

Raport – Laboratorium PTI: ĆW2

Reprezentacja, modelowanie i liczenie informacji

1. Cel realizowanego ćwiczenia

W niniejszej pracy nastąpiło porównanie „losowości”: ciągu liczb pseudolosowych ze zbioru $\{0,1,2,3\}$, wygenerowanych za pomocą programu komputerowego oraz ciągu liczb stworzonego w możliwie jak najbardziej przypadkowy, abstrakcyjny sposób – wybór padł na przekształcanie dowolnych danych pogodowych.

W celu porównania obu ciągów liczbowych, obliczone zostały prawdopodobieństwa poszczególnych przejść (dla modeli bez pamięci oraz z pamięcią rzędu 2. i 3.), a następnie wyznaczona została entropia. Ponadto, dla modeli z pamięcią sporządzono tabele przejść w Excelu, aby lepiej zobrazować zjawisko.

2. Dobór danych testowych

Prace odbywały się na dwóch zbiorach danych: pierwszy „data.txt” utworzony za pomocą własnoręcznie napisanej funkcji, wykorzystującej `randi()`, w Matlabie. Drugi „dane_chmura.txt” zawierał dane ze strony Leuven (podanej w źródłach), przetworzone w celu uzyskania ciągu o alfabecie $\{0,1,2,3\}$. Dane zostały pobrane z kolumny E (Column5) ukazującej stężenie PM_{2,5} zarejestrowane przez odpowiednie sondy. W danych tych kropki zostały zamienione na przecinki w celu dalszej obróbki. Następnie uzyskane dane w formacie liczbowym (dla Excela) zostały zaokrąglone w dół. Następnie utworzony został ciąg składający się z liczb opisanych powyżej, na których zastosowano operację modulo 4. Do dalszej analizy wziętych zostało pierwsze 100000 wyrazów tego ciągu.

3. Kategoria 1: model bez pamięci

Wyniki modelu Markowa bez pamięci zostały uzyskane dzięki funkcji napisanej w Matlabie (dla której zostały podane następujące argumenty: 100000 elementów, alfabet 4, rząd 0).

Uzyskane prawdopodobieństwa na poszczególne wyrazy alfabetu w modelu bez pamięci dla danych „z pogody” (odpowiednio na cyfry 0,1,2,3):

0.24259 0.26337 0.25292 0.24111

Oraz entropia dla danych "z pogody":

0.9995389

Odpowiednie prawdopodobieństwa na poszczególne wyrazy alfabetu w modelu bez pamięci dla danych z Matlaba (odpowiednio na cyfry 0,1,2,3):

0.25061 0.24903 0.25036 0.24999

Oraz entropia dla danych z Matlaba:

0.9999979

Uzyskane wyniki wskazały, że ciąg liczbowy danych "z pogody" jest zauważalnie mniej losowy, co widać w szczególności przy obliczonych prawdopodobieństwach dla cyfry '2'. Jednakże, jeśli spojrzymy na poziomy entropii to dla obu ciągów są one bardzo bliskie 1 co wskazuje na to, że dane są losowe.

4. Kategoria 2: model Markowa rzędu 2

W celu obliczenia prawdopodobieństw oraz entropii źródła dla modelu Markowa rzędu 2, rozróżnionych zostało 18 stanów: 16 ze zbioru $\{(0,0), (0,1), \dots, (3,2), (3,3)\}$, odpowiadających wszystkim możliwościom wyboru 2 elementów z alfabetu 4-elementowego; a także 2 dodatkowe: S – stan pierwotny, "startowy", E – stan oznaczający koniec ciągu. Następnie zliczone zostały przejścia - przykładowo dla fragmentu ciągu liczb: ...012..., inkrementowana jest liczba przejść ze stanu 0,1 (stan początkowy w przejściu definiuje wiersz w poniższych tabelach) do stanu 1,2 (stan końcowy definiuje kolumnę w poniższych tabelach).

Na podstawie tak dokonanych obliczeń, wyznaczone zostały prawdopodobieństwa poszczególnych przejść. Bazując na tych wynikach, policzono entropię źródła. W plikach zamieszczono: arkusz .xlsx, z którego skopiowano poniższe tabele oraz kod programu „model_rzad2.cpp”, gdzie zawarte są szczegółowe obliczenia dla tego modelu. Poniżej przedstawiono wyniki tych obliczeń.

