Politechnika Rzeszowska Wydział Elektrotechniki i Informatyki



Sztuczna Inteligencja Klasyfikator win w Burn _{Projekt}

Spis treści

1. Opis problemu	3
1.1. Istota problemu	3
1.2. Technologiczne aspekty rozwiązania	3
1.3. Cel projektu	4
2. Opis części praktycznej	4
2.1. Uczenie modelu	4
2.2. Nadzorowane uczenie	4
2.3. Ocena skuteczności modelu	5
2.3.1. Dokładność	5
2.3.2. Strata	5
3. Analiza danych wejściowych	6
3.1. Opis zbioru danych	6
3.2. Struktura zbioru danych	6
4. Przygotowanie danych do nauki	7
4.1. Czyszczenie danych	7
4.2. Normalizacja danych	8
4.3. Podział na zbiory treningowe i testowe	8
4.4. Zapis danych do plików	10
4.5. Podsumowanie	10
5. Trenowanie modelu	11
5.1. Przygotowanie danych	11
5.1.1. Implementacja tensorów	11
6. Skrypt	12

1. Opis problemu

Klasyfikacja win, polegająca na przypisaniu konkretnego trunku do właściwej odmiany winogron na podstawie jego parametrów chemicznych, odgrywa ważną rolę w analizie jakości oraz unikalnych cech wina. Zróżnicowanie cech win pod względu składu chemicznego sprawia, że klasyfikacja win jest trudnym zadaniem.

Dzięki nowoczesnym rozwiązaniom, takim jak sztuczna inteligencja i algorytmy uczenia głębokiego, możliwe jest zautomatyzowanie procesu klasyfikacji win. Problem ten można sprowadzić do zadania wieloklasowej klasyfikacji, w którym model, analizując cechy takie jak zawartość alkoholu, poziom kwasowości czy intensywność barwy, jest w stanie rozpoznać odmianę winogron wykorzystaną do produkcji danego wina.

1.1. Istota problemu

Klasyfikacja win na podstawie ich składu chemicznego ma znaczenie zarówno dla producentów, jak i konsumentów. Dla producentów stanowi narzędzie umożliwiające kontrolę jakości, identyfikację charakterystycznych cech produktu oraz optymalizację procesów wytwarzania. Z kolei dla konsumentów wiedza o właściwościach wina pozwala na lepsze dopasowanie do indywidualnych preferencji smakowych.

Automatyzacja tego procesu za pomocą modeli sztucznej inteligencji nie tylko przyspiesza i ułatwia analizę, ale również otwiera nowe możliwości w zakresie badań nad różnorodnością win oraz ich właściwościami.

1.2. Technologiczne aspekty rozwiązania

Zastosowanie sztucznej inteligencji, w szczególności algorytmów uczenia maszynowego, stanowi kluczowy element w procesie automatycznej klasyfikacji win. Nowoczesne techniki, takie jak sieci neuronowe czy głębokie uczenie (Deep Learning), umożliwiają modelom analizowanie złożonych wzorców w danych chemicznych, które są trudne do wychwycenia przez tradycyjne metody.

Do implementacji tych sieci najczęściej wykorzystuje się popularne biblioteki, takie jak PyTorch i TensorFlow, które oferują szeroki zestaw narzędzi do budowy, trenowania i optymalizacji sieci neuronowych. Dodatkowo, Burn, biblioteka stworzona w języku Rust, staje się coraz bardziej popularną alternatywą do bibliotek w językach wysokopoziomowych, dzięki swojej wydajności i możliwościach języka niskopoziomowego.

1.3. Cel projektu

Celem projektu jest opracowanie modelu sztucznej inteligencji, który będzie w stanie automatycznie klasyfikować wina na podstawie ich cech chemicznych. Projekt zakłada stworzenie systemu wykorzystującego algorytmy głębokiego uczenia, który na podstawie danych, precyzyjnie przypisze wino do jednej z określonych odmian. Dążymy do opracowania modelu o wysokiej dokładności, który może zostać wykorzystany zarówno w badaniach naukowych, jak i w przemyśle winiarskim do automatycznej klasyfikacji produktów.

2. Opis części praktycznej

2.1. Uczenie modelu

Uczenie modelu w kontekście Deep Learningu polega na wykorzystaniu zaawansowanych algorytmów sztucznej inteligencji, szczególnie sieci neuronowych, do rozpoznawania wzorców i zależności w danych. *Deep Learning*, będący poddziedziną uczenia maszynowego, różni się od tradycyjnych metod uczenia maszynowego tym, że automatycznie wykonuje proces ekstrakcji cech z surowych danych, bez potrzeby ręcznego selekcjonowania istotnych atrybutów. Proces ten odbywa się dzięki warstwom sieci neuronowych, które są odpowiedzialne za stopniowe przekształcanie danych wejściowych na reprezentacje coraz bardziej złożone.

