej ciągłej lub na należące/nie
138 należące do danej klasy w przypadku zmiennej kategorycznej. Każdy podział, nazywany wę-
139 złem, daje dwie gałęzie, które można dalej niezależnie dzielić aż do ostatniego poziomu (liści).
140 Kolejne podziały wybierane są tak, aby zbiory przykładów wpadające do poszczególnych gałęzi
141 były jak najbardziej jednorodne. Stosuje się różne miary nieporządku takie jak: *indeks Gini*
142 lub entropia.

143 Drzewa decyzyjne są często łączone w komitety klasyfikatorów (ang. *ensemble methods*).
144 Wiele "słabych" klasyfikatorów jest łączonych w jeden "silny" na dwa sposoby: *bagging* oraz
145 *boosting* (wzmacnianie), które są często ze sobą porównywane.

146 *Bagging* – w zastosowaniu dla drzew decyzyjnych nazywany algorytmem lasów losowych
147 (ang. *random forest*) - polega na wytrenowaniu wielu drzew, każdego na podstawie N przykła-
148 dów losowo wylosowanych z powtórzeniami spośród N -licznego zbioru treningowego. Dodatko-
149 wo, do uczenia każdego drzewa używa się tylko podzbioru wszystkich cech uczących. Końcową
150 predykcję algorytmu otrzymuje się poprzez "głosowanie" wszystkich drzew z odpowiednimi

¹ $J = -\frac{1}{N} \sum_{i=1}^N [y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)]$, gdzie y_i to prawidłowa klasa i -tego przykładu a \hat{y}_i to predykcja algorytmu

151 wagami.

152 *Boosting* – wzmacniane drzewa decyzyjne (ang. *boosted decision trees*) – jest metodą po-
153 dobłą do *baggingu*. Główną różnicą jest zwiększanie wag przykładom uczącym, które przez po-
154 przednie drzewo zostały źle zaklasyfikowane – każde kolejne drzewo koncentruje się bardziej na
155 poprawie błędów poprzednich drzew. Widać tu kolejną ważną cechę odróżniającą obie metody:
156 *boosting* jest algorytmem sekwencyjnym podczas gdy *bagging* daje się trywialnie zrównoleglić
157 (każde drzewo trenowane jest w osobnym wątku).

158 Parametry i sposób trenowania drzew decyzyjnych na analizowanych danych

159 W niniejszej pracy wykorzystano wzmacniane drzewa decyzyjne zaimplementowane w wydaj-
160 nej bibliotece **XGBoost** [9]. Szybkość obliczeń jest bardzo ważna, gdyż oprócz komfortu pracy
161 z algorytmem, przekłada się na jakość otrzymanych wyników – krótszy czas obliczeń oznacza
162 możliwość przeprowadzenia większej ilości eksperymentów i lepsze dobranie parametrów oraz
163 danych. Implementacja wzmacnianych drzew decyzyjnych w **XGBoost** wykorzystuje wszystkie
164 rdzenie procesora, pomimo że sam algorytm ma charakter sekwencyjny – jest to możliwe dzięki
165 paralelizacji procesu tworzenia każdego drzewa (przed każdym podziałem konieczne jest sprawd-
166 zenie pewnej ilości możliwych zmiennych i wartości progowych i ten proces jest wykonywany
167 równolegle).

168 Dzięki szybkiemu uczeniu się algorytmu, możliwe było użycie kosztownego obliczeniowo au-
169 tomatycznego przeszukiwania przestrzeni parametrów przy pomocy przeszukiwania losowego
170 (ang. *random search*), które jest zwykle preferowane nad przeszukiwanie sieciowe [10]. W tym
171 celu cały zbiór danych dzielony był na dwie części: trenującą oraz testową (80/20%). Następ-
172 nie algorytm był trenowany i oceniany z użyciem trzy- lub pięciokrotnej walidacji krzyżowej
173 (ang. *cross-validation*) na zbiorze trenującym dla różnych zestawów parametrów. Model z naj-
174 lepszym wynikiem uzyskanym w walidacji krzyżowej był sprawdzany na zbiorze testowym.

175 Parametry optymalizowane w opisanym procesie to:

- 176 • *max_depth* – maksymalna głębokość każdego drzewa (niekoniecznie osiągnana)
- 177 • *n_estimators* – liczba drzew
- 178 • *learning_rate* – parametr szybkości uczenia, komplementarny do *n_estimators*, w prak-
179 tyce można ustalić liczbę drzew i szukać optymalnej szybkości uczenia
- 180 • *subsample*, *colsample_bytree*, *colsample_bylevel* – parametry regularyzacyjne określają-
181 ce ułamek kolejno: danych użytych do trenowania każdego drzewa, kolumn użytych w
182 każdym drzewie (cechy losowane raz dla danego drzewa), kolumn użytych przy każdym
183 podziale (cechy losowane przy każdym podziale)
- 184 • γ – minimalny zysk w postaci zmniejszenia wartości funkcji straty konieczny do wykonania
185 podziału

186 2.2 Sieci neuronowe

187 Sieci neuronowe (ang. *neural networks* – *NN*) są szczególnym algorytmem uczenia maszyno-
188 wego. Występują w bardzo wielu odmianach i są wykorzystywane w rozwiązywaniu szerokiej
189 gamy problemów. Nawet bardzo pobieżny opis sieci neuronowych wymaga dużo więcej miejsca
190 niż może być temu poświęcone w tej pracy. Wprowadzenia do sieci neuronowych na różnym
191 poziomie zaawansowania można znaleźć m.in. w [REF]x100. Tu przedstawione zostaną wy-
192 łącznie wybrane zagadnienia mające ściślejszy związek z pracą. Używane mogą być terminy,
193 których znaczenie wyjaśniane jest w podanych źródłach.

194 W niniejszej pracy, wykorzystane zostały dwa rodzaje sieci neuronowych: sieci w pełni połą-
195 czone (ang. *fully connected NN – FC NN*), nazywane także wielowarstwowymi perceptronami
196 (ang. *multi-layer perceptron – MLP*) oraz sieci konwolucyjne (ang. *convolutional NN – Co-
197 nvNets, CNN*).

