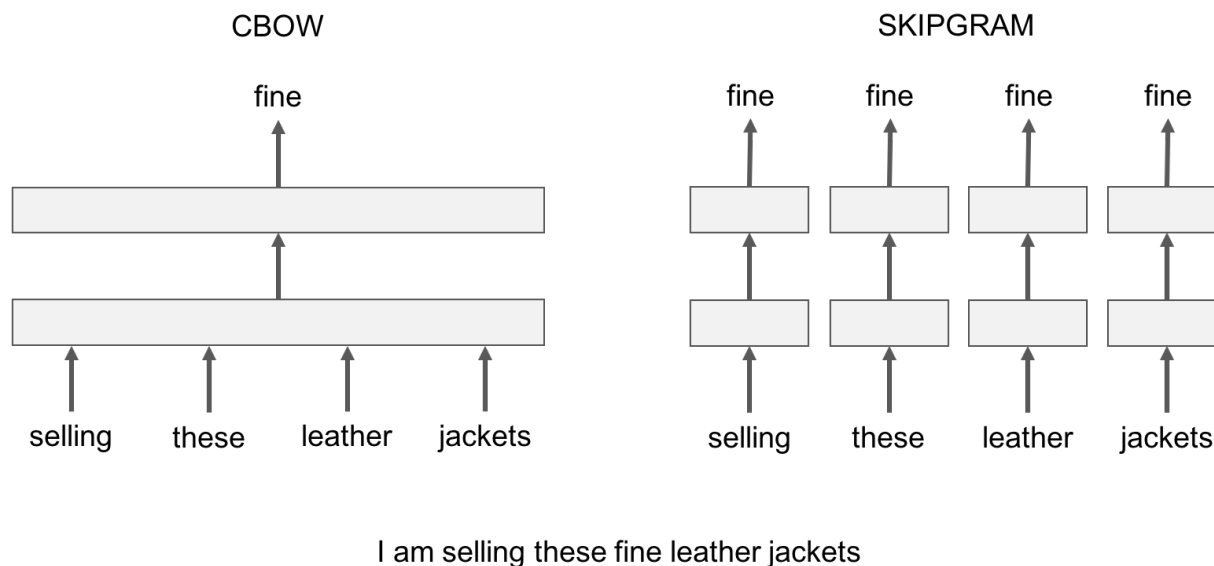


1 Modele dystrybucyjne

Oba użyte przeze mnie modele dystrybucyjne utworzyłem używając biblioteki `fasttext`. Jeden z nich to model `cbow`, a drugi to `skipgram`. Oba modele przetwarzają słowa w 100-wymiarowe wektory i używa rozmiaru okna kontekstowego 5.

`fasttext` bazuje na `word2vec`, a jego charakterystyka jest to że używa podczas uczenia informacje o znakach, z jakich składa się słowo. Działanie opiera się na dwuczęściowej sieci neuronowej składającej się z enkodera i dekodera. Nie jest to autokoder ponieważ jego celem nie jest odbudowywanie cech wejściowych. Dokładność takiego klasyfikatora ostatecznie nie jest ważna, ponieważ największe znaczenie ma jakość tworzonych wektorów reprezentacji ukrytej słowa. Usprawnieniem odróżniającym `fasttext` od `word2vec` jest wykorzystanie informacji nie tylko z kontekstu słowa, ale też ze znaków z jakich się składa. Do reprezentacji słowa dodawane są zahashowane n-gramy znaków z wnętrza danego słowa. Domyślny przedział długości tych n-gramów to od 3 do 6 znaków. Pozwala to na uzyskanie lepszych reprezentacji dla słów rzadko występujących w tekście, ponieważ o ile istnieje mało kontekstów dla takiego słowa, model może przybliżyć jego reprezentację do słów o podobnym zapisie.