Uczenie modelu odbywa się w procesie zwanym treningiem, który polega na minimalizowaniu funkcji błędu. Funkcja ta mierzy, jak bardzo przewidywania modelu różnią się od rzeczywistych etykiet w danych. Podczas treningu sieć neuronowa iteracyjnie dostosowuje swoje parametry (*wagi i biasy*), aby jak najlepiej odwzorować prawdziwe dane wyjściowe. Używa się tutaj algorytmu optymalizacji, najczęściej **spadku gradientu** (*Gradient Descent*), który pozwala na modyfikację wag w taki sposób, by minimalizować błąd predykcji.

W procesie tym kluczową rolę odgrywa funkcja aktywacji, która wprowadza nieliniowość do modelu i pozwala na modelowanie bardziej złożonych zależności. Przykładowe funkcje aktywacji to **ReLU** (*Rectified Linear Unit*), **sigmoida** czy **tanh**, które pomagają w decyzji, czy dane wejściowe mają być uwzględnione w obliczeniach w danej warstwie.

W trakcie treningu model uczy się generalizować, czyli wykrywać ogólne wzorce, które będą skuteczne nie tylko dla danych treningowych, ale także dla nowych, niewidzianych wcześniej przykładów. Aby uniknąć przeuczenia (*overfitting*), stosuje się techniki, takie jak regularizacja, *dropout* czy *early stopping*, które pomagają utrzymać model w równowadze pomiędzy dokładnością a zdolnością do uogólniania.

2.2. Nadzorowane uczenie

Uczenie nadzorowane to metoda uczenia maszynowego, w której model jest trenowany na danych wejściowych, które są już przypisane do odpowiednich klas (czyli wyników). W przypadku klasyfikacji win oznacza to, że dla każdego przykładu (wina) znamy już jego kategorię, czyli odmianę winogron,

z której pochodzi. Model wykorzystuje te informacje, aby nauczyć się, jak przypisać nowe, nieznane dane do właściwej klasy.

W procesie uczenia modelu, sieć neuronowa analizuje cechy chemiczne wina (takie jak zawartość alkoholu, kwasowość, intensywność barwy itp.), a następnie na podstawie tych danych podejmuje decyzję, do której z trzech odmian winogron należy dane wino. Model jest trenowany na wcześniej oznaczonych danych (winach o znanych odmianach), a celem jest minimalizacja błędu klasyfikacji, aby przewidywania modelu były jak najbardziej dokładne.

2.3. Ocena skuteczności modelu

Testowanie modelu polega na sprawdzeniu jego zdolności do poprawnego przypisywania nowych, niewidzianych wcześniej danych do odpowiednich klas. W tym przypadku model, który został wytrenowany na danych zawierających cechy chemiczne win, jest testowany na osobnym zbiorze testowym, który nie był używany w trakcie treningu. Celem jest ocena, jak dobrze model generalizuje swoje umiejętności na nowych danych, a nie tylko na tych, które były obecne w zestawie treningowym.

Podczas testowania modelu ocenia się jego skuteczność przy pomocy odpowiednich metryk, takich jak **dokładność** (*accuracy*), która wskazuje, jaka część wszystkich przewidywań była poprawna oraz **strata** (*loss*), która określa, jak dobrze model radzi sobie z minimalizacją błędów w klasyfikacji.

2.3.1. Dokładność

Dokładność modelu to procent poprawnych przewidywań w stosunku do wszystkich przewidywań. Im wyższa wartość dokładności, tym lepsza jest skuteczność modelu. Na przykład, jeśli model przewidział poprawnie 90 z 100 próbek, to jego dokładność wynosi 90%. Dokładność jest jedną z najczęściej stosowanych metryk w problemach klasyfikacji, ponieważ jest intuicyjna i łatwa do interpretacji.

$$\label{eq:accuracy} Accuracy = \frac{Liczba \ poprawnych \ przewidywań}{Liczba \ wszystkich \ próbek}$$

2.3.2. Strata

Strata (*loss*) to wartość funkcji błędu, która mierzy, jak bardzo przewidywania modelu różnią się od rzeczywistych etykiet w danych.

Dla problemu klasyfikacji, jedną z najczęściej stosowanych funkcji *loss* jest *cross-entropy loss*, która mierzy różnicę między przewidywaną a rzeczywistą rozkładem prawdopodobieństw dla poszczególnych klas. Im mniejsza wartość *loss*, tym lepsza jest skuteczność modelu, ponieważ oznacza to, że przewidywania modelu są bardziej zbliżone do rzeczywistych etykiet.

3. Analiza danych wejściowych

Analiza danych wejściowych to kluczowy etap w każdym projekcie uczenia maszynowego, który pozwala na zrozumienie struktury danych, ich właściwości oraz przygotowanie ich do dalszego przetwarzania i trenowania modelu. W tym przypadku, celem analizy jest dokładne zrozumienie cech chemicznych win, które stanowią podstawę do klasyfikacji różnych odmian winogron.

3.1. Opis zbioru danych

Zbiór danych, który jest wykorzystywany w tym projekcie, pochodzi z analizy chemicznej win pochodzących z tego samego regionu we Włoszech, ale z trzech różnych odmian winogron. Zawiera on 13 cech chemicznych, które zostały zmierzone dla każdego z win, takich jak zawartość alkoholu, kwasowość, intensywność barwy oraz inne właściwości, które wpływają na charakterystykę wina. W zbiorze danych znajdują się również etykiety, które informują o rodzaju odmiany winogron, z której pochodzi dane wino.