198 Sieci w pełni połączone

199 W nierekurencyjnych sieciach neuronowych (tylko takie są używane w tej pracy), informacja
200 jest przekazywana kolejno od warstw wejściowych, poprzez warstwy ukryte aż do wyjściowej. W
201 sieciach typu *FC* wszystkie warstwy składają się z identycznych neuronów – każdy neuron do-
202 staje na wejściu wektor, natomiast zwraca skalar – wartość pewnej zadanej, nieliniowej funkcji,
203 jako argument podając średnią ważoną z elementów wektora wejściowego. Wartości zwraca-
204 ne przez neurony w danej warstwie składają się na wektor wejściowy dla neuronów kolejnej
205 warstwy. Wejściami dla pierwszej warstwy są natomiast kolejne przykłady ze zbioru uczącego.
206 Trenowanie sieci neuronowych polega na zmienianiu wag (parametrów) w liczonej w każdym
207 neuronie średniej, każdy neuron posiada własny, niezależny zestaw wag.

208 Istnieje twierdzenie o sieciach neuronowych jako uniwersalnych aproksymatorach funkcji
209 (ang. *universal approximation theorem*) [11], mówiące, że już sieć neuronowa o jednej warstwie
210 ukrytej jest zdolna do przybliżenia dowolnej funkcji z dowolną dokładnością. Twierdzenie to
211 nie podaje niestety liczby potrzebnych neuronów a przede wszystkim – sposobu ich trenowania.
212 Trenowanie jest prostsze w przypadku zastosowania wielu warstw, które odpowiadają kolejnym
213 poziomą abstrakcji jednak nadal jest dużym wyzwaniem ze względu na fakt, że nawet stosun-
214 kowo niewielka sieć może posiadać bardzo dużą liczbę parametrów, przykładowo sieć o czterech
215 warstwach, w każdej po 128 neuronów ma ich ponad 65 tysięcy.

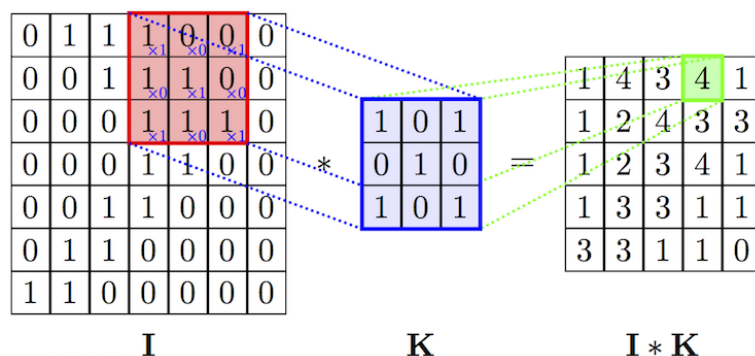
216 Sieci konwolucyjne

217 Jednym ze sposobów na ograniczenie liczby trenowanych parametrów jest użycie konwolucyj-
218 nych sieci neuronowych (bardziej poprawną choć rzadko używaną nazwą w języku polskim jest
219 sieć splotowa). Są one inspirowane połączeniami w korze wzrokowej zwierząt i wywodzą się z ba-
220 dań w obszarze widzenia komputerowego, gdzie liczby parametrów są szczególnie duże (wektor
221 wejściowy ma wymiar równy liczbie pikseli w obrazie), na takim przykładzie również najłatwiej
222 zrozumieć ich działanie.

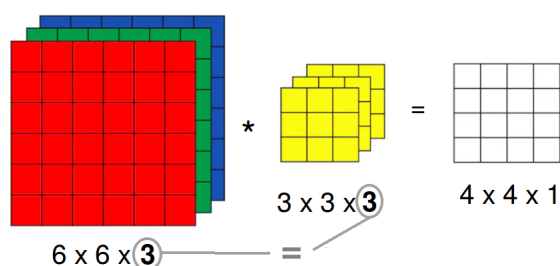
223 Sieci konwolucyjne różnią się od sieci typu *FC* tym, że część wag połączeń między warstwami
224 jest dzielona. Występuje w nich nowy rodzaj warstwy, nazywany warstwą konwolucyjną. Każda
225 jednostka w warstwie konwolucyjnej (filtr) ma pewną stałą (niewielką) liczbę wag. Połączenie z
226 dużym wejściem realizowane jest przez powielanie tych samych wag w połączeniach z kolejnymi
227 fragmentami wektora wejściowego (por. Rys. 1). Rezultatem działania filtra na macierz jest wy-
228 nik operacji splotu. Liczba parametrów przypadająca na każdy filtr jest równa jego rozmiarowi
229 i nie zależy od wielkości wektora wejściowego.

230 W przypadku gdy zamiast wejścia dwuwymiarowego (jak np. obraz czarno-biały), mamy do
231 czynienia z wejściem trójwymiarowym (np. trzeci wymiar to kolejne kolory w kodowaniu RGB),
232 filtry również muszą mieć trzy wymiary, przy czym rozmiar w ostatnim wymiarze musi być
233 równy rozmiarowi w tym kierunku wektora wejściowego. Wynik operacji splotu jest ponownie
234 dwuwymiarowy, gdyż filtr przesuwany jest tylko w dwóch pierwszych wymiarach. Trzeci wymiar
235 powstaje przez składanie kolejnych filtrów. Widać zatem, że również w przypadku gdy na
236 wejścia podawany jest obraz czarno-biały, filtry w kolejnych warstwach konwolucyjnych (oprócz
237 pierwszej) mają po trzy wymiary.

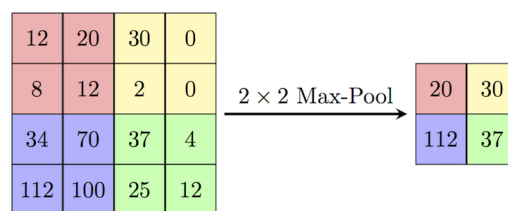
238 Oprócz warstw konwolucyjnych, w sieciach tego typu stosowane są także tzw. warstwy typu
239 *max-pooling*. Zasada jej działania jest bardzo prosta: wykonuje funkcję *maksimum* na zadanym



Rysunek 1: Schemat działania pojedynczego filtra z warstwy konwolucyjnej (operacja splotu).



Rysunek 2: Działanie pojedynczego filtra (3D) na wejście o trzech wymiarach.



