

### Streszczenie

W poniższym dokumencie przedstawiam moją drogę prób i błędów podczas pisania projektu regresji liniowej z Metod Probabilistycznych w Uczeniu Maszynowym. Ten dokument będzie się składał z dwóch części - w pierwszej wycisnę do maksimum możliwości model nie używając funkcji bazowych - czyli traktując wynik jako kombinację liniową cech. W drugiej części dokumentu spróbuję przeanalizować zależności między cechami i na tej podstawie jak najlepiej dopasować naszą krzywą do punktów

## Przygotowania i pierwsze kroki

Projekt zacząłem od napisania funkcji wczytującej dane. Za każdym razem funkcja losuje kolejność danych, które następnie dzieli na zbiory: treningowy, walidacyjny i testowy. Stwierdziłem, że na dobry początek odpowiednim podziałem będzie 0.6 : 0.2 : 0.2, zatem w takiej proporcji początkowo testowałem model. Również początkowo zdecydowałem się użyć kwadratowej funkcji błędu, która dla ustalonych funkcji bazowych liczy połowę średniego kwadratu odległości predykcji od faktycznej wartości.

$$l(\theta) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - \sum_j \theta_j \phi_j(x^{(i)}))^2$$

Żeby przeskalować dane zdecydowałem się użyć standaryzacji. Niestety nie jestem wybitnym programistą...

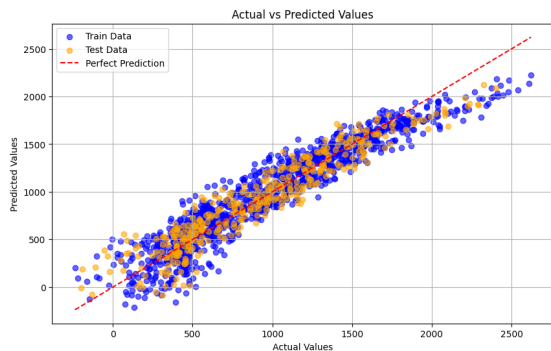
## Standaryzacja danych

Smutny dla mnie rozdział. Bardzo bolesny. W poniedziałek, 7 kwietnia, po prawie tygodniu prac nad projektem odkryłem przyczynę moich dotychczasowych trudności. Dlaczego nie mogę zbić MSE? Dlaczego **learning rate** jest tak niestabilny numerycznie przy doborze różnych funkcji bazowych? Słuchając znajomych chwalaących się swoimi MSE zauważyłem ciekawą zależność – ich MSE na zbiorze **treningowym** było **dużo niższe**, niż moje rozwiązanie analityczne, mimo użycia tej samej funkcji straty. Przegrzebałem cały ten kod i znalazłem – literówka funkcji odpowiedzialnej za standaryzację danych. Skutkowało to tym, że jak przed standaryzacją model miał trudności z nauczeniem się na danych, tak po niej przestawały one mieć jakikolwiek sens.

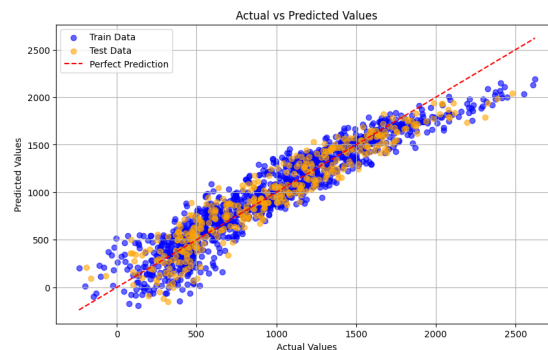
Dlaczego standaryzacja danych jest tak ważna ?

Zamiast mierzyć dane wejściowe przez ich wartość, w momencie gdy mamy wystarczająco duży zbiór treningowy lepiej jest mierzyć inną wartość – jaką wielokrotność odchylenia standardowego przyjmuje nasza cecha od jej średniej wartości? Dzięki temu ustaliśmy wspólną miarę dla każdej z cech.