Dokładnym źródłem danych jest "PARVUS - An Extendible Package for Data Exploration, Classification and Correlation". Instytut Analiz Farmaceutycznych i Żywnościowych oraz Technologii, Genua, Włochy. To źródło dostarcza narzędzie o nazwie PARVUS, które jest pakietem rozszerzalnym do eksploracji danych, klasyfikacji i analizy korelacji. Zostało opracowane przez zespół badaczy w Instytucie, który specjalizuje się w analizach farmaceutycznych i żywnościowych. Pakiet ten umożliwia dokładną analizę chemiczną win i może być używany do różnych celów badawczych, w tym klasyfikacji różnych rodzajów win.

3.2. Struktura zbioru danych

Każdy wiersz w zbiorze danych reprezentuje jedno wino, a 13 kolumn odpowiada różnym cechom chemicznym, które zostały zmierzone w próbce. Wartości w tych kolumnach są liczbowe, co oznacza, że dane są już w formie numerycznej, gotowe do dalszego przetwarzania. Również warto wspomnieć, że nie ma brakujących danych w zbiorze, co ułatwia analizę.

Zbiór danych podzielony jest na 178 próbek, a liczby próbek dla poszczególnych odmian winogron to:

```
• Klasa 1: 59 próbek - (33.1%)
```

- Klasa 2: 71 próbek (39.9%)
- Klasa 3: 48 próbek (27.0%)

Atrybuty chemiczne, które zostały zmierzone dla każdego wina, to:

- 1. **Alkohol** Zawartość alkoholu w winie, wyrażona jako procent objętości alkoholu w stosunku do całkowitej objętości płynu. Jest to jeden z kluczowych wskaźników wpływających na smak i charakterystykę wina.
- 2. **Kwas jabłkowy** Zawartość kwasu jabłkowego w winie. Jest to naturalny kwas organiczny, który wpływa na smak wina, nadając mu kwaskowatość. Wina o wyższej zawartości tego kwasu mogą być bardziej kwasowe i orzeźwiające.

- 3. **Popiół** Zawartość popiołu w winie, który jest pozostałością po spaleniu organicznych składników w winie. Zawartość popiołu może mieć wpływ na mineralność wina.
- 4. **Zasadowość popiołu** Określa poziom zasadowości popiołu w winie. Zasadowość wpływa na pH wina i może oddziaływać na smak, zwłaszcza w kontekście równowagi kwasowości i słodyczy.
- 5. **Magnez** Zawartość magnezu w winie. Magnez jest jednym z minerałów, który wpływa na smak wina, a jego obecność może wpływać na ogólną jakość i strukturę wina.
- 6. **Całkowita zawartość fenoli** Fenole są naturalnymi związkami organicznymi występującymi w winie, które odpowiadają za smak, zapach oraz właściwości zdrowotne wina. Ich zawartość może wpływać na goryczkę oraz astringencję wina.
- 7. **Flawanoidy** Jest to grupa fenoli, która wpływa na smak, zapach i kolor wina. Flawanoidy są również antyoksydantami, co ma znaczenie dla trwałości wina.
- 8. **Fenole nieflawanoidowe** Inna grupa fenoli, która nie należy do flawanoidów, ale również wpływa na smak i zapach wina, nadając mu specyficzne cechy organoleptyczne.
- 9. **Proantocyjaniny** Związki odpowiedzialne za kolor wina, zwłaszcza czerwonych win. Są to silne antyoksydanty, które również wpływają na strukturę smaku i astringencję.
- 10. Intensywność koloru Określa głębokość koloru wina, co jest ważnym atrybutem w przypadku win czerwonych i białych. Intensywność koloru może wskazywać na dojrzałość wina oraz jego skład chemiczny.
- 11. **Odcień** Dotyczy barwy wina, który może się różnić w zależności od odmiany winogron, procesu fermentacji i starzenia wina. Odcień jest kluczowym elementem oceny jakości wina.
- 12. **OD280/OD315 rozcieńczonych win** Stosunek absorbancji przy długości fali 280 nm do 315 nm, który jest wykorzystywany do oceny jakości i zawartości substancji organicznych w winie, takich jak fenole.
- 13. **Prolina** Aminokwas występujący w winie, który wpływa na strukturę wina oraz jego stabilność. Zawartość proliny może mieć wpływ na właściwości smakowe i chemiczne wina.

4. Przygotowanie danych do nauki

Odpowiednie przygotowanie danych ma na celu zapewnienie, że model będzie w stanie uczyć się skutecznie i osiągać jak najlepsze wyniki. W tym etapie dokonuje się kilku istotnych operacji, takich jak czyszczenie danych, normalizacja, dodanie brakujących wartości, podział na zbiory treningowe i testowe, a także inne techniki, które pomagają w dostosowaniu danych do wymagań modelu.