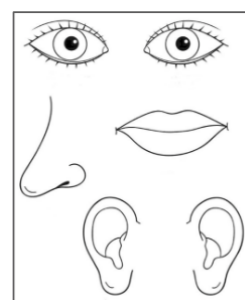
Rysunek 3: Działanie warstwy typu *max-pool*

240 fragmencie obrazu (por. Rys. 3). Ich rolą jest zmniejszanie rozmiaru przekazywanej w sieci
241 informacji.

242 Typowa architektura stosowana w przypadku sieci konwolucyjnych jest następująca: naj-
243 pierw warstwy konwolucyjne (pomiędzy nimi czasem warstwy typu *max-pool*), następnie wszyst-
244 kie filtry są rozwijane i składane w długi jednowymiarowy wektor, który przekazywany jest do
245 warstw typu *FC*. W przypadku problemu klasyfikacji, na końcu znajduje się jeszcze warstwa
246 typu *softmax* normalizująca wyjście z sieci do jedynki (por. Rys. 4).

247 Sieci konwolucyjne posiadają dwie właściwości odróżniające je od
248 *MLP*:

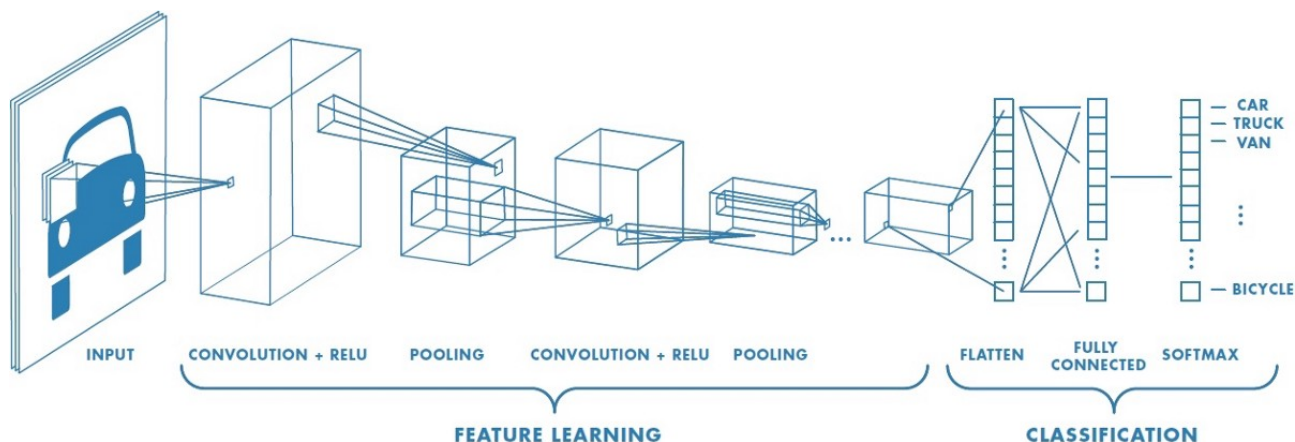
- 249 • niezmienniczość względem przesunięcia (ang. *translation invariance*) – głównie za sprawą dzielenia wag oraz obecności warstw
250 typu *max-pool*, położenie danej cechy na obrazie jest niemal
251 bez znaczenia (obraz po prawej stronie byłby rozpoznany ja-
252 ko twarz)
253
- 254 • lokalność połączeń – filtry obejmują tylko kilka sąsiednich pikseli
255 (tam zwykle występują najsilniejsze zależności) nie są w stanie
256 dostrzec cechy rozciągniętej na obszar większy od rozmiaru filtra



257 Hiperparametry i trenowanie sieci neuronowych na analizowanych danych

258 Parametry sieci, których wartości są określane przez projektanta sieci, takie jak liczba warstw
259 ukrytych są nazywane hiperparametrami (dla odróżnienia od parametrów - wag połączeń).

260 Testowane były trzy architektury: sieci w pełni połączone, sieci konwolucyjne oraz sieć
261 złożona z dwóch gałęzi, osobnych dla wtórnych wierzchołków i cząstek tworzących dżet (por.
262 Rozdz. 3) przedstawione schematycznie na Rys. 5. Do trenowania sieci wykorzystano wysoko-

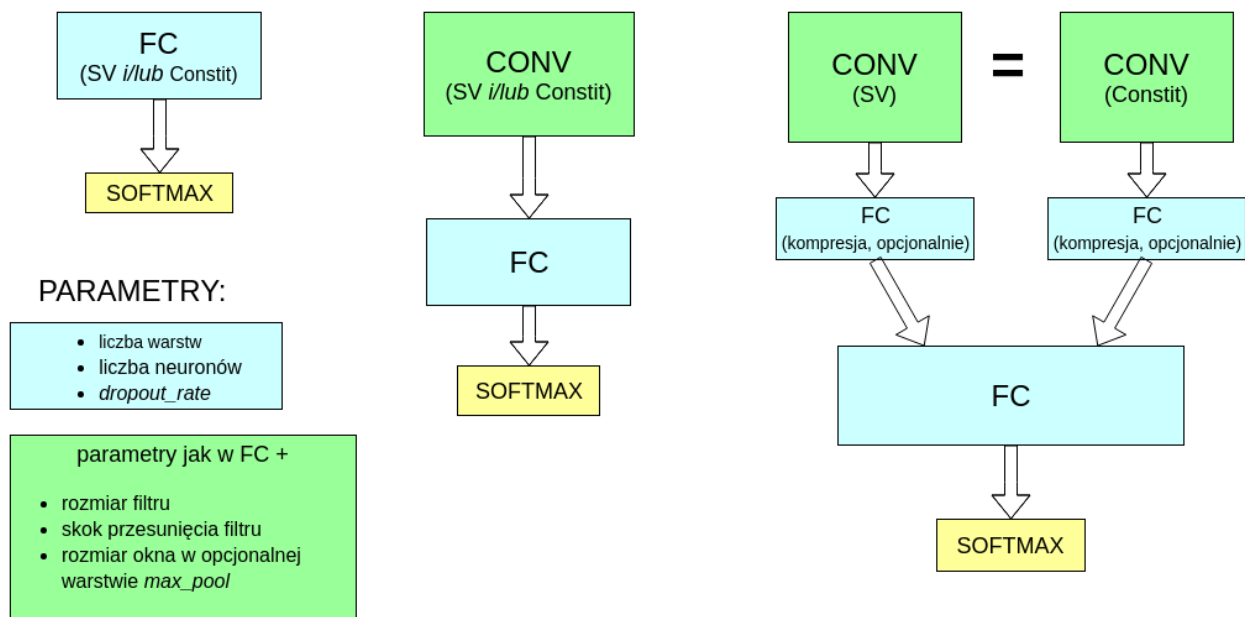


Rysunek 4: Typowa struktura stosowana w sieciach konwolucyjnych. Warstwy konwolucyjne mają za zadanie wydobywać cechy, na podstawie których późniejsze warstwy dokonują klasyfikacji. Widoczna jest charakterystyczna stopniowa zmiana rozmiaru przekazywanej macierzy: rozmiar poprzeczny maleje kosztem głębokości, co odpowiada rosnącej liczbie filtrów i malejącemu rozmiarowi obszaru po jakim są one przesuwane.

