

MA4: Sprawozdanie z projektu

Denys Moldovan 335967

January 18, 2025

Contents

1	Wprowadzenie teoretyczne	2
1.1	Badane hipotezy	2
1.2	Prawdopodobieństwo a priori i a posteriori	2
1.3	Funkcja wiarygodności	2
2	Wyniki dla przypadków I i II	3
2.1	Przypadek I: Równe prawdopodobieństwa a priori	3
2.1.1	Opis wyników	3
2.1.2	Wykres zmian prawdopodobieństw posteriori	3
2.1.3	Wnioski z analizy	3
2.2	Przypadek II: Preferencja dla jednego języka	4
2.2.1	Opis wyników	4
2.2.2	Wykresy wyników	5
2.2.3	Wnioski	6
3	Metody aktualizacji Bayesa	7
3.1	Standardowa metoda aktualizacji	7
3.1.1	Opis działania	7
3.1.2	Kod implementacji	7
3.1.3	Zalety i ograniczenia	7
3.2	Metoda z kryterium stopu	8
3.2.1	Opis działania	8
3.2.2	Kod implementacji	8
3.2.3	Zalety i ograniczenia	9
3.3	Porównanie obu metod	9
4	Modyfikacja wiadomości	11
4.1	Przypadek IV: Zmodyfikowana wiadomość	11
4.1.1	Założenia i cel analizy	11
4.1.2	Wyniki analizy	11
4.1.3	Wnioski	13

1 Wprowadzenie teoretyczne

Klasyfikacja Bayesowska opiera się na założeniach teorii prawdopodobieństwa, gdzie celem jest określenie prawdopodobieństwa posteriori dla każdej z hipotez na podstawie danych wejściowych. W kontekście tego projektu, hipotezy reprezentują języki, w których może być napisana wiadomość.

1.1 Badane hipotezy

Dla problemu przyjęto trzy hipotezy:

- H_W : Wiadomość została napisana w języku wakandyjskim (W).
- H_L : Wiadomość została napisana w języku łatweryjskim (L).
- H_S : Wiadomość została napisana w języku symkariańskim (S).

1.2 Prawdopodobieństwo a priori i a posteriori

- **Prawdopodobieństwo a priori** ($P(H)$): Określa nasze pierwotne założenia na temat hipotez przed analizą danych. W przypadku I przyjęto równomierne rozkłady dla wszystkich hipotez, tj. $P(H_W) = P(H_L) = P(H_S) = \frac{1}{3}$.
- **Prawdopodobieństwo a posteriori** ($P(H|D)$): Obliczane jest na podstawie danych D i określa prawdopodobieństwo hipotezy po uwzględnieniu tych danych, zgodnie z wzorem Bayesa:

$$P(H|D) = \frac{P(D|H) \cdot P(H)}{P(D)},$$

gdzie $P(D|H)$ to funkcja wiarygodności, a $P(D)$ to normalizator.

1.3 Funkcja wiarygodności

Funkcja wiarygodności $P(D|H)$ określa, jak dobrze dane D pasują do danej hipotezy. W tym projekcie wyznacza ona częstość występowania symboli w poszczególnych językach, na podstawie danych wejściowych.

2 Wyniki dla przypadków I i II

2.1 Przypadek I: Równe prawdopodobieństwa a priori

Dla tego przypadku, założono równomierne rozkłady początkowe:

$$P(H_W) = P(H_L) = P(H_S) = \frac{1}{3}.$$

2.1.1 Opis wyników

Na podstawie analizy uzyskano następujące obserwacje:

- Wykres zmian prawdopodobieństwa posteriori pokazuje, że kolejne symbole wiadomości prowadzi do sukcesywnej aktualizacji prawdopodobieństw dla każdego języka.
- Ostateczne wartości prawdopodobieństw posteriori wskazują na wyraźną dominację jednego języka po analizie pełnej wiadomości.
- Szczegółowe wartości prawdopodobieństw posteriori dla każdej iteracji zostały zapisane w pliku CSV, który znajduje się w katalogu `wyniki/csv/history_uniform.csv`.