Niech  $T = \{[x_1^{(i)}, x_2^{(i)}, \dots, x_7^{(i)}, y^{(i)}]\} = \{(x^{(i)}, y^{(i)})\}$  będzie naszym zbiorem treningowym, Definiuję następujące przekształcenie standaryzujące  $s : T \rightarrow S$  w taki sposób, że  $s(x^{(i)})_j = \frac{x_k^{(i)} - \mathbb{E}[x_j]}{\sqrt{\text{Var}(x_j)}}$  oraz  $s(y^{(i)}) = y^{(i)}$ , gdzie wartość oczekiwana i wariancja jest brana względem rozkładu jednostajnego nad  $T$ . W ten sposób tworzę zbiór ustandaryzowanych danych  $S$ , na którym będę trenować mój model. W momencie gdy będę mierzył skuteczność mojego modelu na zbiorze testowym (lub wykorzystywał zbiór walidacyjny), to dane znajdujące się w nim będę mierzył miarą mojego **zbioru treningowego** - będę przekształcał używając policzonych na nim średnich i odchyłeń standardowych.



Rysunek 1: Dane oryginalne, MSE: 13607.09



Rysunek 2: Dane ustandaryzowane, MSE: 13675.66

Standaryzacja jednak w klasycznym modelu regresji liniowej nie ma aż tak dużego znaczenia - ta której użyłem jest przekształceniem liniowym danych wejściowych, zatem nasz model może się nauczyć innych parametrów co doprowadzi go do podobnych rezultatów. Powyżej przedstawiam porównanie wyników rozwiązania analitycznego na zbiorze testowym danych oryginalnych i ustandaryzowanych.

Standaryzacja danych użyteczna jest jednak w momencie, gdy będziemy chcieli szerzej je analizować. Gdy ustandaryzujemy nasze dane, możemy na podstawie rozwiązania problemu regresji liniowej wnioskować (po rozmiarze parametrów) od jakich cech nasz wynik zależy najbardziej. Dodatkowo na ustandaryzowane dane możemy swobodnie aplikować funkcje bazowe i regularizacje, jednocześnie nie przejmując się tym że skala konkretnych cech przyćmi pozostałe.

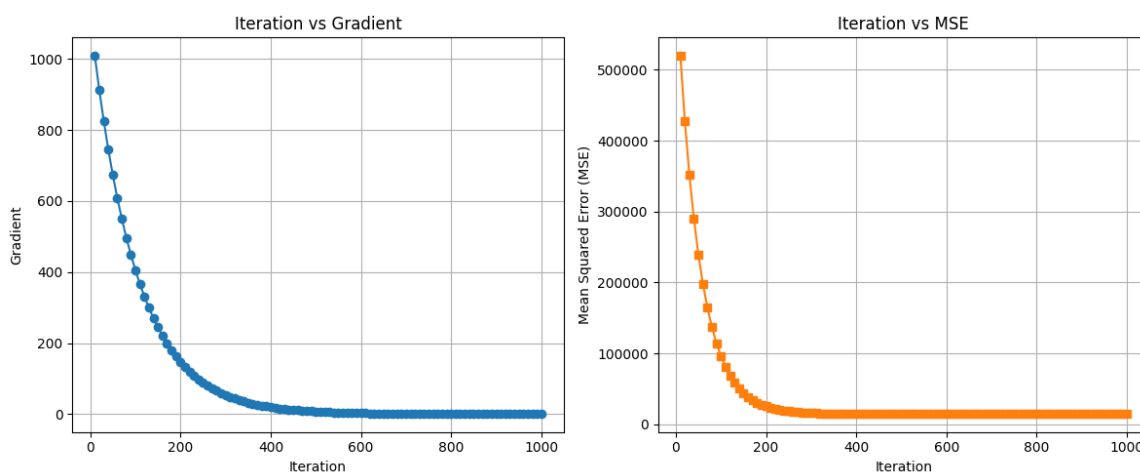
## Algorytmy

Podczas tego projektu zaimplementowałem dwie wersje algorytmu spadku wzdłuż gradientu. Pierwsza używała pełnego zestawu danych treningowych - jedna iteracja to jedna epoka. W drugiej użyłem metody gradientu mini-batch. Przy **poprawnie** ustandaryzowanych danych odpowiednim **learning rate** okazał się być  $\eta = 0.01$  - po prostu gradient przy doborach różnych funkcji bazowych zbiegał odpowiednio szybko, ale nie był na tyle duży żeby algorytm przeskakiwał minimum czy wystrzeliwał do nieskończoności. Algorytm (de facto jeden, w dwóch trybach) miał różne warunki stopu - można ustalić czy uzna model wytrenowany po  $k$  iteracjach, czy w momencie gdy norma gradientu spadnie poniżej ustalonej liczby  $\mu$ .