263 poziomową bibliotekę `Keras` [12] korzystającą z silnika obliczeniowego zaimplementowanego w
 264 `TensorFlow` [13].

265 Zestaw hiperparametrów definiujący działanie sieci w pełni połączonej:

- 266 • liczba warstw ukrytych
- 267 • liczba neuronów w każdej warstwie
- 268 • funkcja aktywacji – nieliniowa funkcja aplikowana przed zwróceniem wartości w każdym
 269 neuronie, najpopularniejsze to \tanh , $ReLU$ ($f(x) = \max(0, x)$) oraz funkcja sigmoidalna
 270 ($f(x) = \frac{1}{1+\exp(x)}$)
- 271 • algorytm optymalizacyjny – spadek gradientowy lub jego wariacje
- 272 • parametr szybkości uczenia i jego modyfikacje w trakcie uczenia
- 273 • liczba przykładów trenujących przetwarzanych w jednym kroku uczenia (ang. *batch_size*)
 274 – im większy tym szybsze jest trenowanie sieci (dzięki wydajnym operacjom macierzowym), natomiast może się to odbywać kosztem precyzji
- 275
- 276 • liczba epok uczenia – ile razy będzie pokazywany sieci każdy przykład
- 277 • opcjonalnie: warunki stopu (ang. *early stopping*) w razie osiągnięcia *plateau* (wysycenia
 278 procesu uczenia)
- 279 • opcjonalnie: regularyzacja przy pomocy różnych technik (zwykle konieczna)
- 280 • ponadto dla sieci konwolucyjnych:
 - 281 – liczba warstw konwolucyjnych i liczba filtrów w każdej warstwie
 - 282 – obecność lub brak warstw *max-pool* i rozmiar ich okna
 - 283 – rozmiar filtrów i długość skoku przy ich przesuwaniu



Rysunek 5: Schematyczne przedstawienie trzech testowanych rodzin architektur sieci. Każdy blok odpowiada kilku warstwom danego typu. Bloki warstw typu *FC* i opisane jako "kompresja" składały się z kilku neuronów i miały za zadanie zredukować całą informację z danej gałęzi do wektora kilku liczb. Obie gałęzie miały taką samą strukturę.

Same dwie pierwsze wielkości dają nieograniczoną liczbę konfiguracji. Czas trenowania sieci neuronowych jest rzędu wielkości większy niż drzew decyzyjnych, dlatego przyjęto szereg kroków mających na celu zmniejszenie przeszukiwanej przestrzeni hiperparametrów. Na podstawie wstępnych testów oraz różnych wskazówek dostępnych w literaturze przyjęto:

- batch_size* zawsze równy 64 (inne testowane wartości: 16, 32, 128)
- za algorytm optymalizacyjny przyjęto algorytm o nazwie *Nadam* [14], tj. rozwinięcie algorytmu *Adam* [15] o parametr Nesterova (inne testowane to zwykły spadek gradientowy oraz *Adam*)
- funkcję aktywacji: *ReLU*
- liczba epok równa 100 lub 200 (lub mniej do testów), zrezygnowano z *early stopping*
- stałe w trakcie treningu wartości parametru szybkości uczenia
- spośród technik regularyzacyjnych testowano wyłącznie *dropout* [16] z parametrami 0.1, 0.2, 0.5
- kilka wybranych kombinacji dla zestawu parametrów: rozmiar filtra, długość skoku i rozmiar okna w warstwach *max-pool* – takie same w kolejnych warstwach
- liczby neuronów/filtrów w warstwach będące zawsze potęgami dwójki oraz stałą liczbę w kolejnych warstwach lub zmieniającą się o stały czynnik, np. 256-128-64, 128-128-128 lub 16-32-64
- liczba warstw *FC*: 2-8, konwolucyjnych 2-6

303 Nawet po przyjęciu powyższych uproszczeń nie sposób sprawdzić wszystkich możliwych
304 zestawów hiperparametrów, dlatego sposób ich dobierania w kolejnych testach był mocno em-
305 piryczny. Dostępne dane dzielone były na trzy zbiory: trenujący, walidacyjny i testowy. Wobec
306 braku warunków stopu, zbiór walidacyjny użyty był wyłącznie do porównywania różnych ze-
307 stawów parametrów, tak aby wynik testowy pozostał nieobciążony.

308 Zgodnie z zasadą ortogonalizacji działań, proces doboru hiperparametrów dzielono na dwie
309 części: najpierw starano się uzyskać jak najlepsze wyniki na zbiorze uczącym, a dopiero później
310 zmusić algorytm do lepszej generalizacji na zbiorze testowym przez zwiększoną regularyzację i
311 modyfikację parametru szybkości uczenia.

312 2.3 Dyskusja użycia dwóch algorytmów

313 Użycie więcej niż jednego algorytmu ma wiele zalet. Po pierwsze daje możliwość porównania
314 wyników. Pozwala to na oszacowanie błędu *Bayesowskiego* (najniższego możliwego do osiągnię-
315 cia przez jakikolwiek algorytm błędu). Jest to bardzo ważne w sytuacji, gdy nie dysponuje
316 się innym oszacowaniem tego błędu (w wielu problemach naturalnych dla człowieka jak roz-
317 poznawanie obiektów na obrazkach jest nim błąd ludzki lub też błąd popełniany przez zespół
318 ekspertów w bardziej zaawansowanych zastosowaniach).