2.1.2 Wykres zmian prawdopodobieństw posteriori

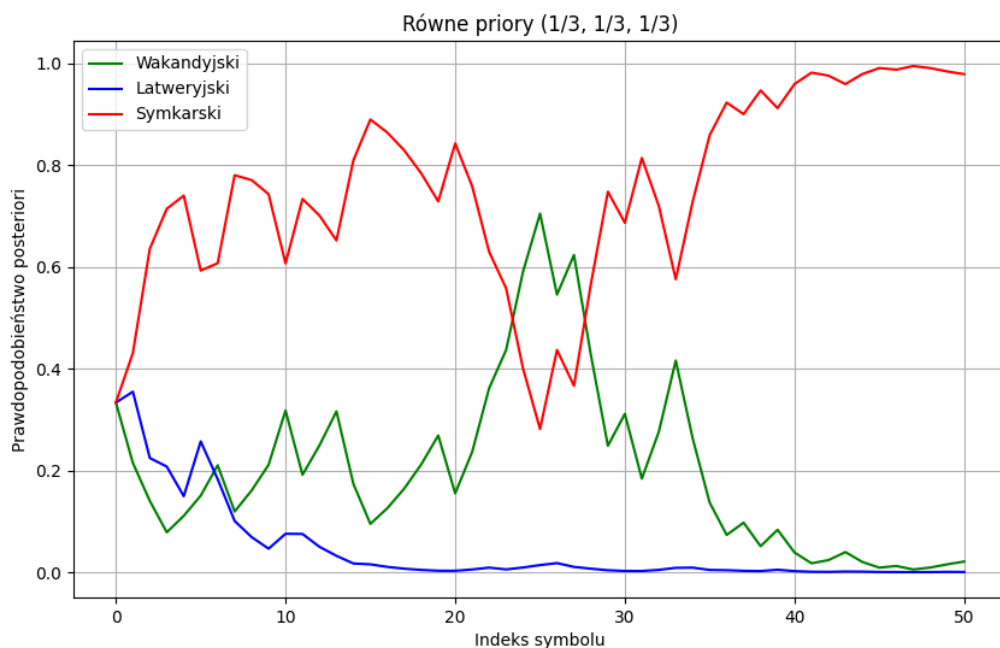


Figure 1: Zmiany prawdopodobieństw posteriori w Przypadku I: Równe prawdopodobieństwa a priori.

2.1.3 Wnioski z analizy

Na podstawie wyników uzyskanych dla równego rozkładu a priori można wyciągnąć następujące wnioski:

- Początkowe prawdopodobieństwa są równe, co oznacza maksymalną niepewność w rozpoznawaniu języka.
- W miarę przetwarzania kolejnych symboli wiadomości, prawdopodobieństwo posteriori dla jednego z języków rośnie, wskazując na jednoznaczny wybór języka.
- Model poprawnie aktualizuje rozkład posteriori w sposób zgodny z założeniami teorii Bayesa, co oznacza, że:
 - Prawdopodobieństwa posteriori są obliczane na podstawie zależności pomiędzy danymi wejściowymi (symbole wiadomości) a wcześniejszymi założeniami (rozkład a priori).
 - Przy każdym kroku, nowe informacje (symbole wiadomości) są uwzględniane, a wcześniejsze założenia są aktualizowane, prowadząc do coraz bardziej precyzyjnych oszacowań.
 - Proces aktualizacji uwzględnia zarówno dane empiryczne (obserwowane symbole), jak i początkowe założenia, zapewniając równowagę pomiędzy wpływem wcześniejszych danych a nowymi obserwacjami.
 - Ostateczne wyniki wskazują na dominację jednego języka, co świadczy o konwergencji modelu i jego zdolności do adaptacji do danych wejściowych.