Poniżej przedstawię kroki, które ja podjąłem w celu przygotowania danych do trenowania modelu klasyfikacji win.

4.1. Czyszczenie danych

Pierwszym krokiem w przygotowaniu danych jest ich czyszczenie, czyli usunięcie zbędnych informacji, brakujących wartości czy duplikatów. W przypadku zbioru danych, który został wykorzystany w tym projekcie, nie było potrzeby usuwania żadnych próbek. Każdy atrybut może potencjalnie wpłynąć na klasyfikację wina, dlatego nie ma potrzeby eliminacji.

4.2. Normalizacja danych

Normalizacja danych jest ważnym krokiem w przypadku modeli, które wykorzystują algorytmy uczenia maszynowego, ponieważ pozwala na ujednolicenie skali wartości atrybutów. W przypadku zbioru danych z cechami chemicznymi win, wartości poszczególnych atrybutów różnią się znacząco, co może wpłynąć na proces uczenia modelu.

Np. zawartość alkoholu w winie mieści się w zakresie od 11.0% do 14.8%, podczas gdy zawartość kwasu jabłkowego waha się od 1.74 do 4.23. Aby uniknąć problemów związanych z różnicą w skali wartości, zastosowałem normalizację.

W procesie normalizacji dane są przeskalowane w taki sposób, aby mieściły się w określonym zakresie. W moim przypadku od -1 do 1.

$$x_{
m norm} = rac{2*(x-x_{
m min})}{x_{
m max}-x_{
m min}} - 1$$

Wzór 1: Normalizacji danych do zakresu [-1, 1]

```
® Rust
1 for (i, row) in x.iter().enumerate() {
      for (j, &value) in row.iter().enumerate() {
3
       if \max[j] - \min[j] = 0.0  {
          x_norm[i][j] = 0.0; // Avoid division by zero
4
       } else {
5
          // Normalize to range [-1, 1]
6
7
          x_{norm[i][j]} = -1.0 + 2.0 * (value - min[j]) / (max[j] - min[j]);
8
       }
9
      }
10 }
```

Program 1: Normalizacja danych

4.3. Podział na zbiory treningowe i testowe

Podział danych na zbiory treningowe i testowe jest kluczowym elementem w procesie uczenia maszynowego, ponieważ pozwala na ocenę skuteczności modelu na nowych, nie widzianych wcześniej danych. W moim przypadku zdecydowałem się na podział danych w stosunku 80% do 20%, gdzie 80% danych zostało wykorzystane do treningu, a 20% do testowania. W późniejszych etapach zmieniałem ten podział na 70% do 30% w celach eksperymentów.

Podział danych na zbiory treningowe i testowe pozwala na ocenę skuteczności modelu na nowych, nie widzianych wcześniej danych. W ten sposób można sprawdzić, jak dobrze model generalizuje swoje umiejętności na nowych danych, a nie tylko na tych, które były obecne w zestawie treningowym.

W tym celu napisałem funkcję split_data, która dokonuje losowego podziału danych na zbiory treningowe i testowe. Po podziale dane są dodatkowo sortowane według klasy, aby zapewnić równomierne rozłożenie próbek w zbiorach treningowym i testowym.

```
® Rust
1 fn split_data(
     x: Vec<Vec<f32>>,
3 y_t: Vec<i32>
4 ) \rightarrow (Vec<Vec<f32>>, Vec<i32>, Vec<Vec<f32>>, Vec<i32>) {
5 let mut rng = thread_rng();
      let mut combined: Vec<(Vec<f32>, i32)> = x.into_iter().zip(y_t.into_iter()).collect();
7
    combined.shuffle(&mut rng);
9
     let num_records = combined.len();
10
     let num_train = (num_records as f32 * 0.8) as usize;
11
12
     // Split data into training and testing sets
    let (train_data, test_data): (Vec<_>, Vec<_>) = combined
13
14
      .into_iter()
15
      .enumerate()
       .partition(|&(i, _)| i ≥ num_train);
16
17
     // Sort by y_t class and unzip to separate x and y_t
18
19
     let mut train_sorted: Vec<_> = train_data.into_iter().map(|(_, data)| data).collect();
20
      train\_sorted.sort\_by\_key(|&(\_, y)| y);
     let (x_train, y_t_train): (Vec<_>, Vec<_>) = train_sorted.into_iter().unzip();
21
22
23
     let mut test_sorted: Vec<_> = test_data.into_iter().map(|(_, data)| data).collect();
      test\_sorted.sort\_by\_key(|&(\_, y)| y);
     let (x_test, y_t_test): (Vec<_>, Vec<_>) = test_sorted.into_iter().unzip();
25
27
      (x_train, y_t_train, x_test, y_t_test)
28 }
```