1.1 cbow

W tym modelu wejściem sieci neuronowej są słowa z kontekstu dla danego słowa, a wyjściem jest samo słowo. Ten sposób reprezentacji ukryta jest uzyskiwana przez zakodowywanie kontekstu, a nie samego słowa. Mimo to wyuczenie sieci wymaga, żeby dekodery były w stanie z reprezentacji ukrytej jednoznacznie odtworzyć dane słowo.

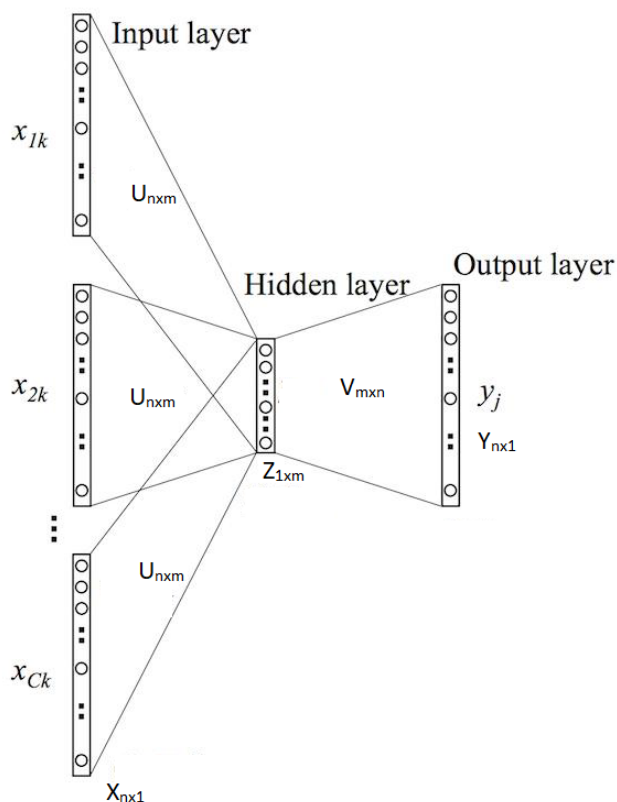


Figure 1: cbow

1.2 skipgram

skipgram działa odwrotnie do cbow. Zakodowuje on słowo do postaci ukrytej, a następnie odtwarza jego kontekst. Wymusza to podobna reprezentacje dla słów używanych w podobnych kontekstach.

2 Regresja

Do regresji użyłem własnej implementacji rekurencyjnej sieci neuronowej. Wybrałem taką architekturę sieci ze względu na zmienną długość danych wejściowych. Sieć składa się z 3 warstw rekurencyjnych (z warstwami ReLU pomiędzy) i jednej warstwy liniowej. Warstwy rekurencyjne na wejściu dostają wektory długości 100 i zwracają wektor i stan ukryty też długości 100. Warstwa liniowa przyjmuje wektor długości 100 i zwraca pojedynczą liczbę będącą przewidzianą punktacją posta/komentarza na Reddicie. Sieć uczona była przez 50 epok. Poniżej są wykresy miar MSE i R2 po każdej epoce. Model m1 to cbow, a m2 to skipgram.

Wykresy miary MSE są pokazane osobno, ponieważ początkowy stan MSE jest mocno zależny od początkowych wartości parametrów regresora, które są losowe. Duże różnice uniemożliwiają bezpośrednie porównanie miary MSE w trakcie uczenia. Można tylko patrzeć na kształt wykresów.