2.2 Przypadek II: Preferencja dla jednego języka

W tym przypadku przyjęto różne wartości rozkładów początkowych, aby nadać preferencje jednemu z języków. Dla przykładu, dla języka łatweryjskiego założono:

$$P(H_L) = 0.6, \quad P(H_W) = P(H_S) = 0.2.$$

2.2.1 Opis wyników

- Na podstawie analizy prawdopodobieństwa posteriori widać, że przyjęcie początkowej preferencji dla języka łatweryjskiego (H_L) sprawia, że proces konwergencji w kierunku tego języka jest szybszy.
- Początkowe prawdopodobieństwa mają wpływ na dynamikę aktualizacji — preferowany język szybciej osiąga dominację w rozkładzie posteriori.
- Jednak niezależnie od początkowego rozkładu a priori, przy dostatecznie długiej wiadomości model zawsze konwerguje do języka najlepiej zgodnego z danymi (czyli z wiadomością).
- Wyniki wskazują, że ostateczny wynik analizy jest determinowany przez dane wejściowe, a nie przez priorytety początkowe.
- Szczegółowe wartości prawdopodobieństw posteriori zostały zapisane w plikach CSV znajdujących się w katalogu `wyniki/csv/`.

2.2.2 Wykresy wyników

Poniższy wykres przedstawia zmiany prawdopodobieństw posteriori w przypadku preferencji dla języka latweryjskiego:

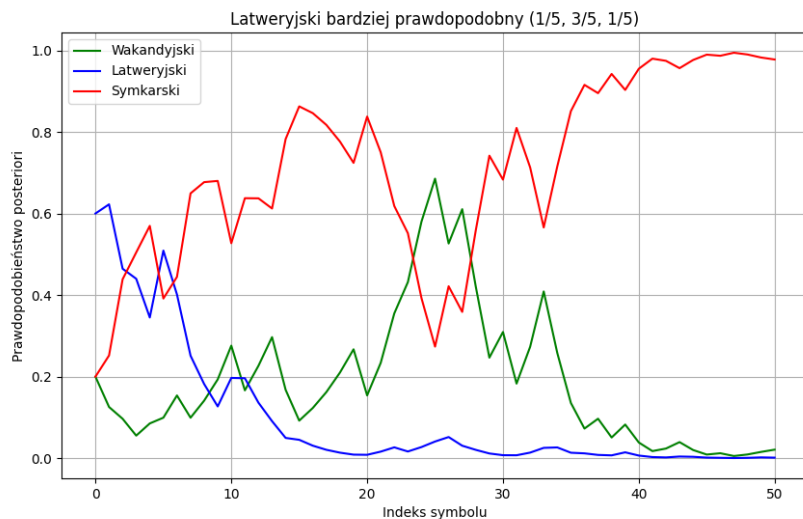


Figure 2: Prawdopodobieństwo posteriori dla preferencji języka latweryjskiego.

Analogiczne wyniki uzyskano dla innych przypadków preferencji:

- Preferencja dla języka symkarskiego ($P(H_S) = 0.6, P(H_W) = P(H_L) = 0.2$).
- Preferencja dla języka wakandyjskiego ($P(H_W) = 0.6, P(H_L) = P(H_S) = 0.2$).

Przykładowe wykresy:

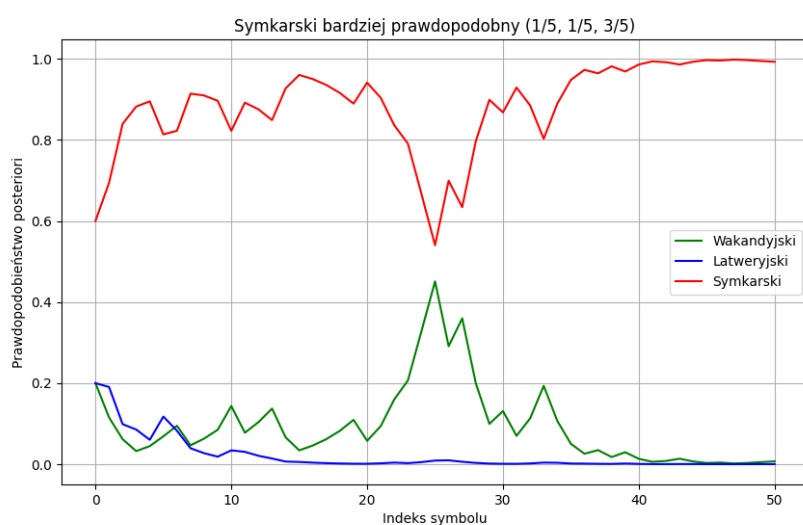


Figure 3: Prawdopodobieństwo posteriori dla preferencji języka symkarskiego.