Program 2: Podział danych na zbiory treningowe i testowe

4.4. Zapis danych do plików

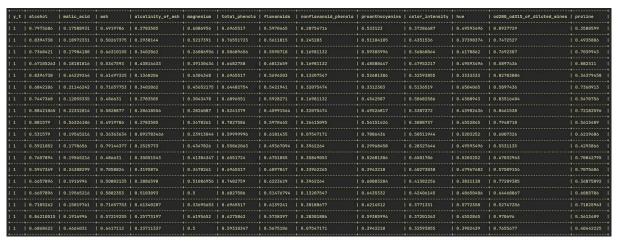
Tak przygotowane dane zostały zapisane do plików w formacie PKL. Do zserializowania danych wykorzystałem bibliotekę serde-pickle z pakietu serde, która pozwala na serializację i deserializację danych w formacie PKL. Dane zostały zapisane w postaci czterech plików: x_train.pkl, y_t_train.pkl, x_test.pkl i y_t_test.pkl.

```
1 /// Save `x` and `y_t` to disk in pickle format.
                                                                                        ® Rust
2 fn dump_to_pkl(x: Vec<Vec<f32>>, y_t: Vec<i32>, prefix: &str) {
3 // Create BTreeMaps to serialize
     let mut x_map = BTreeMap::new();
    let mut y_map = BTreeMap::new();
     x_map.insert(String::from("x"), x);
7
     y_map.insert(String::from("y_t"), y_t);
     // Serialize to pickle format
     let x_serialized = serde_pickle::to_vec(&x_map, Default::default()).unwrap();
     let y_serialized = serde_pickle::to_vec(&y_map, Default::default()).unwrap();
11
12
13
     // Save to disk
      std::fs::write(format!("./data/x_{f}.pkl", prefix), &x_serialized).unwrap();
14
     std::fs::write(format!("./data/y_t_{}.pkl", prefix), &y_serialized).unwrap();
15
16 }
```

Program 3: Funkcja zapisująca dane do plików PKL

4.5. Podsumowanie

Przygotowanie danych do trenowania modelu klasyfikacji win wymagało kilku istotnych kroków, takich jak czyszczenie, normalizacja, podział na zbiory treningowe i testowe oraz zapis do plików. Dzięki tym operacjom dane są gotowe do dalszego przetwarzania i trenowania modelu.



Rysunek 1: Tabela z częścią danych przygotowanych do trenowania modelu

5. Trenowanie modelu

Po przygotowaniu danych przyszedł czas na trenowanie modelu. W tym projekcie zdecydowałem się na wykorzystanie biblioteki Burn do implementacji modelu, jak i jego trenowania.

Przed samym przystąpieniem do trenowania modelu, podtrzebne jest zdefiniowanie architektury sieci neuronowej, przygotowanie danych treningowych i testowych oraz wybór odpowiednich hiperparametrów, takich jak współczynnik uczenia, liczba epok czy rozmiar wsadu (*batch size*).

5.1. Przygotowanie danych

Trenowanie modelu w bibliotece Burn wymaga tzw. *Batcherów* (plik batcher.rs w części skryptowej), które są odpowiedzialne za dostarczanie danych do modelu w postaci wsadów (*batches*). *Batcher* przyjmuje strukture *Dataset* (plik dataset.rs, w części skryptowej) jako dane wejściowe i dzieli je na wsady o określonym rozmiarze.

W moim przypadku, struktura WineDataset przechowuje dane treningowe i testowe, które zostały wcześniej przygotowane i zapisane do plików PKL. Następnie, WineBatcher dzieli te dane na wsady WineBatch, które są przekazywane do modelu podczas treningu.

5.1.1. Implementacja tensorów

Ważnym elementem w trenowaniu modelu jest implementacja tensorów, które są podstawowymi strukturami danych w bibliotece Burn. Tensor reprezentuje wielowymiarowy wektor lub macierz, który jest przechowywany w pamięci urządzenia. W przpadku mojego programu, tensory tworzone są w funkcji batch w strukturze WineBatcher. Przechowują one dane wejściowe i wyjściowe dla modelu, które są przekazywane do sieci neuronowej.

Pierwsze tensor jest tworzony dla każdego wina z danych treningowych, a drugi tensor przechowuje etykietę klasy dla tego wina. Następnie, te tensory są łączone w jeden tensor dla danych wejściowych i wyjściowych, który jest przekazywany bezpośrednio do modelu.

```
impl<B: Backend> Batcher<WineItem, WineBatch<B>> for WineBatcher<B> {
                                                                                             ® Rust
2
      fn\ batch(\&self,\ data:\ Vec<WineItem>) \rightarrow WineBatch<B> {
3
        let x = data
          .iter()
5
          .map(|wine| {
            Tensor::<B, 2>::from_data(
6
7
8
                wine.alcohol,
9
                wine.malic_acid,
10
                wine.ash,
11
                wine.alcalinity_of_ash,
                wine.magnesium,
13
                wine.total_phenols,
14
                wine.flavanoids,
                wine.nonflavanoid_phenols,
15
                wine.proanthocyanins,
16
17
                wine.color_intensity,
                wine.hue,
18
19
                wine.od280_od315_of_diluted_wines,
20
                wine.proline,
21
              ]],
              &self.device,
22
23
24
          7)
          .collect::<Vec<_>>();
26
27
        let y = data
          .iter()
28
          .map(|wine| Tensor::<B, 1, Int>::from_data([wine.class], &self.device))
30
          .collect::<Vec<_>>();
31
        let x = Tensor :: <B, 2> :: cat(x, 0).to_device(&self.device);
32
33
        let y = Tensor::<B, 1, Int>::cat(y, 0).to_device(&self.device);
35
        WineBatch { x, y }
36
     }
37 }
```