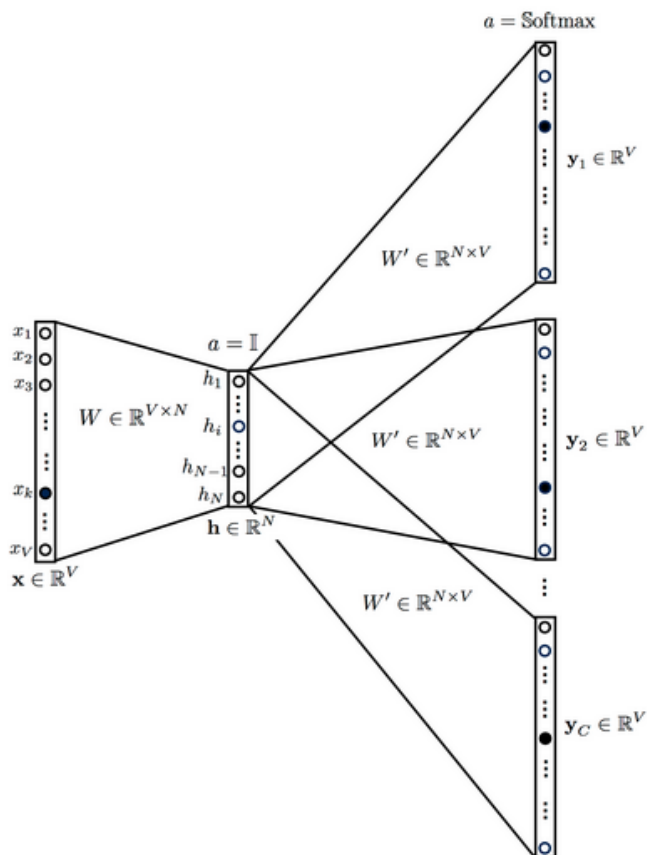


Figure 2: skipgram

Jak widać sieć ma dużą skłonność do przeuczania się. Wartości metryk dla zbioru walidacyjnego pogarszają się mniej więcej w takim tempie, jak polepszają się dla zbioru treningowego. Dla zbioru treningowego znacznie lepiej