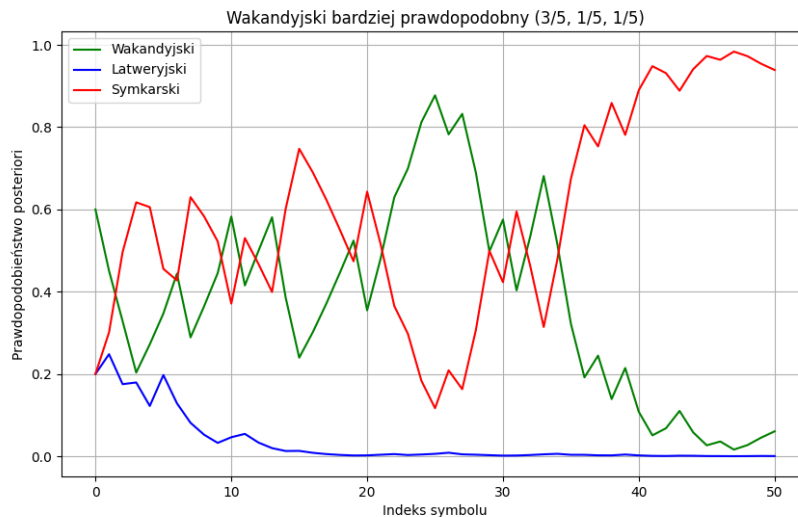


Figure 4: Prawdopodobieństwo posteriori dla preferencji języka wakandyjskiego.

2.2.3 Wnioski

Na podstawie uzyskanych wyników można sformułować następujące wnioski:

- Początkowe priorytety wpływają na szybkość konwergencji modelu w kierunku dominacji jednego języka.
- Dla krótkich wiadomości ($n < 10$ symboli) początkowe rozkłady a priori mogą znacząco wpłynąć na wynik, ponieważ model nie ma jeszcze wystarczającej ilości danych, aby "przewycieżyć" założenia początkowe.
- Dla długich wiadomości ($n \gg 10$ symboli) priorytety początkowe tracą na znaczeniu, a ostateczny wynik jest zdeterminowany przez dane wejściowe — rozkład prawdopodobieństwa symboli dla języków w wiadomości.
- Model poprawnie aktualizuje rozkład posteriori w sposób zgodny z założeniami teorii Bayesa, co potwierdzają obserwowane zmiany prawdopodobieństw na wykresach.

3 Metody aktualizacji Bayesa

3.1 Standardowa metoda aktualizacji

3.1.1 Opis działania

Standardowa metoda aktualizacji Bayesa iteracyjnie aktualizuje prawdopodobieństwa posteriori na podstawie kolejnych symboli wiadomości. Jest to realizowane zgodnie z regułą Bayesa, gdzie priorytetowe prawdopodobieństwa są aktualizowane w miarę przetwarzania nowych danych.

Działanie metody można opisać następująco:

1. Inicjalizacja wartości prawdopodobieństw priorytetowych $P(H_W)$, $P(H_L)$, $P(H_S)$.
2. Iteracyjne przetwarzanie każdego symbolu wiadomości, w którym następuje:
 - Aktualizacja prawdopodobieństw dla każdego języka na podstawie zgodności symbolu z danymi języka.
 - Normalizacja, aby zachować sumę prawdopodobieństw równą 1.
3. Zapisywanie wartości posteriori dla każdej iteracji.

3.1.2 Kod implementacji

Główne linie kodu standardowej metody aktualizacji są przedstawione poniżej:

```
def bayesian_update(message, pW, pL, pS):
    history = {"W": [], "L": [], "S": []}
    p_MW, p_ML, p_MS = 1, 1, 1

    history["W"].append(pW)
    history["L"].append(pL)
    history["S"].append(pS)

    for symbol in message:
        if symbol != "N":
            p_MW *= dwak_probs.get(symbol, 0.01)
            p_ML *= dlatver_probs.get(symbol, 0.01)
            p_MS *= dsymk_probs.get(symbol, 0.01)

            pM = (p_MW * pW) + (p_ML * pL) + (p_MS * pS)
            history["W"].append((p_MW * pW) / pM)
            history["L"].append((p_ML * pL) / pM)
            history["S"].append((p_MS * pS) / pM)

    return history
```

3.1.3 Zalety i ograniczenia

Zalety:

- Pełna analiza całej wiadomości gwarantuje dokładność wyników.