319 Po drugie, wykorzystane zostały dwa algorytmy mocno różniące się w swojej naturze, co po-
320 zwala wykorzystać cechy każdego z nich w analizie: przykładowo sieci neuronowe dobrze radzą
321 sobie z nieustrukturyzowanymi danymi – potrafią tworzyć wysoko poziomowe cechy na podsta-
322 wie niskopoziomowego wejścia (np. położenia oka na zdjęciu twarzy na podstawie pixeli). Są
323 natomiast trudne w interpretacji i często traktowane są jako tzw. "czarne skrzynki" (ang. *black*
324 *box*). Oprócz tego, liczba możliwych konfiguracji sieci jest ogromna i przez to niemożliwe jest
325 stwierdzenie czy wykorzystane zostały pełne ich możliwości.

326 Z kolei drzewa decyzyjne posiadają stosunkowo niewielką liczbę parametrów, a ich trenowa-
327 nie jest bardzo szybkie co pozwala na ich ekstensywne przeszukiwanie i otrzymanie wyników,
328 które można uznać za optymalne dla tego algorytmu. Ponadto, w przypadku drzew istnieją
329 niewymagające dodatkowych obliczeń miary użyteczności poszczególnych zmiennych, co daje
330 wgląd w działanie algorytmu i poprawia intuicyjne zrozumienie jego predykcji.

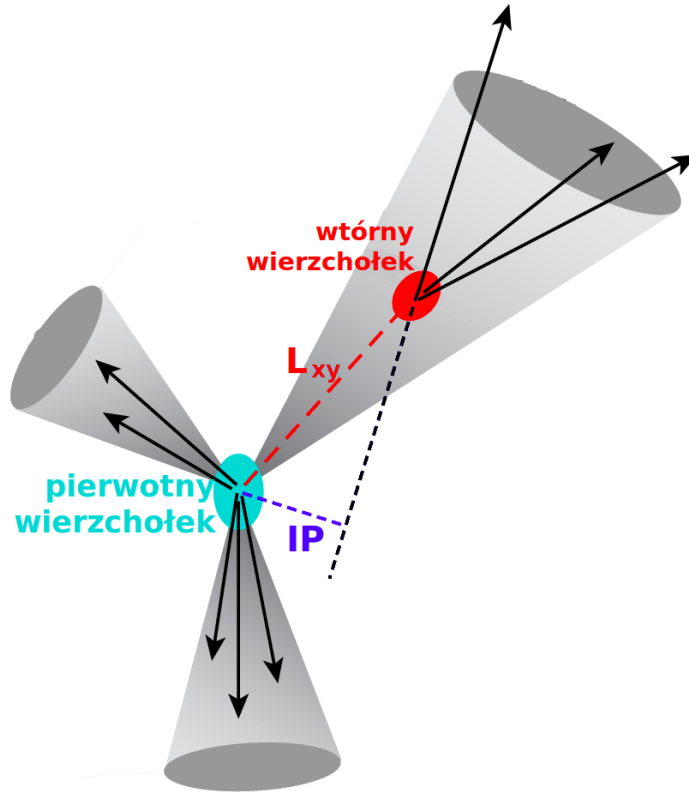
3 Dane

Dane użyte w analizie pochodzą z symulacji Monte Carlo zderzeń proton-proton przy energii w układzie środka masy równej $\sqrt{s} = 13$ TeV dostępnych na serwerach eksperymentu ALICE. W symulacjach zderzeń użyty był symulator `Pythia8` (*tune: Pythia8Jets_Monash2013*), następnie przy pomocy `Geant3` symulowany był transport cząstek i odpowiedź detektora ALICE.

Do rekonstrukcji dżetów wykorzystany został algorytm *anti-kt* z parametrem $R = 0.4$ zaimplementowany w pakiecie `FASTJET`. Dżetów poszukiwano tylko wyłącznie wśród cząstek naładowanych (ang. *charged jets*).

Do analizy wybrano dżety o p_T większym niż 15 GeV i mieszczących się w całości w akceptancji detektora *TPC*, tj. $|\eta| < 0.9$, co przy użytym parametrze rozmiaru dżetu $R = 0.4$, daje ograniczenie na pseudospieszość $|\eta| < 0.5$ dla osi dżetu.

Dla każdego dżetu obliczony został szereg wielkości, które można podzielić na zmienne na poziomie dżetu, związane z wtórnymi wierzchołkami oraz cząstkami tworzącymi dżet. Za potencjalne wtórne wierzchołki uznaje się wszystkie kombinacje trzech cząstek spełniających pewne dosyć luźne wymagania jak $p_T > 1\text{GeV}$ (rozważane są wyłącznie trzy-cząstkowe wtórne wierzchołki), stąd ich liczba może być dużo większa od liczby cząstek tworzących dżet.



Rysunek 6: Rysunek ilustrujący znaczenie używanych wielkości: L_{xy} oraz parametru zderzenia IP .

Lista używanych zmiennych:

- Zmienne na poziomie dżetu:

- η, ϕ – pseudospieszość (ang. *pseudorapidity*) i kąt azymutalny
- p_T – pęd poprzeczny dżetu

– masa dżetu

– powierzchnia dżetu – liczona w płaszczyźnie (η, ϕ) , do powierzchni dżetu zaliczany jest element w powierzchni w którym dodanie cząstki o nieskończenie małym pędzie poprzecznym sprawi, że zostanie ona zaliczona do tego dżetu [REF <https://arxiv.org/pdf/0707>]

– gęstość tła (w danym zdarzeniu)

– N_{SV} – liczba wtórnych wierzchołków

– $N_{Constit}$ – liczba cząstek tworzących dżet

• Zmienne opisujące cząstki tworzące dżet:

– η, ϕ – pseudopospieschność i kąt azymutalny cząstki względem osi dżetu

– p_T – pęd poprzeczny cząstki

– IP_D – rzut na kierunek poprzeczny wektora parametru zderzenia

– IP_Z – rzut na oś z wektora parametru zderzenia

• Zmienne opisujące wtórne wierzchołki:

– L_{xy} – odległość między pierwotnym a wtórnym wierzchołkiem (ang. *decay length*)