radzi sobie model **cbow**. Oba modele dają lepsze MSE jeżeli są wyuczone na korpusie wzorcowym. Oba zjawiska widać zarówno na wykresie R2, jak i wykresach MSE.

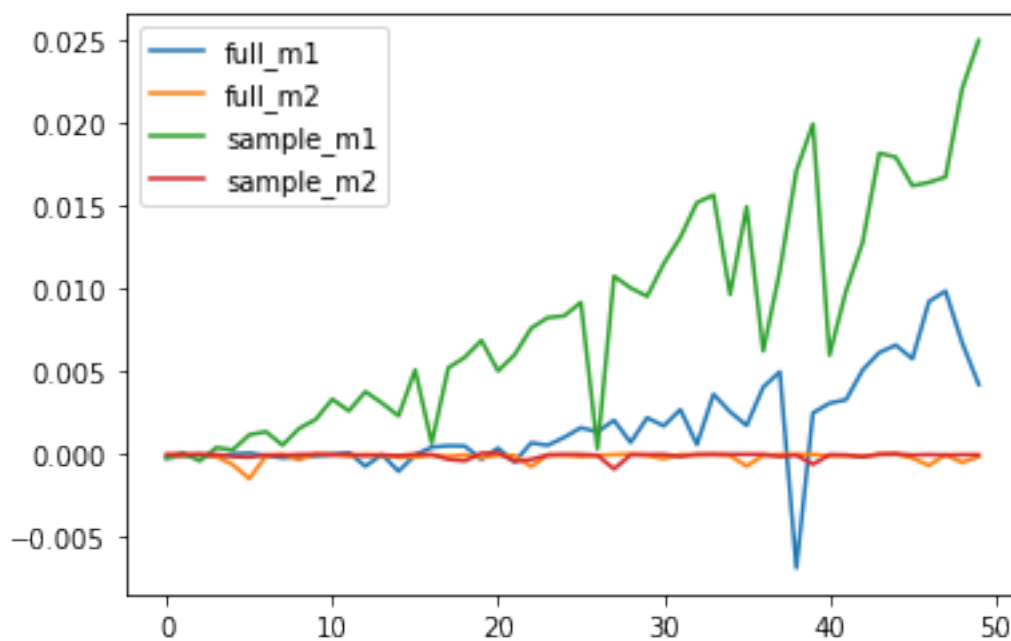
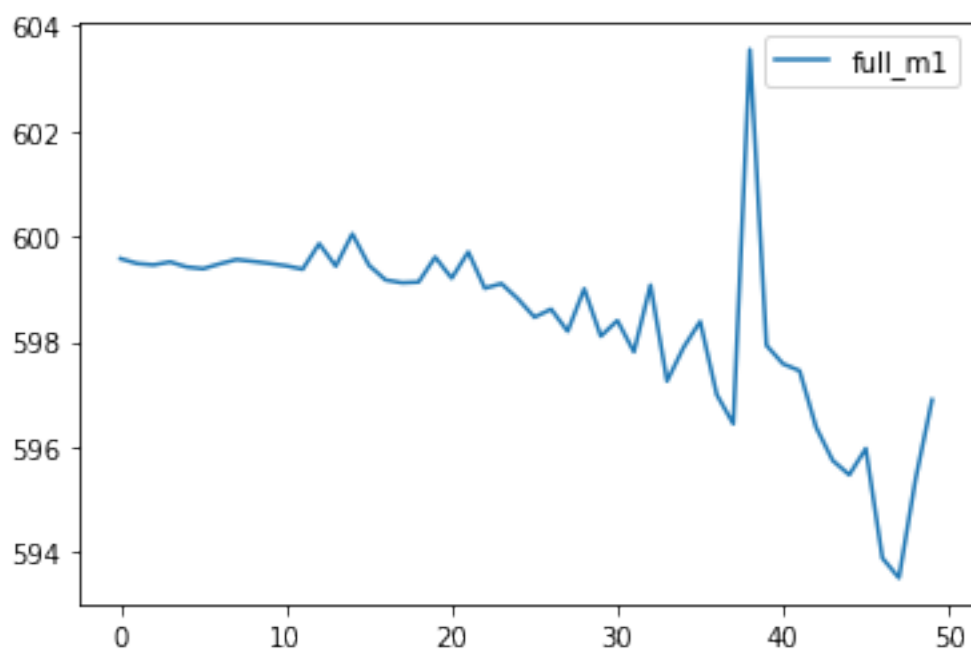
Figure 3: r^2 score dla zbioru treningowego