- Implementacja jest prosta i w pełni zgodna z teorią Bayesa.

Ograniczenia:

- Wymaga przetwarzania całej wiadomości, co może być czasochłonne przy dużych danych.
- Aktualizacje wykonywane są nawet wtedy, gdy model osiągnął już stabilność.

3.2 Metoda z kryterium stopu

3.2.1 Opis działania

Metoda z kryterium stopu wprowadza dodatkowy mechanizm zatrzymania procesu aktualizacji w momencie, gdy zmiany w prawdopodobieństwach posteriori stają się wystarczająco małe. Działanie metody opiera się na wprowadzeniu dwóch warunków:

- Kryterium konwergencji — proces zatrzymuje się, jeśli zmiany wartości posteriori pomiędzy kolejnymi iteracjami są mniejsze niż zadany próg ϵ .
- Maksymalna liczba iteracji — proces zatrzymuje się, jeśli liczba przetworzonych symboli osiągnie zadana wartość maksymalna.

3.2.2 Kod implementacji

Implementacja metody stopu:

```
def bayesian_stop(message, pW, pL, pS, epsilon=0.01, max_iterations=100):
    history = {"W": [], "L": [], "S": []}
    p_MW, p_ML, p_MS = 1, 1, 1

    history["W"].append(pW)
    history["L"].append(pL)
    history["S"].append(pS)

    for i, symbol in enumerate(message):
        if i >= max_iterations:
            print("Osiągnięto maksymalną liczbę iteracji.")
            break

        if symbol != "N":
            p_MW *= dwak_probs.get(symbol, 0.01)
            p_ML *= dlatver_probs.get(symbol, 0.01)
            p_MS *= dsymk_probs.get(symbol, 0.01)

    pM = (p_MW * pW) + (p_ML * pL) + (p_MS * pS)
    new_pW = (p_MW * pW) / pM
    new_pL = (p_ML * pL) / pM
    new_pS = (p_MS * pS) / pM

    history["W"].append(new_pW)
```



```

history["L"].append(new_pL)
history["S"].append(new_pS)

# Sprawdzenie kryterium stopu
if i > 0:
    delta_W = abs(history["W"][-1] - history["W"][-2])
    delta_L = abs(history["L"][-1] - history["L"][-2])
    delta_S = abs(history["S"][-1] - history["S"][-2])

    if delta_W < epsilon and delta_L < epsilon and delta_S < epsilon:
        print(f"Kryterium konwergencji osiagniete po {i+1} iteracjach.")
        break

return history

```

3.2.3 Zalety i ograniczenia

Zalety:

- Znacznie szybsze działanie dzięki wcześniejszemu zakończeniu procesu.
- Efektywne wykorzystanie zasobów obliczeniowych.

Ograniczenia:

- Wymaga ustawienia odpowiednich wartości ϵ i maksymalnej liczby iteracji.
- Może zakończyć działanie zbyt wcześnie, jeśli ϵ jest zbyt duże.

3.3 Porównanie obu metod

Porównanie obu metod przedstawiono na poniższym wykresie:

Porównanie pełnej iteracji i metody stopu

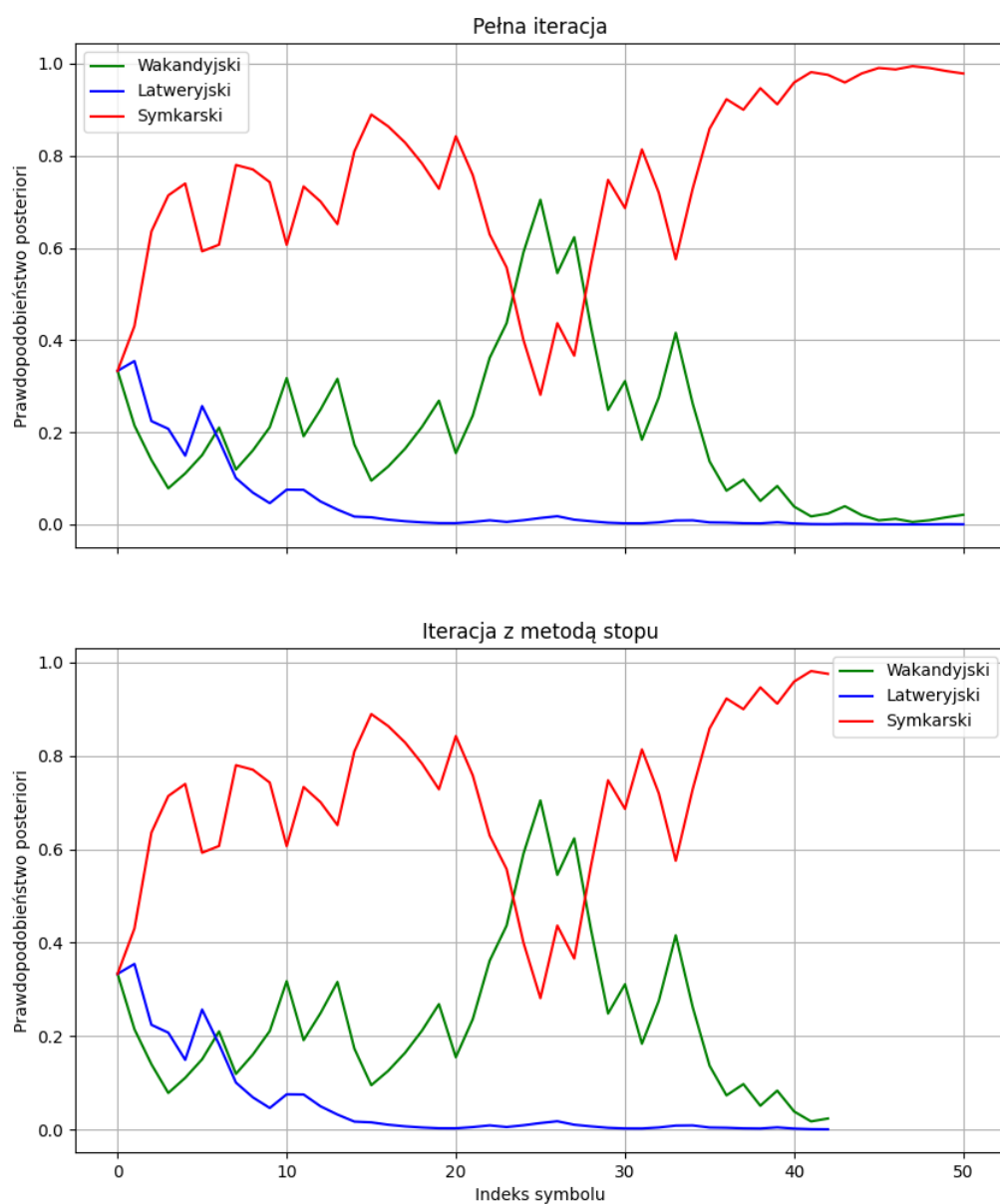


Figure 5: Porównanie pełnej iteracji i metody stopu.

Wnioski:

- Obie metody prowadzi do takich samych wartości prawdopodobieństw posteriori w przypadku długich wiadomości.
- Metoda stopu pozwala na znaczna oszczędność iteracji, szczególnie przy konwergencji modelu w pierwszych etapach analizy.
- Standardowa metoda jest preferowana w sytuacjach, gdy pełna dokładność jest

kluczowa, a czas obliczeń nie stanowi ograniczenia.

- Metoda stopu jest bardziej efektywna w praktycznych zastosowaniach, gdzie liczba danych jest duża, a konwergencja jest wystarczającym kryterium analizy.

