

Regresja liniowa – metoda najmniejszych kwadratów

Techniki Optymalizacji

29-30.10.2018

1 Problem regresji

Przewidywania/wyjaśnienie zmian jednej zmiennej (Y) pod wpływem zmian innych zmiennych ($\mathbf{X} = (X_1, \dots, X_m)$). W problemie regresji zmienna Y jest *ciągła*. Po co nam regresja? Zmienne \mathbf{X} zwykle łatwe do pozyskania, Y – trudne lub niemożliwe do pozyskania. Przykładowo: \mathbf{X} – ceny akcji w ostatnim tygodniu, Y – cena akcji jutro; lub: \mathbf{X} – wyniki testów medycznych, Y – poziom zaawansowania choroby.

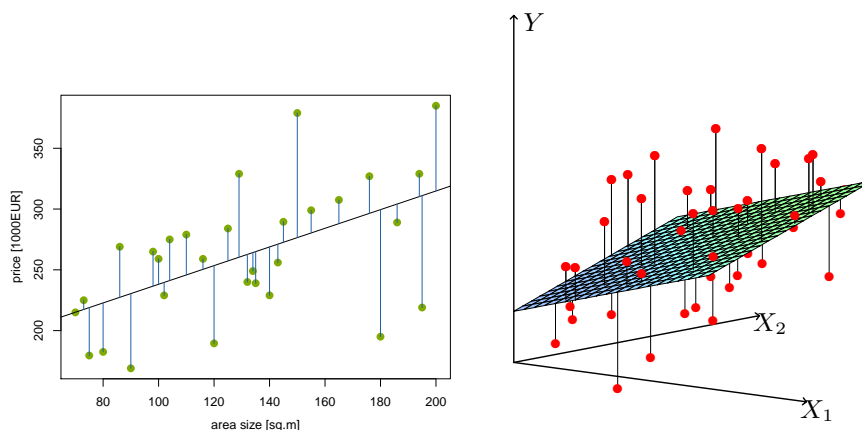
Dzisiaj zajmiemy się *regresją liniową*, kiedy modelujemy zmienną Y jako funkcję liniową \mathbf{X} :

$$\hat{Y} = w_0 + w_1X_1 + w_2X_2 + \dots + w_mX_m = \mathbf{w}^\top \mathbf{X}, \quad \text{gdzie } \mathbf{X} = (1, X_1, \dots, X_m).$$

Zwykle otrzymujemy zbiór danych składający się n obserwacji, dla których znamy wartość zmiennej Y :

$$\begin{array}{ll} (x_{11}, x_{12}, \dots, x_{1m}, y_1) & (\mathbf{x}_1, y_1) \\ (x_{21}, x_{22}, \dots, x_{2m}, y_2) & (\mathbf{x}_2, y_2) \\ \dots & \dots \\ (x_{n1}, x_{n2}, \dots, x_{nm}, y_n) & (\mathbf{x}_n, y_n) \end{array} \quad \text{lub w skrócie}$$

Próbujemy dopasować funkcję liniową (czyli wektor \mathbf{w}) do zbioru uczącego, tak aby wartości przewidywane przez model \hat{y}_i były jak najbliżej zaobserwowanych wartości y_i .



Rysunek 1: Odchylenia, tj. różnice między wartością zaobserwowaną a przewidywaną (zaznaczone liniami pionowymi). Przypadek jednej (lewy obrazek) i dwóch zmiennych wejściowych (prawy obrazek). Prawy obrazek z Hastie, Tibshirani, Friedman: *Elements of Statistical Learning*.

Odchyłeń jest n , a my chcemy mieć jedną funkcję celu. Najpopularniejszą funkcją celu jest *suma kwadratów odchyłeń*:

$$L(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2.$$

Chcemy minimalizować $L(\mathbf{w})$. W tym celu możemy użyć metody Newtona-Raphsona, która rozwiąże ten problem w jednym kroku (dlaczego? przypomnę potem na wykładzie), rozpoczynając od dowolnego \mathbf{w}_0 . Liczymy gradient i hesjan:

$$\nabla L(\mathbf{w}) = \sum_{i=1}^n -2(y_i - \mathbf{w}^\top \mathbf{x}_i)\mathbf{x}_i, \quad \mathbf{H}(\mathbf{w}) = \sum_{i=1}^n 2\mathbf{x}_i\mathbf{x}_i^\top,$$

i używamy metody Newtona-Raphsona:

$$\mathbf{w} = \mathbf{w}_0 - (\mathbf{H}(\mathbf{w}_0))^{-1} \nabla L(\mathbf{w}_0),$$

z czego po podstawieniu $\mathbf{w}_0 = \mathbf{0}$ (dla wygody, dowolne \mathbf{w}_0 by dało ten sam wynik) otrzymujemy:

$$\mathbf{w} = \left(\sum_{i=1}^n \mathbf{x}_i\mathbf{x}_i^\top \right)^{-1} \left(\sum_{i=1}^n y_i\mathbf{x}_i \right)$$

Osobliwość hesjanu. Metoda Levenberga-Marquada. Jeśli hesjan jest osobliwy, możemy dodać do niego macierz jednostkową przemnożoną przez stałą λ . Jest to metoda Levenberga-Marquada. Tak, jak metoda Newtona-Raphsona odpowiada zwykłej regresji liniowej z minimalizacją sumy kwadratów odchyłeń, tak metoda Levenberga-Marquada odpowiada *regresji grzbietowej (ridge regression)*. W regresji grzbietowej wychodzimy z problemu minimalizacji:

$$L(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|^2.$$

Liczymy gradient i hesjan:

$$\nabla L(\mathbf{w}) = \sum_{i=1}^n -2(y_i - \mathbf{w}^\top \mathbf{x}_i)\mathbf{x}_i + 2\lambda\mathbf{w}, \quad \mathbf{H}(\mathbf{w}) = \sum_{i=1}^n 2\mathbf{x}_i\mathbf{x}_i^\top + \lambda\mathbf{I},$$

i używając metody Newtona-Raphsona po podstawieniu $\mathbf{w}_0 = \mathbf{0}$ otrzymujemy:

$$\mathbf{w} = \left(\sum_{i=1}^n \mathbf{x}_i\mathbf{x}_i^\top + \lambda\mathbf{I} \right)^{-1} \left(\sum_{i=1}^n y_i\mathbf{x}_i \right).$$

Dostajemy krok metody Levenberga-Marquada. Tak jak poprzednio, rozwiązujemy problem w jednym kroku.