Figure 4: mse score dla zbioru treningowego

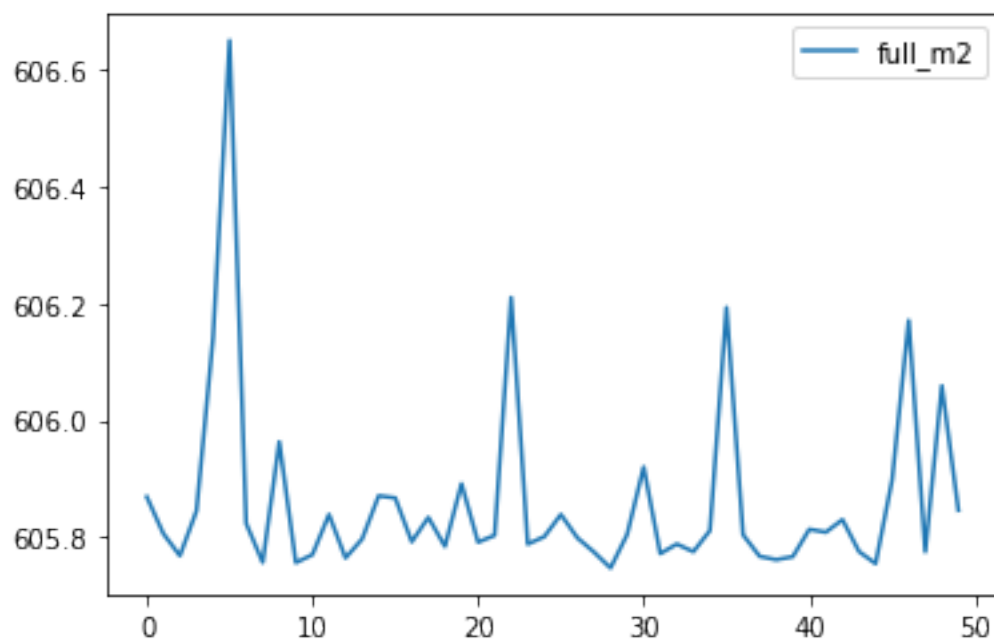


Figure 5: mse score dla zbioru treningowego

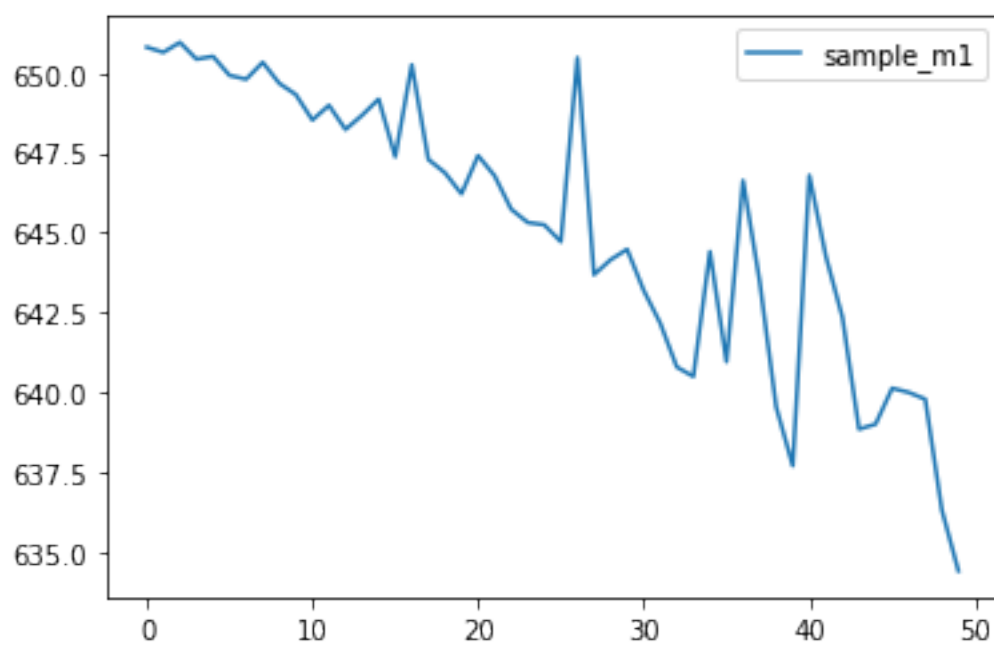