– σ_{Lxy} – niepewność wyznaczenia L_{xy}

– $\sigma_{vertex} = \sqrt{d_1^2 + d_2^2 + d_3^2}$ – rozrzut śladów (ang. *tracks*) wokół wtórnego wierzchołka, gdzie d_i to odległość najbliższego zbliżenia śladu / odległość najbliższego przelotu do wtórnego wierzchołka (ang. *distance of closest approach – DCA*)

– $M_{inv} = \sqrt{(E_1 + E_2 + E_3)^2 - (\vec{p}_1 + \vec{p}_2 + \vec{p}_3)^2}$ – masa niezmiennicza wierzchołka, gdzie E_i, p_i to energia i pęd i -tej cząstki tworzącej wierzchołek

– χ^2/Ndf dopasowania wtórnego wierzchołka

TBD: tabelaryczna postać, liczba SV i Nconstit, jak posortowane, wypieniania zerami korespondencja między użytymi danymi a algorytmami z sec. "identyfikacja dżetów b"

375 Dobór metryki

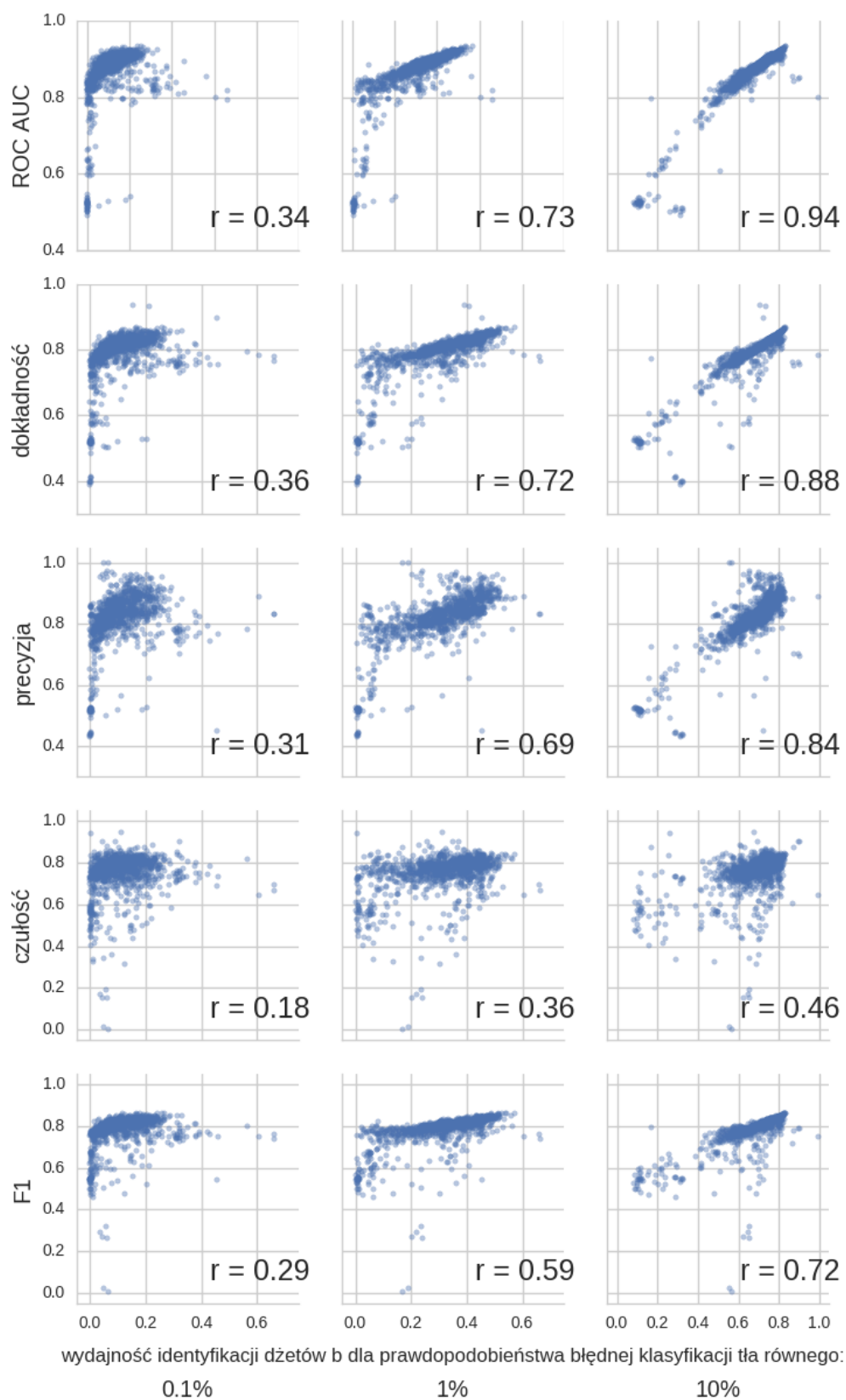
376 Bardzo ważnym elementem w trenowaniu algorytmów uczenia maszynowego jest dobór odpo-
 377 wiedniej metryki – klasycznym złym przykładem jest używanie dokładności (ang. *accuracy*) do
 378 oceniania klasyfikacji binarnej w przypadku dużego niezrównoważenia klas – algorytm przewi-
 379 dujący zawsze klasę większościową może osiągnąć dużą wartość dokładności będąc jednocześnie
 380 bardzo słabym modelem.

381 Kilka najczęściej używanych metryk wymieniono w Tab. ?? . Używanie i porównywanie
 382 kilku miar efektywności jest często niepraktyczne dlatego dobrze jest wybrać jedną metrykę.
 383 Przy jej wyborze należy kierować się potencjalnymi zastosowaniami modelu. W tym przypadku
 384 są to analizy fizyczne, które mogą mieć różne wymagania dotyczące czystości i liczebności
 385 otrzymywanych próbek a co za tym idzie, preferować inne punkty pracy zdefiniowane jako pary
 386 liczb: wydajność poprawnej klasyfikacji dżetów b (ang. *tagging efficiency = true positive rate*
 387 $= recall$), i ułamek niepoprawnie zaklasyfikowanych przypadków tła (ang. *mistagging rate =*
 388 *false positive rate*).