Na zajęciach wyprowadziliśmy fakt, że jeśli  $\theta$  minimalizuje funkcję błędu z normalizacją grzbietową o parametrze  $\lambda$ , to  $\theta = (X^T X + \lambda I)^{-1} X^T y$ , przy czym odwrotność tej macierzy zawsze istnieje dla  $\lambda > 0$ . Dlatego napisałem również algorytm wyznaczający rozwiązanie analityczne, przy czym macierz  $X^T X$  nie zawsze jest odwracalna więc żeby wyznaczyć rozwiązanie bez normalizacji biorę  $\lambda \rightarrow 0$ . Poniżej zamieszczam wyniki otrzymane na zbiorze testowym za pomocą rozwiązania analitycznego dla różnych  $\lambda$  (biorąc jako zbiór treningowy 60% danych wejściowych)

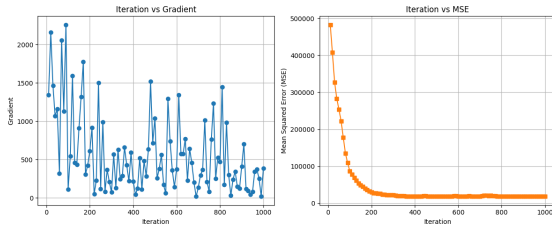
W dalszej części tej sekcji chciałbym porównać skuteczność algorytmu **mini-batch** w porównaniu do klasycznego algorytmu regresji liniowej. Na początku sprawdzimy, jak sprawuje się klasyczny algorytm. Poniższe wykresy przedstawiają zmianę gradientu oraz MSE w kolejnych iteracjach. W tym wywołaniu dostałem MSE na zbiorze treningowym: 14124.78, na testowym: 14630.03, a gradient zbiegł do 0.04.

$\lambda$	Train Set MSE	Test Set MSE
$1 \times 10^{-6}$	14358.6917	14980.0649
0.1	14358.6957	14980.8088
0.2	14358.7079	14981.5606
0.5	14358.7931	14983.8635
1	14359.0972	14987.8599
2	14360.3111	14996.4448
5	14368.7628	15026.8995
10	14398.6450	15093.0902
20	14515.9055	15281.6319
50	15294.8474	16264.2287
100	17821.7851	19099.6548
1000	135659.7145	139351.7539
100000	590708.4282	595069.9019

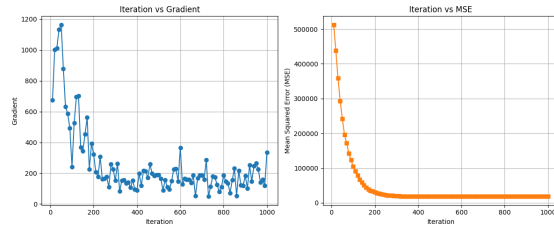
Tabela 1: MSE rozwiązania analitycznego dla różnych wartości  $\lambda$ Rysunek 3: Gradient i MSE dla kolejnych iteracji, bez **mini-batch**

Jeśli chodzi o algorytm **mini-batch**, to ma on pewne ciekawe cechy które odkryłem podczas jego używania. Przez to że wybieramy dane w losowy sposób, to próbka którą dostajemy **względnie dobrze** symuluje średnią po wszystkich danych testowych. Jednak jest on ekstremalnie szybszy (miałem sytuacje gdy mini-batch trenował model podobnie efektywnie co zwykły gradient descent, ale robił to w 10 sekund zamiast w 2 minuty). Poniżej porównanie kilku przebiegów algorytmu **mini-batch** dla różnych ustawień rozmiaru pakietu.

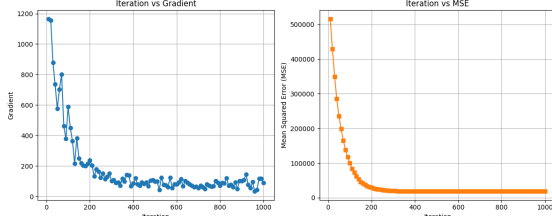
Przez to, że w każdym przebiegu iteracji losujemy podzbiór danych na których się uczymy, to gradient w odróżnieniu do wersji bez **mini-batch** nie zbiega do 0, tylko drga zależnie od wylosowanego podzbioru. Nie przeszkadza to jednak w tym, aby nasz model poprawnie dopasowywał dane. Żeby zrównoważyć efektywność **mini-batch**, wybrałem `batch-size= 32` do dalszych testów.