 $Program~4: Implementacji~funkcji~{\tt batch},~kt\'orej~wymaga~struktura~{\tt WineBatcher}$

6. Skrypt

```
12
      let (x, y_t) = extract_x_and_y_t(&wines);
13
14
      // 2. Normalize `x` → `x_norm
15
      let x_norm = normalize_x(&x);
16
17
      // Sorting by `y_t` is unnecessary for this dataset. Data is already sorted by class.
18
19
     // 3. Split data into training and testing sets
      let (x_train, y_t_train, x_test, y_t_test) = split_data(x_norm, y_t);
21
      // 3. Display data in a table
22
23
     println!("Data test:");
24
      table_data(&x_train, &y_t_train);
25
     println!("\n\n\nData train:");
      table_data(&x_test, &y_t_test);
27
28
      // 4. Save everything to an PKL file
29
     dump_to_pkl(x_train, y_t_train, "train");
30
      dump_to_pkl(x_test, y_t_test, "test");
31
32
      println!("Data processing and saving completed.");
33 }
34
35 /// Load data file into a `String`.
36 fn load_data_file() \rightarrow String {
37     read_to_string("./data/wine.data").unwrap()
38 }
39
40 /// Parse raw data into a Vector of `Wine` structs.
41 fn parse_data(data: String) → Vec<WineItem> {
42
      data.lines()
43
      .map(|line| {
          let mut wine_data = line.split(',').map(|s| s.parse::<f32>().unwrap());
44
45
         WineItem {
46
            class: wine data.next().unwrap() as i32.
47
           alcohol: wine_data.next().unwrap(),
            malic acid: wine data.next().unwrap(),
48
49
            ash: wine_data.next().unwrap(),
50
            alcalinity_of_ash: wine_data.next().unwrap(),
51
            magnesium: wine_data.next().unwrap(),
52
            total_phenols: wine_data.next().unwrap(),
53
            flavanoids: wine_data.next().unwrap(),
54
            nonflavanoid_phenols: wine_data.next().unwrap(),
55
            proanthocyanins: wine_data.next().unwrap(),
56
            color_intensity: wine_data.next().unwrap(),
            hue: wine_data.next().unwrap(),
57
            od280_od315_of_diluted_wines: wine_data.next().unwrap(),
            proline: wine_data.next().unwrap(),
59
60
61
        })
62
         .collect()
63 }
65 /// Extract `x` (Wine attributes 1-13) and `y_t` (class attribute).
    fn \ extract_x\_and\_y\_t(wines: \&[WineItem]) \rightarrow (Vec<Vec<f32>>, Vec<i32>) {
    let num_records = wines.len();
      let mut x = vec![vec![0.0; 13]; num_records];
69
     let mut y_t = vec![0; num_records];
70
71
     for (i, wine) in wines.iter().enumerate() {
72
        x[i][0] = wine.alcohol;
      x[i][1] = wine.malic_acid;
73
```

```
x[i][2] = wine.ash;
      x[i][3] = wine.alcalinity_of_ash;
76
        x[i][4] = wine.magnesium;
77
        x[i][5] = wine.total_phenols;
        x[i][6] = wine.flavanoids;
       x[i][7] = wine.nonflavanoid_phenols;
        x[i][8] = wine.proanthocyanins;
       x[i][9] = wine.color_intensity;
81
        x[i][10] = wine.hue;
83
      x[i][11] = wine.od280_od315_of_diluted_wines;
        x[i][12] = wine.proline;
85
86
        y_t[i] = wine.class;
87
88
89 (x, y_t)
90 }
91
92 /// Normalize `x` (array of wine attributes 1-13).
93 fn normalize_x(x: &Vec<Vec<f32>>) \rightarrow Vec<<math>Vec<f32>> {
94
      let num_columns = x[0].len();
95
96
      let mut min = vec![f32::MAX; num_columns];
97
     let mut max = vec![f32::MIN; num_columns];
98
99
      for row in x.iter() {
100
        for (j, &value) in row.iter().enumerate() {
101
         if value < min[j] {</pre>
102
            min[j] = value;
103
104
          if value > max[j] {
          max[j] = value;
105
106
          }
107
108
      }
109
110
      let mut x_norm = x.clone();
111
      for (i, row) in x.iter().enumerate() {
        for (j, &value) in row.iter().enumerate() {
112
113
        if \max[j] - \min[j] = 0.0 {
114
            x_norm[i][j] = 0.0; // Avoid division by zero
115
          } else {
116
            // Normalize to range [-1, 1]
117
            x_norm[i][j] = -1.0 + 2.0 * (value - min[j]) / (max[j] - min[j]);
118
119
120
121
122
      x_norm
123 }
124
125 /// Display data in a table.
127  let mut table = Vec::new();
128
129
     // Add data to table
      for (i, row) in x.iter().enumerate() {
      let mut row_data = Vec::new();
131
132
        row_data.push(y_t[i].to_string());
133
        for &value in row.iter() {
134
            row_data.push(value.to_string());
135
```

```
136
        table.push(row_data);
137 }
138
139 // Add header to table
      let table = table.table().title(vec![
140
141
     "y_t".cell().bold(true),
142
         "alcohol".cell().bold(true),
       "malic_acid".cell().bold(true),
143
         "ash".cell().bold(true),
145
      "alcalinity_of_ash".cell().bold(true),
         "magnesium".cell().bold(true),
        "total_phenols".cell().bold(true),
         "flavanoids".cell().bold(true),
149
      "nonflavanoid_phenols".cell().bold(true),
         "proanthocyanins".cell().bold(true),
     "color_intensity".cell().bold(true),
152
         "hue".cell().bold(true),
153
       "od280_od315_of_diluted_wines".cell().bold(true),
154
         "proline".cell().bold(true),
155
156
157
      print_stdout(table).unwrap();
158 }
159
160 /// Save `x` and `y_t` to disk in pickle format.
161 fn dump_to_pkl(x: Vec<Vec<f32>>, y_t: Vec<i32>, prefix: &str) {
     // Create BTreeMaps to serialize
162
163    let mut x_map = BTreeMap::new();
164    let mut y_map = BTreeMap::new();
165  x_map.insert(String::from("x"), x);
166
     y_map.insert(String::from("y_t"), y_t);
167
168
      // Serialize to pickle format
169
      let x_serialized = serde_pickle::to_vec(&x_map, Default::default()).unwrap();
170
      let y_serialized = serde_pickle::to_vec(&y_map, Default::default()).unwrap();
171
172
      // Save to disk
173
     std::fs::write(format!("./data/x_{\frac{1}{2}}.pkl", prefix), &x_serialized).unwrap();
174
      std::fs::write(format!("./data/y_t_{}.pkl", prefix), &y_serialized).unwrap();
175 }
176
177 /// Split data into training and testing sets.
178 use rand::seq::SliceRandom;
179 use rand::thread_rng;
180
181 fn split_data(
182
        x: Vec<Vec<f32>>,
183
        y_t: Vec<i32>,
184 ) \rightarrow (Vec<Vec<f32>>, Vec<i32>, Vec<Vec<f32>>, Vec<i32>) {
185    let mut rng = thread_rng();
      let mut combined: Vec<(Vec<f32>, i32)> = x.into_iter().zip(y_t.into_iter()).collect();
187
      combined.shuffle(&mut rng);
189
      let num_records = combined.len();
190
      let num_train = (num_records as f32 * 0.7) as usize;
191
192
      // Split data into training and testing sets
193
      let (train_data, test_data): (Vec<_>, Vec<_>) = combined
194
        .into_iter()
195
        .enumerate()
196
        .partition(|&(i, _)| i ≥ num_train);
197
```

```
// Sort by y_t class

// Sort by y_t class

// Let mut train_sorted: Vec<_> = train_data.into_iter().map(|(_, data)| data).collect();

// train_sorted.sort_by_key(|&(_, y)| y);

// Let (x_train, y_t_train): (Vec<_>, Vec<_>) = train_sorted.into_iter().unzip();

// Let mut test_sorted: Vec<_> = test_data.into_iter().map(|(_, data)| data).collect();

// test_sorted.sort_by_key(|&(_, y)| y);

// Let (x_test, y_t_test): (Vec<_>, Vec<_>) = test_sorted.into_iter().unzip();

// (x_train, y_t_train, x_test, y_t_test)

// (x_train, y_t_train, x_test, y_t_test)

// (X_train, y_t_train, x_test, y_t_test)
```