Figure 6: mse score dla zbioru treningowego

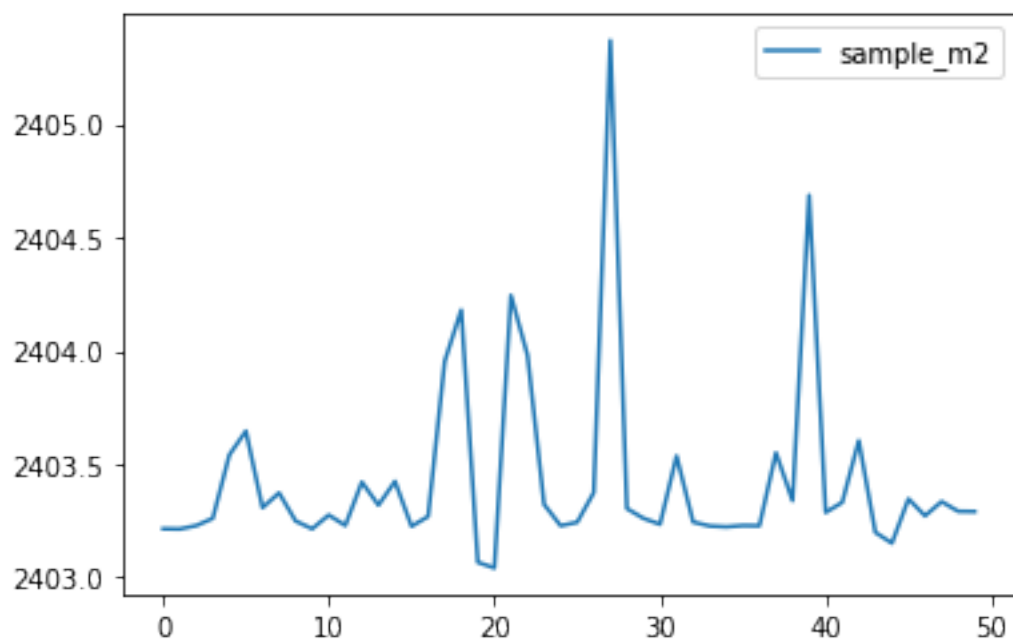


Figure 7: mse score dla zbioru treningowego

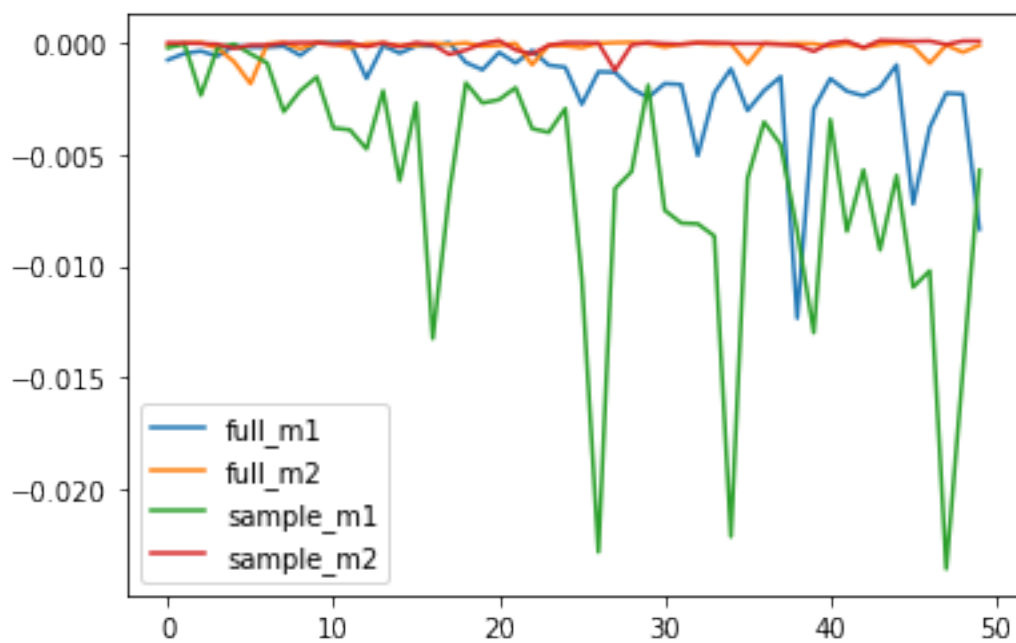


Figure 8: r2 score dla zbioru walidacyjnego