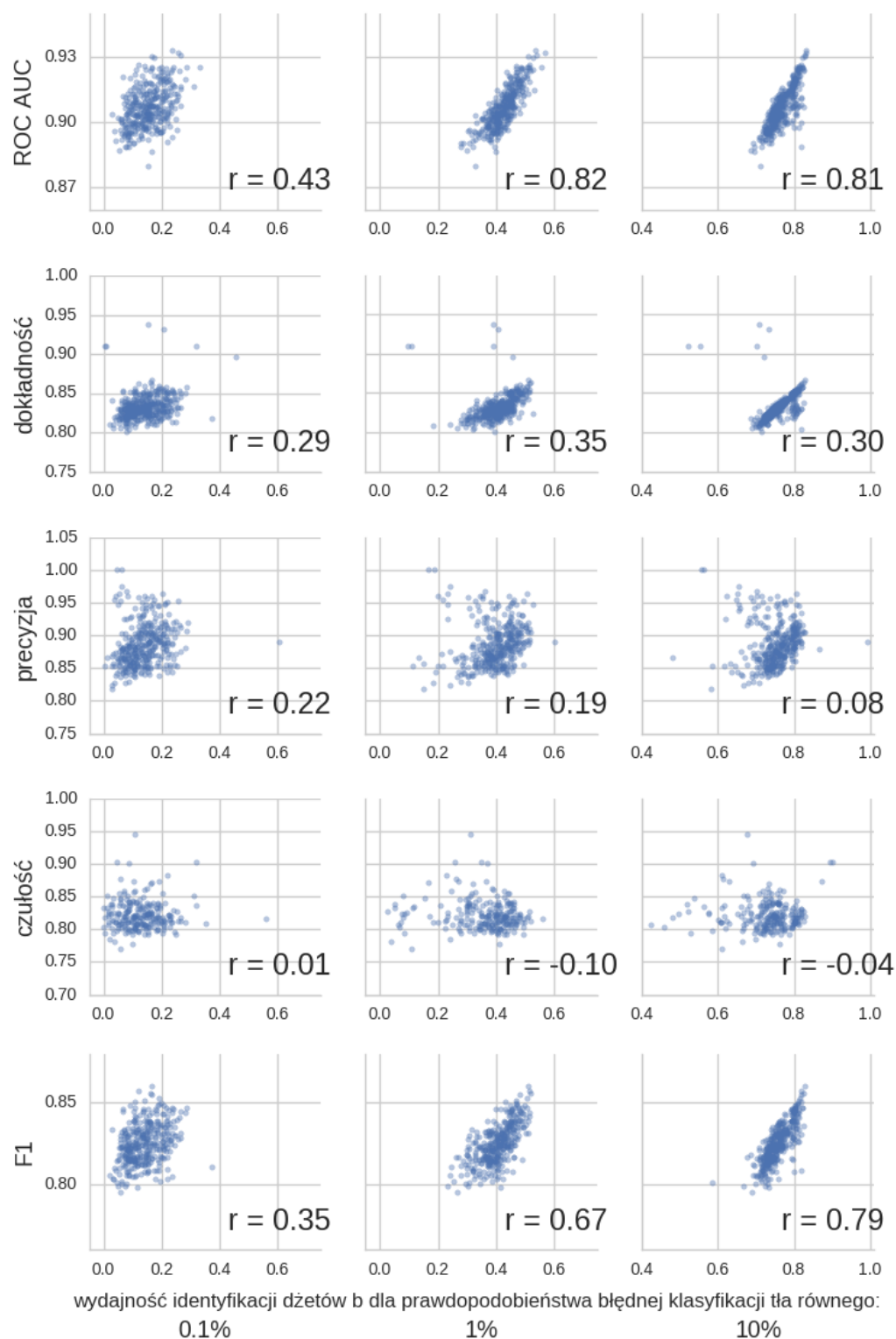
389 Naturalnym wyborem wydaje się pole pod powierzchnią krzywej ROC (ang. *ROC Area*
 390 *Under Curve – ROC AUC*). Potencjalną przeszkodą może być zakres rozsądnych wartości
 391 prawdopodobieństwa błędnej klasyfikacji przypadków tła: dżety b stanowią tylko kilka procent
 392 liczby wszystkich dżetów, zatem z punktu widzenia analizy dopuszczalne będą punkty pracy
 393 zapewniające wydajność identyfikacji dżetów b ok. 10 – 100 razy większą niż częstość niepo-
 394 prawnego zaklasyfikowania przypadków tła. Oznacza to, że zdecydowana większość punktów
 395 pracy znajdujących się na krzywej ROC jest nie do zaakceptowania – interesujące są tylko te o
 396 najniższych częstościach.

397 Aby ilościowo porównywać różne algorytmy wprowadzone zostaną trzy punkty pracy: o
 398 prawdopodobieństwie błędnej klasyfikacji tła równej 0.1%, 1% oraz 10%. Na Rys. 7 przedsta-
 399 wione zostały zależności poszczególnych metryk od wydajności identyfikacji dżetów b w tych
 400 trzech punktach pracy. Każdy punkt odpowiada jednemu eksperymentowi (dla dowolnego al-
 401 gorytmu) przeprowadzonemu w trakcie przygotowywania analizy. Daje to pogląd, używanie
 402 której metryki zapewni jednocześnie wysokie wartości wydajności na identyfikację dżetów b w
 403 wybranych punktach pracy. Najwyższe korelacje występują dla punktu pracy o najwyższym
 404 prawdopodobieństwie błędnej klasyfikacji tła – jak będzie to pokazane później wyniki dla tego
 405 punktu pracy są najbardziej stabilne. Spośród analizowanych metryk najwyższe wartości współ-
 406 czynnika Pearsona otrzymano dla pola pod powierzchnią krzywej ROC, dosyć wysokie również
 407 dla dokładności i precyzji. Widać, że wartości czułości są najsłabiej skorelowane z wydajnością
 408 identyfikacji dżetów b – jest zrozumiałe, że są one dużo niższe niż dla precyzji: wartości czułości
 409 maleją gdy błędnie klasyfikowane są dżety b stanowiące sygnał, podczas gdy wartości precyzji
 410 maleją, gdy błędnie klasyfikowane są przypadki tła. Z tych dwóch błędów, drugi jest bardziej
 411 kosztowny, gdyż błędna klasyfikacja nawet niewielkiej części tła znacznie pogarsza czystość
 412 otrzymywanej próbki.