Jak zaobserwować można w poniższych tabelach, zdecydowanie bardziej równomierny rozkład liczby przejść uzyskano dla danych wygenerowanych w Matlabie, dla których większość wartości oscylowała w okolicy średniej - czyli ok. 1550. Dla danych wygenerowanych „z pogody”, zaobserwowane odchylenia od średniej są już większe, np. różnice takie jak między liczbami przejść wynoszącymi 1379 oraz 2994. Jednakże należy również zaznaczyć, że mimo tych odchylen, dane wygenerowane „z pogody” wciąż wykazały zróżnicowanie:

Tabela prawdopodobieństw przejść dla danych wygenerowanych „z pogody”:

Stan pocz \ Stan końc	S	0,0	0,1	0,2	0,3	1,0	1,1	1,2	1,3	2,0	2,1	2,2	2,3	3,0	3,1	3,2	3,3	E	
S	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	
0,0	0	0,274845	0,254874	0,223173	0,247107	0	0	0	0	0	0	0	0	0	0	0	0	0	
0,1	0	0	0	0	0	0,257184	0,276234	0,235117	0,231465	0	0	0	0	0	0	0	0	0	
0,2	0	0	0	0	0	0	0	0	0	0,256608	0,257658	0,24243	0,243305	0	0	0	0	0	
0,3	0	0	0	0	0	0	0	0	0	0	0	0	0	0,265914	0,240317	0,236948	0,25682	0	
1,0	0	0,258147	0,275879	0,235304	0,230671	0	0	0	0	0	0	0	0	0	0	0	0	0	
1,1	0	0	0	0	0	0,220873	0,374672	0,220373	0,184082	0	0	0	0	0	0	0	0	0	
1,2	0	0	0	0	0	0	0	0	0	0,229511	0,288815	0,259399	0,222275	0	0	0	0	0	
1,3	0	0	0	0	0	0	0	0	0	0	0	0	0	0,253622	0,254844	0,245069	0,246465	0	
2,0	0	0,246649	0,254482	0,253264	0,245605	0	0	0	0	0	0	0	0	0	0	0	0	0	
2,1	0	0	0	0	0	0,231722	0,274783	0,268897	0,224597	0	0	0	0	0	0	0	0	0	
2,2	0	0	0	0	0	0	0	0	0	0,203152	0,233816	0,327797	0,235236	0	0	0	0	0	
2,3	0	0	0	0	0	0	0	0	0	0	0	0	0	0,225897	0,225236	0,277824	0,271043	0	
3,0	0	0,259378	0,252649	0,231623	0,25635	0	0	0	0	0	0	0	0	0	0	0	0	0	
3,1	0	0	0	0	0	0,246737	0,265153	0,246558	0,241373	0	0	0	0	0	0	0	0	0,000178795	
3,2	0	0	0	0	0	0	0	0	0	0,224705	0,242998	0,275376	0,256921	0	0	0	0	0	
3,3	0	0	0	0	0	0	0	0	0	0	0	0	0	0,241832	0,209942	0,263561	0,284665	0	
E	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1

Ostatecznie, po obliczeniu entropii źródła, dla modelu Markowa rzędu 2, otrzymano następujące wyniki:

Entropia źródła (oraz dla poszczególnych stanów) dla danych wygenerowanych w Matlabie:

```
Entropie dla poszczególnych stanów:
0.999974 0.999907 0.999978 0.999969 0.999956 0.999945 0.999915 1.00092 0.999992 0.99974 0.999894 0.99983 0.999906 0.99982 0.999972 0.999539
Entropia źródła wynosi: 0.999953
```

Entropia źródła (oraz dla poszczególnych stanów) dla danych wygenerowanych „z pogody”:

```
Entropie dla poszczególnych stanów:
0.998014 0.998132 0.999705 0.999188 0.998099 0.971084 0.99603 0.999894 0.999912 0.997181 0.987963 0.996519 0.999307 1.0006 0.998005 0.995534
Entropia źródła wynosi: 0.995342
```

Widocznym jest, że entropia źródła dla danych wygenerowanych w Matlabie, wynosząca 0,999953, jest większa od entropii źródła obliczonej dla danych wygenerowanych „z pogody”, wynoszącej 0,995342. Zatem dane otrzymane w wyniku działania funkcji randi() są bardziej losowe niż dane obliczone z „pogody”. Jednak, co zaskakujące, stwierdzono, że wartość entropii źródła dla tych drugich danych wciąż jest bardzo dobra, gdyż jest stosunkowo bliska wartości 1, zatem dane te również są dobrym przybliżeniem danych losowych – oczywiście jednak słabszym niż pierwszy ciąg liczbowy.