Program 5: prepare_data.rs

```
1    use burn::data::dataset::{Dataset, InMemDataset};
                                                                                                       ® Rust
2  use serde::{Deserialize, Serialize};
3     use std::collections::BTreeMap;
4  use std::fs::File;
5  use std::io;
7 #[derive(Debug, PartialEq, Clone, Serialize, Deserialize)]
8   pub struct WineItem {
9 #[serde(rename = "class")]
10 pub class: i32,
11
12
      #[serde(rename = "Alcohol")]
13 pub alcohol: f32,
14
15
     #[serde(rename = "Malic acid")]
16
      pub malic_acid: f32,
17
18
      #[serde(rename = "Ash")]
19
     pub ash: f32,
20
21
     #[serde(rename = "Alcalinity of ash")]
22
      pub alcalinity_of_ash: f32,
23
24
      #[serde(rename = "Magnesium")]
25
     pub magnesium: f32,
26
27
      #[serde(rename = "Total phenols")]
28
      pub total_phenols: f32,
29
      #[serde(rename = "Flavanoids")]
30
31
      pub flavanoids: f32,
32
33
      #[serde(rename = "Nonflavanoid phenols")]
34
      pub nonflavanoid_phenols: f32,
35
36
      #[serde(rename = "Proanthocyanins")]
37
      pub proanthocyanins: f32,
38
39
      #[serde(rename = "Color intensity")]
      pub color_intensity: f32,
40
41
      #[serde(rename = "Hue")]
42
43
      pub hue: f32,
      #[serde(rename = "OD280/OD315 of diluted wines")]
      pub od280_od315_of_diluted_wines: f32,
47
48
      #[serde(rename = "Proline")]
49
      pub proline: f32,
```

```
50 }
52  pub struct WineDataset {
53    pub(crate) dataset: InMemDataset<WineItem>,
55
56 impl WineDataset {
57 pub fn subset(&self, indices: &[usize]) → Self {
         let data: Vec<WineItem> = indices
59
          .map(|&i| self.dataset.get(i).unwrap().clone())
          .collect();
         Self {
63
        dataset: InMemDataset::new(data),
64
65
66
   }
67
68 impl WineDataset {
69
   pub fn from_pkl(split: &str) → Result<Self, io::Error> {
70
         let x_file = File::open(format!("./data/x_{}.pkl", split))?;
71
        let y_file = File::open(format!("./data/y_t_{}.pkl", split))?;
72
73
         let x_map: BTreeMap<String, Vec<Vec<f32>>> =
74
          serde\_pickle::from\_reader(x\_file, \ \textit{Default}::default())
75
            .map_err(|e| io::Error::new(io::ErrorKind::InvalidData, e))?;
76
         let y_map: BTreeMap<String, Vec<i32>> =
77
         serde_pickle::from_reader(y_file, Default::default())
78
             .map_err(|e| io::Error::new(io::ErrorKind::InvalidData, e))?;
79
80
         let mut data = Vec::new();
81
         for (key, x_values) in x_map.get("x").unwrap().iter().enumerate() {
82
83
          if let Some(y_values) = y_map.get("y_t").unwrap().get(key) {
84
             let wine = WineItem {
85
              class: *y_values,
86
               alcohol: x_values[0],
87
              malic_acid: x_values[1],
88
               ash: x_values[2],
89
              alcalinity_of_ash: x_values[3],
90
               magnesium: x_values[4],
91
               total_phenols: x_values[5],
               flavanoids: x_values[6],
93
              nonflavanoid_phenols: x_values[7],
               proanthocyanins: x_values[8],
95
               color_intensity: x_values[9],
96
               hue: x_values[10],
97
               od280_od315_of_diluted_wines: x_values[11],
               proline: x_values[12],
99
100
             data.push(wine);
101
102
103
104
         let dataset = InMemDataset::new(data);
105
106
         Ok(Self { dataset })
107
108
      pub fn train() \rightarrow Self {
109
110
         Self::from_pkl("train").unwrap()
111
```

```
112
113 pub fn test() \rightarrow Self {
114
        Self::from_pkl("test").unwrap()
115 }
116 }
117
118 impl Dataset<WineItem> for WineDataset {
119 fn \ get(\&self, index: usize) \rightarrow Option<WineItem> {
        self.dataset.get(index)
121
122
123 fn len(\&self) \rightarrow usize {
      self.dataset.len()
125 }
126 }
```