4 Modyfikacja wiadomości

4.1 Przypadek IV: Zmodyfikowana wiadomość

W przypadku IV dokonano modyfikacji wiadomości. Spośród liter alfabetu $\{A, B, C, D, E, F\}$ wybrano dwie litery: C oraz D . Wszystkie pozostałe symbole w wiadomości zostały zastąpione symbolem N , co oznaczało dowolną literę spoza wybranych.

4.1.1 Założenia i cel analizy

Celem analizy było sprawdzenie wpływu redukcji alfabetu wiadomości na przebieg procesu aktualizacji Bayesowskiej oraz na ostateczne wartości prawdopodobieństw posteriori. Rozważono różne priory początkowe:

- Równy rozkład a priori: $P(H_W) = P(H_L) = P(H_S) = \frac{1}{3}$,
- Preferencja dla języka łatweryjskiego: $P(H_L) = 0.6, P(H_W) = P(H_S) = 0.2$,
- Preferencja dla języka wakandyjskiego: $P(H_W) = 0.6, P(H_L) = P(H_S) = 0.2$,
- Preferencja dla języka symkarskiego: $P(H_S) = 0.6, P(H_W) = P(H_L) = 0.2$.

4.1.2 Wyniki analizy

Dla każdego z priory przeprowadzono aktualizacje prawdopodobieństw posteriori, a wyniki przedstawiono na wykresach.

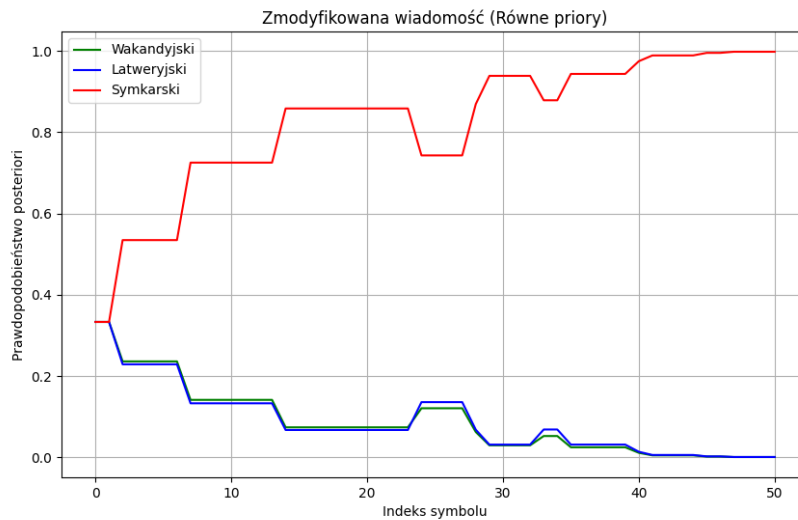


Figure 6: Zmodyfikowana wiadomość (Równe priory).

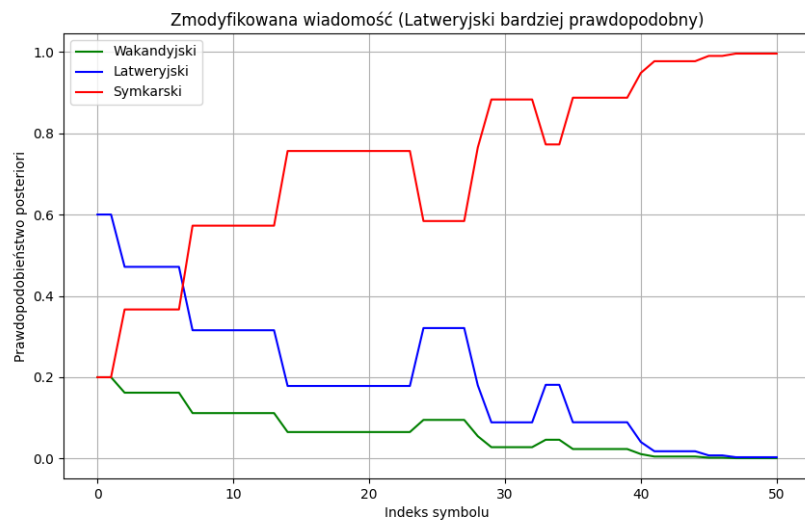


Figure 7: Zmodyfikowana wiadomość (Latweryjski bardziej prawdopodobny).