5. Kategoria 3: model Markowa rzędu 3

Obliczenia przebiegały analogicznie do modelu Markowa rzędu 2. Najpierw zostały określone 64 możliwe stany ze zbioru $\{(0\ 0\ 0), (0\ 0\ 1) \dots (3\ 3\ 2), (3\ 3\ 3)\}$ i 2 dodatkowe – pierwotny (S) i zamykający ciąg (E). Po zliczeniu wszystkich możliwych przejść z jednego stanu w drugi, obliczone zostały prawdopodobieństwa warunkowe. Część wyników przedstawiają poniższe tabelki:

Dla danych wygenerowanych w Matlabie:

	S	000	001	002	003	...	330	331	332	333	E
S	0	0	0	0	0	...	1	0	0	0	0
000	0	0,2275	0,2492	0,2562	0,2670	...	0	0	0	0	0
001	0	0	0	0	0	...	0	0	0	0	0
002	0	0	0	0	0	...	0	0	0	0	0
003	0	0	0	0	0	...	0	0	0	0	0
...
330	0	0	0	0	0	...	0	0	0	0	0
331	0	0	0	0	0	...	0	0	0	0	0
332	0	0	0	0	0	...	0	0	0	0	0
333	0	0	0	0	0	...	0,2574	0,2593	0,2341	0,2492	0
E	0	0	0	0	0	...	0	0	0	0	1

Dla danych wygenerowanych z pogody:

	S	000	001	002	003	...	330	331	332	333	E
S	0	0	0	0	0	...	0	0	0	0	0
000	0	0,3154	0,2416	0,1920	0,2508	...	0	0	0	0	0
001	0	0	0	0	0	...	0	0	0	0	0
002	0	0	0	0	0	...	0	0	0	0	0
003	0	0	0	0	0	...	0	0	0	0	0
...
330	0	0	0	0	0	...	0	0	0	0	0
331	0	0	0	0	0	...	0	0	0	0	0,0007
332	0	0	0	0	0	...	0	0	0	0	0
333	0	0	0	0	0	...	0,2185	0,1916	0,2734	0,3163	0
E	0	0	0	0	0	...	0	0	0	0	1

Chociaż różnice między wynikami są niewielkie, wartości takie jak 0,3154 i 0,3163, widoczne w tabeli powyżej, oznaczają, że niektóre przejścia występują znacznie częściej. Stąd można więc wywnioskować, że dane z drugiego źródła są mniej losowe niż te wygenerowane w Matlabie.

Wyznaczone prawdopodobieństwa warunkowe zostały wykorzystane do obliczenia entropii dla każdego stanu:

	Dane wygenerowane w Matlabie	Dane wygenerowane z pogody		Dane wygenerowane w Matlabie	Dane wygenerowane z pogody		Dane wygenerowane w Matlabie	Dane wygenerowane z pogody
H(S 0 0 0)	0.998783	0.988923	H(S 1 1 0)	0.999508	0.985809	H(S 2 2 0)	0.999561	0.99951
H(S 0 0 1)	0.999061	0.997862	H(S 1 1 1)	0.999121	0.87711	H(S 2 2 1)	0.999584	0.993515
H(S 0 0 2)	0.999568	0.999295	H(S 1 1 2)	0.999428	0.986212	H(S 2 2 2)	0.999632	0.958568
H(S 0 0 3)	0.999225	0.998626	H(S 1 1 3)	0.999939	0.999911	H(S 2 2 3)	0.999487	0.987322
H(S 0 1 0)	0.999479	0.998729	H(S 1 2 0)	0.999837	0.999601	H(S 2 3 0)	0.999541	0.997089
H(S 0 1 1)	0.999941	0.991553	H(S 1 2 1)	0.997884	0.991259	H(S 2 3 1)	0.999829	0.996835
H(S 0 1 2)	0.998789	0.999971	H(S 1 2 2)	0.999719	0.995785	H(S 2 3 2)	0.999408	0.992079
H(S 0 1 3)	1.00279	0.999626	H(S 1 2 3)	0.999911	0.99919	H(S 2 3 3)	0.998943	0.990302
H(S 0 2 0)	0.999851	0.999442	H(S 1 3 0)	0.999433	0.998466	H(S 3 0 0)	0.999693	0.997886
H(S 0 2 1)	0.999457	0.998285	H(S 1 3 1)	0.999621	0.997108	H(S 3 0 1)	0.999666	0.998353
H(S 0 2 2)	0.999482	0.998712	H(S 1 3 2)	0.999831	0.997877	H(S 3 0 2)	0.999863	0.999952
H(S 0 2 3)	0.997893	0.99968	H(S 1 3 3)	0.998507	0.999509	H(S 3 0 3)	0.999782	0.996246
H(S 0 3 0)	0.999746	0.995894	H(S 2 0 0)	0.999269	0.999832	H(S 3 1 0)	0.999732	0.999668
H(S 0 3 1)	0.999585	0.998497	H(S 2 0 1)	0.999858	0.996872	H(S 3 1 1)	0.998531	0.999361
H(S 0 3 2)	0.99994	0.999903	H(S 2 0 2)	0.999061	0.999666	H(S 3 1 2)	0.999795	0.99668
H(S 0 3 3)	0.999211	0.99447	H(S 2 0 3)	0.99964	0.999926	H(S 3 1 3)	0.999398	0.999087
H(S 1 0 0)	0.999762	0.997503	H(S 2 1 0)	0.999323	0.999878	H(S 3 2 0)	0.999552	0.999856
H(S 1 0 1)	0.999222	0.993413	H(S 2 1 1)	0.999238	0.990431	H(S 3 2 1)	0.999152	0.999294
H(S 1 0 2)	0.999417	0.998724	H(S 2 1 2)	0.999841	0.994046	H(S 3 2 2)	0.998835	0.985235
H(S 1 0 3)	0.999928	0.998978	H(S 2 1 3)	0.999666	0.999426	H(S 3 2 3)	0.999648	0.987854
						H(S 3 3 0)	0.99941	0.99896
						H(S 3 3 1)	0.998969	1.00183
						H(S 3 3 2)	0.999324	0.991028
						H(S 3 3 3)	0.999429	0.986546