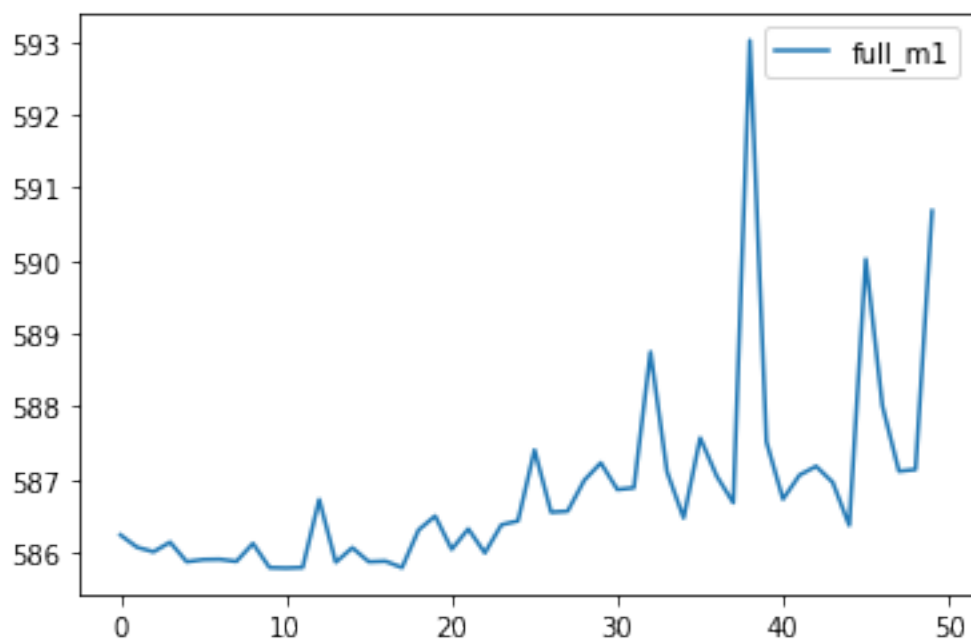


Figure 9: mse score dla zbioru walidacyjnego

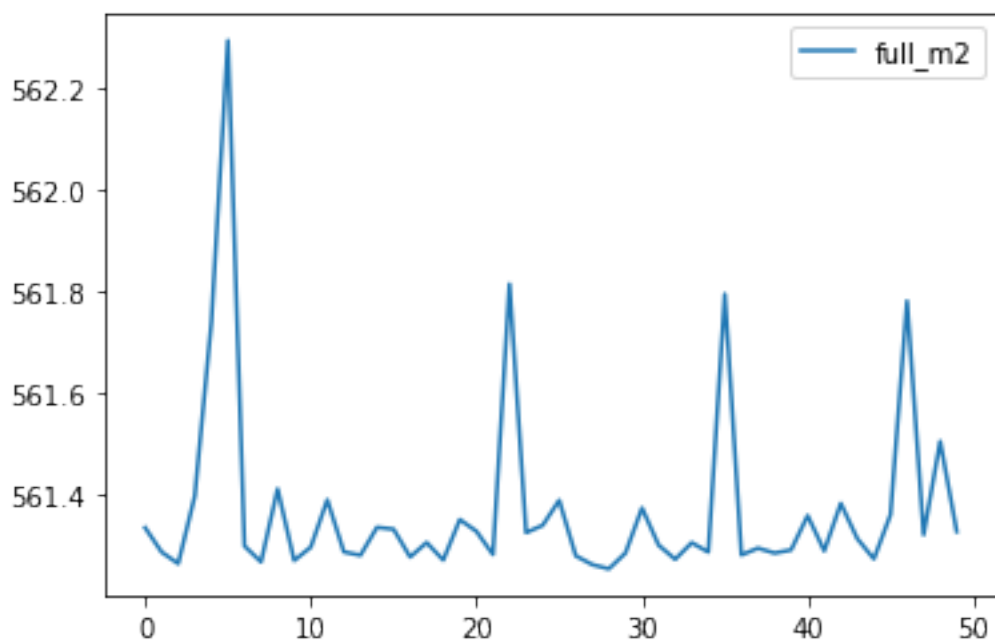


Figure 10: mse score dla zbioru walidacyjnego