413 Na Rys. 8 przedstawiono te same wykresy, ale tym razem wybrano tylko 25% najlepszych
 414 wartości dla każdej metryki – te punkty są bardziej znaczące, gdyż ostatecznie modele z ekspe-
 415 rymentów dających najlepsze wyniki będą używane. Dla tych wykresów otrzymano zdecydowa-
 416 ną dominację *ROC AUC* – wybór modeli dających najwyższe pole pod powierzchnią krzywej
 417 ROC zapewnia jednocześnie otrzymanie wysokich wartości wydajności identyfikacji dżetów b
 418 dla wybranych punktów pracy.



Rysunek 7: Zależność podstawowych metryk od wydajności identyfikacji dżetów b dla punktów pracy o prawdopodobieństwie błędnej klasyfikacji tła równej 0.1%, 1% oraz 10%. Dla każdego wykresu przedstawiono współczynnik korelacji r Pearsona.



Rysunek 8: Rysunek podobny do Rys. 7, ale przedstawione zostały tylko punkty odpowiadające eksperymentom o wartościach metryki będących w górnym kwartylu wartości danej metryki dla wszystkich eksperymentów.

References

- 420 [1] Anton Andronic. “An overview of the experimental study of quark-gluon matter in high-
421 energy nucleus-nucleus collisions”. In: *Int. J. Mod. Phys. A* 29 (2014), p. 1430047. DOI:
422 10.1142/S0217751X14300476. arXiv: 1407.5003 [nucl-ex].
- 423 [2] Vardan Khachatryan et al. “Evidence for collectivity in pp collisions at the LHC”. In:
424 *Phys. Lett. B* 765 (2017), pp. 193–220. DOI: 10.1016/j.physletb.2016.12.009. arXiv:
425 1606.06198 [nucl-ex].
- 426 [3] Jaroslav Adam et al. “Enhanced production of multi-strange hadrons in high-multiplicity
427 proton-proton collisions”. In: *Nature Phys.* 13 (2017), pp. 535–539. DOI: 10.1038/nphys4111.
428 arXiv: 1606.07424 [nucl-ex].
- 429 [4] Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. “The Anti-k(t) jet clustering al-
430 gorithm”. In: *JHEP* 04 (2008), p. 063. DOI: 10.1088/1126-6708/2008/04/063. arXiv:
431 0802.1189 [hep-ph].
- 432 [5] Carlos A. Salgado and Urs Achim Wiedemann. “Calculating quenching weights”. In: *Phys.*
433 *Rev. D* 68 (2003), p. 014008. DOI: 10.1103/PhysRevD.68.014008. arXiv: hep-ph/0302184
434 [hep-ph].
- 435 [6] Yuri L. Dokshitzer and D. E. Kharzeev. “Heavy quark colorimetry of QCD matter”. In:
436 *Phys. Lett. B* 519 (2001), pp. 199–206. DOI: 10.1016/S0370-2693(01)01130-3. arXiv:
437 hep-ph/0106202 [hep-ph].
- 438 [7] Georges Aad et al. “Performance of *b*-Jet Identification in the ATLAS Experiment”. In:
439 *JINST* 11.04 (2016), P04008. DOI: 10.1088/1748-0221/11/04/P04008. arXiv: 1512.
440 01094 [hep-ex].
- 441 [8] A. M. Sirunyan et al. “Identification of heavy-flavour jets with the CMS detector in pp
442 collisions at 13 TeV”. In: *JINST* 13.05 (2018), P05011. DOI: 10.1088/1748-0221/13/
443 05/P05011. arXiv: 1712.07158 [physics.ins-det].
- 444 [9] Tianqi Chen and Carlos Guestrin. “XGBoost: A Scalable Tree Boosting System”. In:
445 *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery*
446 *and Data Mining. KDD '16*. San Francisco, California, USA: ACM, 2016, pp. 785–794.
447 ISBN: 978-1-4503-4232-2. DOI: 10.1145/2939672.2939785. URL: [http://doi.acm.org/](http://doi.acm.org/10.1145/2939672.2939785)
448 [10.1145/2939672.2939785](http://doi.acm.org/10.1145/2939672.2939785).
- 449 [10] James Bergstra and Yoshua Bengio. “Random Search for Hyper-parameter Optimization”.
450 In: *J. Mach. Learn. Res.* 13 (Feb. 2012), pp. 281–305. ISSN: 1532-4435. URL: [http://dl.](http://dl.acm.org/citation.cfm?id=2188385)
451 [acm.org/citation.cfm?id=2188385](http://dl.acm.org/citation.cfm?id=2188385).2188395.
- 452 [11] Kurt Hornik. “Approximation capabilities of multilayer feedforward networks”. In: *Neural*
453 *Networks* 4.2 (1991), pp. 251–257. ISSN: 0893-6080. DOI: [https://doi.org/10.1016/](https://doi.org/10.1016/0893-6080(91)90009-T)
454 [0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T). URL: [http://www.sciencedirect.com/science/article/](http://www.sciencedirect.com/science/article/pii/089360809190009T)
455 [pii/089360809190009T](http://www.sciencedirect.com/science/article/pii/089360809190009T).
- 456 [12] François Chollet et al. *Keras*. <https://keras.io>. 2015.
- 457 [13] Martín Abadi et al. *TensorFlow: Large-Scale Machine Learning on Heterogeneous Sys-*
458 *tems*. Software available from [tensorflow.org](https://www.tensorflow.org). 2015. URL: <https://www.tensorflow.org/>.
- 459 [14] Timothy Dozat. *Incorporating Nesterov Momentum into Adam*. 2015. URL: [{http://](http://cs229.stanford.edu/proj2015/054_report.pdf)
460 cs229.stanford.edu/proj2015/054_report.pdf}.

- 461 [15] Diederik P. Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In:
462 *CoRR* abs/1412.6980 (2014). arXiv: 1412.6980. URL: [http://arxiv.org/abs/1412.](http://arxiv.org/abs/1412.6980)
463 6980.
- 464 [16] Nitish Srivastava et al. “Dropout: A Simple Way to Prevent Neural Networks from Over-
465 fitting”. In: *J. Mach. Learn. Res.* 15.1 (Jan. 2014), pp. 1929–1958. ISSN: 1532-4435. URL:
466 <http://dl.acm.org/citation.cfm?id=2627435.2670313>.