Uśredniając wyniki po wszystkich stanach otrzymujemy wartość entropii źródła:

- Dla danych wygenerowanych w Matlabie: 0,999477
- Dla danych wygenerowanych z pogody: 0,991606

Chociaż obie wartości można uznać za dobre, mniejsza entropia dla danych wygenerowanych z pogody oznacza, że są one mniej losowe. Zgadza się to z wyciągniętymi wcześniej wnioskami oraz wynikami otrzymanymi dla modelu bez pamięci i modelu Markowa rzędu 2.

6. Zestawienie wyników i dyskusja całości badań. Formułowanie wniosków końcowych.

Uzyskano różniące się wyniki, jednak nie są to różnice drastyczne. Dla badanych rzędów modeli, nie uzyskano dużych różnic w poziomach entropii pomiędzy oboma ciągami znaków.

Początkowo, wraz ze wzrostem rzędu pamięci, różnica między entropiami rosła, aczkolwiek nie uzyskała ona znacznej wartości. Najmniejsza różnica jest widoczna dla modelu bez pamięci - wtedy, chociaż entropia obliczona dla danych wygenerowanych w Matlabie jest większa, to można przyjąć, że są one w przybliżeniu równe.

Przyjmując do obliczeń modele większych rzędów, tzn. w badanym przypadku rzędów większych niż 8, można zaobserwować drastyczny spadek entropii ze względu na zbyt dużą liczbę możliwych kontekstów w stosunku do ilości danych, jednak wciąż różnica między entropiami nie jest znacząca.

rzęd modelu	0	1	2	3	4	5	6	7	8
dane z matlaba	0,999998	0,999980	0,999893	0,999427	0,997301	0,989020	0,953663	0,785655	0,401451
dane z "chmury"	0,999539	0,997577	0,995281	0,991556	0,986116	0,973813	0,935615	0,765955	0,400820

Pomimo widocznych różnic poziomy entropii wskazują na to, że ilość informacji w obu zbiorach jest podobna.

7. Opis wykorzystanych narzędzi, źródeł itp.:

Narzędzia:

- Matlab – program generujący liczby losowe ze zbioru $\{0, 1, 2, 3\}$, program liczący entropię dla modelu bez pamięci i modelu Markowa dowolnego rzędu
- C++ - programy liczące entropie warunkową dla modelu Markowa rzędu 2 i 3
- Excel – wizualizacja obliczonych wyników

Źródła:

- <https://leuvenair.be/data/data-l.html> - dane pogodowe wykorzystane do wygenerowania drugiego ciągu liczbowego (Complete datadump 2018 (from 2018-01-30))
- Wikipedia
- Prezentacje z wykładów

Stworzone pliki i programy:

- *data_z_matlaba.txt*
- *data_z_pogody.txt*

8. Opis innowacyjnych osiągnięć własnych realizacji ćwiczenia:

Realizacja ćwiczenia odbyła się w oparciu o ideę sprawdzenia losowości danych “z pogody”, które zostały wyszukane, pobrane oraz przetworzone.

Podczas realizacji projektu, zostały napisane programy w C++ obliczające entropię, prawdopodobieństwa oraz liczby przejść. Uzyskane dane następnie zostały przetworzone.

Ponadto, napisany został również program w Matlabie, który wyznacza entropię dla dowolnego rzędu, alfabetu i danych.