Program 6: dataset.rs

```
1  use crate::dataset::WineItem;
2  use burn::data::dataloader::batcher::Batcher;
3 use burn::prelude::{Backend, Int, Tensor};
5 #[derive(Clone)]
6  pub struct WineBatcher<B: Backend> {
7 device: B::Device,
8 }
9
10 impl<B: Backend> WineBatcher<B> {
11 pub fn new(device: B::Device) \rightarrow Self {
      Self { device }
13 }
14 }
15
16 #[derive(Clone, Debug)]
17 pub struct WineBatch<B: Backend> {
19 pub y: Tensor<B, 1, Int>,
20 }
21
22 impl<B: Backend> Batcher<WineItem, WineBatch<B>> for WineBatcher<B> {
24
     let x = data
25
     .iter()
26
        .map(|wine| {
27
        Tensor::<B, 2>::from_data(
28
          ГΓ
29
           wine.alcohol,
30
              wine.malic_acid,
31
             wine.ash,
32
              wine.alcalinity_of_ash,
33
             wine.magnesium,
34
              wine.total_phenols,
35
              wine.flavanoids,
              wine.nonflavanoid_phenols,
36
37
              wine.proanthocyanins,
              wine.color_intensity,
39
              wine.hue,
              wine.od280_od315_of_diluted_wines,
41
              wine.proline,
            ]],
43
            &self.device,
44
45
        })
```

```
46
        .collect::<Vec<_>>();
47
48
       let y = data
49
      .iter()
50
         .map(|wine| Tensor::<B, 1, Int>::from_data([wine.class], &self.device))
51
       .collect::<Vec<_>>();
52
     let x = Tensor::<B, 2>::cat(x, 0).to_device(&self.device);
53
       let y = Tensor::<B, 1, Int>::cat(y, 0).to_device(&self.device);
55
       WineBatch { x, y }
57 }
58 }
```

Program 7: batcher.rs