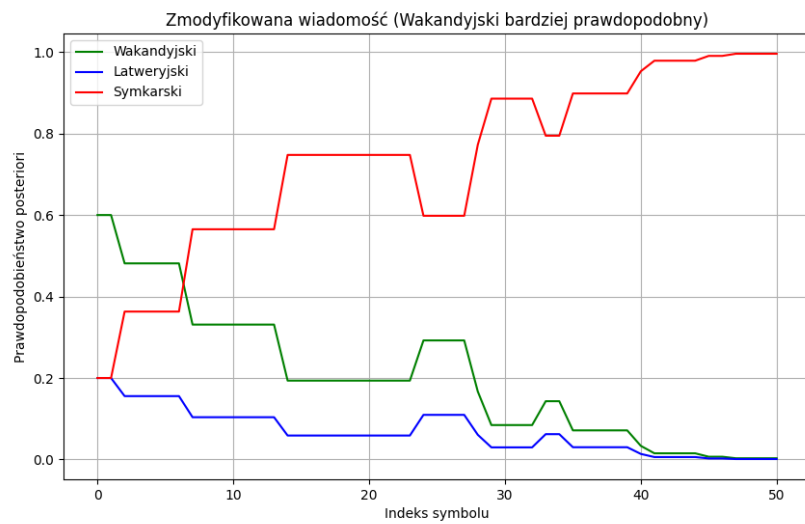


Figure 8: Zmodyfikowana wiadomość (Wakandyjski bardziej prawdopodobny).

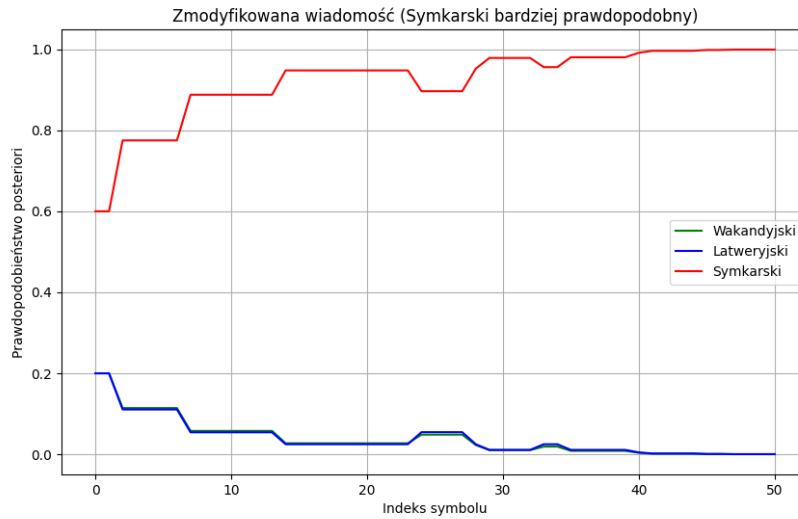


Figure 9: Zmodyfikowana wiadomość (Symkarski bardziej prawdopodobny).

4.1.3 Wnioski

Na podstawie wyników można zauważyć następujące prawidłowości:

- Redukcja alfabetu prowadzi do większej stabilności wyników. Prawdopodobieństwa posteriori szybciej konvergują do ostatecznych wartości.
- Przy równym rozkładzie a priori ($P(H_W) = P(H_L) = P(H_S) = \frac{1}{3}$), model początkowo wykazuje dużą niepewność, ale ostatecznie identyfikuje dominujący język.
- W przypadkach preferencji dla konkretnego języka (np. latweryjski, wakandyjski, symkarski), model szybciej konverguje do prawdopodobieństwa posteriori wspierającego ten język. Ostateczne wyniki są zbliżone do tych uzyskanych w przypadku pełnej wiadomości.
- Zmodyfikowana wiadomość dostarcza istotnych informacji nawet przy redukcji alfabetu, co świadczy o efektywności metody Bayesowskiej.

Warto zaznaczyć, że szczegółowe dane dla każdej iteracji znajdują się w katalogu `wyniki/csv`.