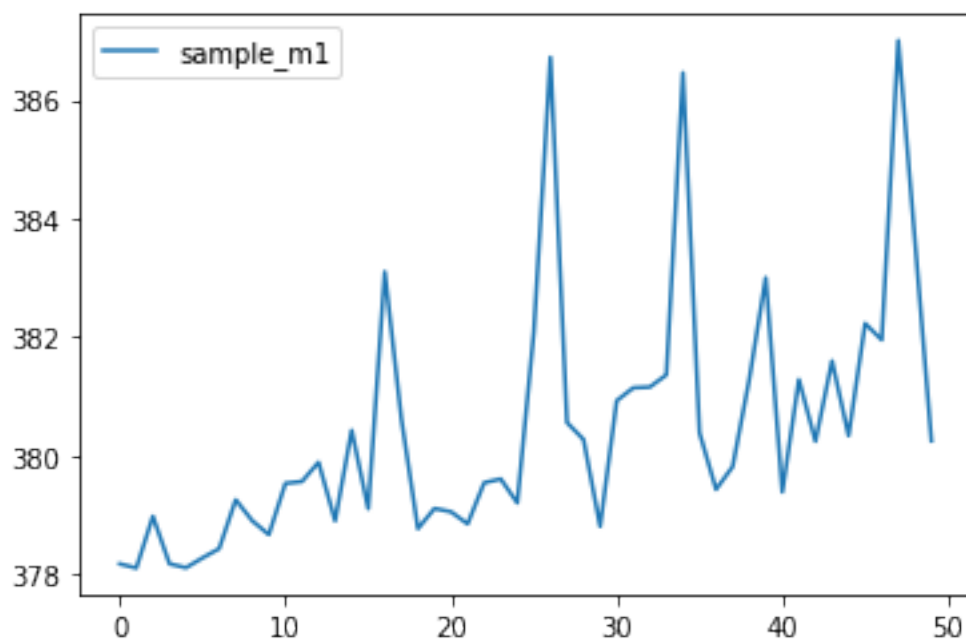


Figure 11: mse score dla zbioru walidacyjnego

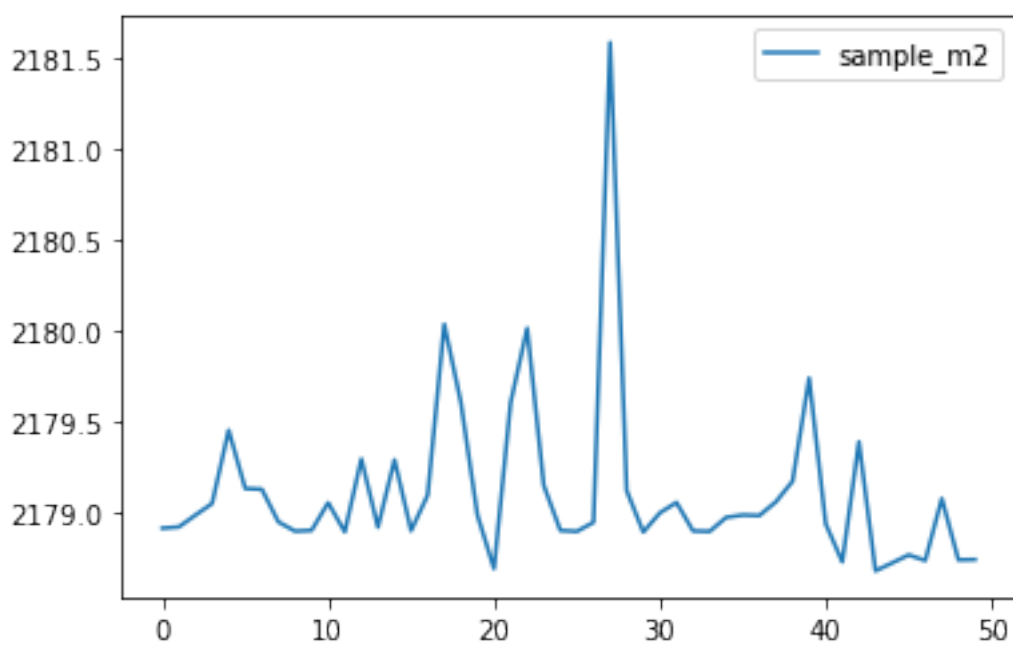


Figure 12: mse score dla zbioru walidacyjnego