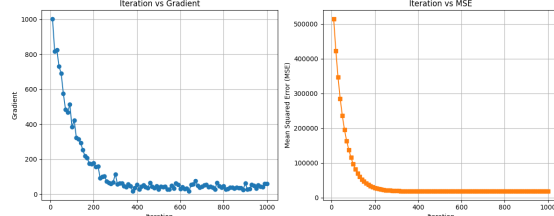
Rysunek 4: batch size= 1, MSE= 14957.47



Rysunek 5: batch size= 8, MSE= 14740.81



Rysunek 6: batch size= 32, MSE= 14584.44

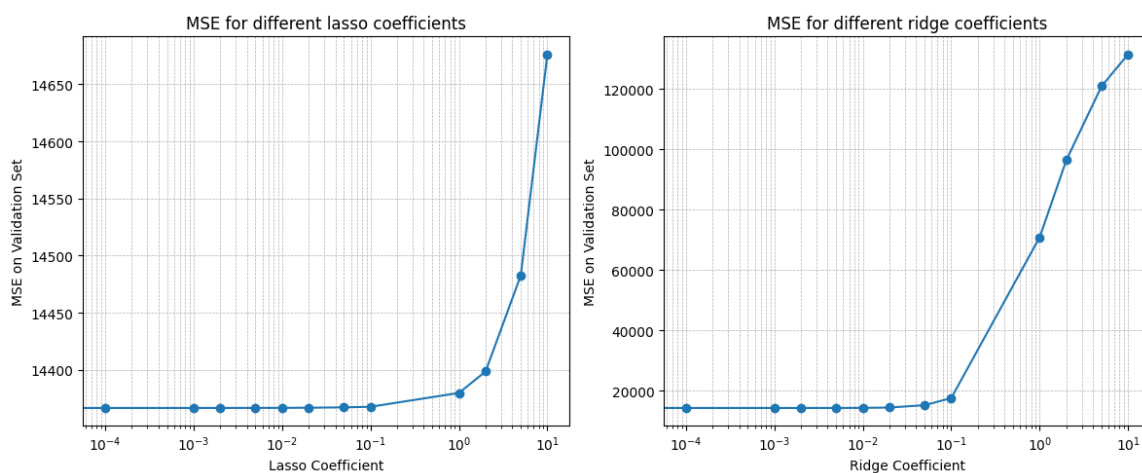


Rysunek 7: batch size= 128, MSE= 14609.96

## Regularyzacja i wyznaczanie hiperparametrów

Podczas trenowania modelu uwzględniłem również regularyzację L1 i L2 oraz regularyzację z siecią elastyczną. Aby wyznaczyć parametry regularyzacji, używam metody gradientu w wersji bez **mini-batch**, gdyż różnice dla różnych parametrów są często marginalne. Po rozdzieleniu danych na zbiór treningowy, walidacyjny i testowy porównuję wyniki na zbiorze walidacyjnym dla różnych parametrów  $\lambda_1, \lambda_2 \in \{0, 0.0001, 0.001, 0.002, 0.005, 0.01, 0.02, 0.05, 0.1, 1, 2, 5, 10\}$ .

Okazuje się że gdy nie rozważamy funkcji bazowych to regularyzacja nie daje prawie żadnej poprawy gdy chodzi o przewidywanie naszych danych, gdyż nawet na danych treningowych ma problemy z dopasowaniem odpowiedniej płaszczyzny i nie zachodzi przeuczenie. Gdy będziemy jednak trenować nasz model dla wybranych funkcji bazowych, będziemy z niej korzystać żeby algorytm mógł dopasowywać się do danych bez overfittingu. Na wykresach poniżej widzimy, że model póki co osiąga najlepsze wyniki, gdy w ogóle nie korzystamy z regularyzacji.



Rysunek 8: Wyniki dla różnych parametrów regularyzacji na zbiorze walidacyjnym

**Podział danych**

**Szersza analiza danych**