2 Uwagi o implementacji

Kolejne zadania laboratorium będą wykorzystywały napisaną przez Was implementację metody regresji liniowej poprzez minimalizację sumy kwadratów odchyłeń, zgodnie z tym co opisano powyżej. Przy implementacji należy zwrócić uwagę na następujące rzeczy:

- Należy napisać program tak, aby przyjmował na wejściu dowolny zbiór danych w określonym formacie.
- **Tu** znajduje się prosty skrypt napisane w Octave, który można użyć do odczytania dowolnego pliku używanego na tych zajęciach.

- Naturalne jest wczytanie danych jako macierz \mathbf{X} o rozmiarze $n \times m$, której poszczególne wiersze to obserwacje $\mathbf{x}_1, \dots, \mathbf{x}_n$ (tak wczytuje to powyższy skrypt). Czyli $(\mathbf{X})_{i,j}$ to wartość j -tej zmiennej wejściowej dla i -tej obserwacji. Wtedy macierz $\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$ (główny element hesjanu) można zapisać jako $\mathbf{X}^\top \mathbf{X} \dots$
- ... ale uwaga: do macierzy \mathbf{X} po wczytaniu należy dodać kolumnę jedynek aby uwzględnić wyraz wolny!
- Podobnie, naturalne jest wczytanie wartości zmiennej wyjściowej jako wektora \mathbf{y} (tak wczytuje to powyższy skrypt). Wtedy $\sum_{i=1}^n y_i \mathbf{x}_i$ można zapisać jako $\mathbf{X}^\top \mathbf{y}$ (główny element gradientu).
- Trzeba się upewnić, czy hesjan nie jest osobliwy, jeśli jest, należy użyć metody Levenberga-Marquada. Jak sprawdzić osobliwość macierzy?
- Można użyć **tego zbioru danych** do przetestowania swojego programu. Hesjan nie jest tu osobliwy, a wartości współczynników to: $w_0 = 0.33795, w_1 = 0.97572, w_2 = 0.60207$. Z kolei **ten zbiór danych** powinien dać osobliwy hesjan, a wartości współczynników dla $\lambda = 10^{-5}$ to $w_0 = 0.036364, w_1 = -0.020053, w_2 = -0.033422$.

3 Sztuczne zbiory danych

Celem zadania jest implementacja metody regresji liniowej poprzez minimalizację sumy kwadratów odchyłeń (zgodnie z uwagami wymienionymi powyżej), a następnie użycie jej na dwóch małych, sztucznych zbiorach danych: **pierwszy zbiór** oraz **drugi zbiór**. Format obu plików jest taki sam: pierwszy wiersz zawiera nazwy zmiennych oddzielone spacją, a każdy kolejny wiersz to opis jednej obserwacji z listą wartości zmiennych (również oddzielone spacją).

Odpowiedz na następujące pytania:

1. Jakie są wartości współczynników \mathbf{w} ?
2. Czy hesjan był osobliwy? Jeśli tak to dlaczego i jak można sobie z tym poradzić?

4 Rzeczywisty zbiór danych

Celem zadania jest użycie metody regresji liniowej zaimplementowanej w poprzednim zadaniu na zbiorze danych wyceny domów w Den Bosch. Zbiór zawiera 119 domów, opisanych za pomocą 10 zmiennych (9 wejściowych i jedna wyjściowa – cena domu). Więcej o problemie, w tym opis zmiennych, można znaleźć na slajdach z wykładu. Zbiór został podzielony na dwie części – treningową (80 domów) i testową (39 domów), które można pobrać z: **zbiór treningowy** oraz **zbiór testowy**.

Format obu plików jest tak sam: pierwszy wiersz zawiera nazwy zmiennych oddzielone spacją, a każdy kolejny wiersz to opis jednego domu, z listą wartości zmiennych (również oddzielone spacją). Zmienna wyjściowa, którą będziemy modelować ("price") znajduje się na samym końcu listy, pozostałe zmienne są zmiennymi wejściowymi.

Należy zastosować metodę regresji liniowej poprzez minimalizację sumy kwadratów odchyłeń *do zbioru treningowego*, uzyskując wartości współczynników. Następnie, należy użyć wartości współczynników do *predykcji* wartości cen domów na zbiorze testowym. Odpowiedz na następujące pytania:

1. Jakie są wartości współczynników \mathbf{w} ? Jak można je interpretować?
2. Czy hesjan był osobliwy?
3. Jaka jest średnia wartość błędu na zbiorze treningowym, a jaka na zbiorze testowym? Spróbuj zinterpretować różnicę.

Uwaga: jako średnią wartość błędu na danym zbiorze przyjmij tak zwany *pierwiastkowany średni błąd kwadratowy (root mean squared error)*, wyznaczony jako:

$$\text{err} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}.$$

Jeśli obserwacje znajdują się jako wiersze w macierzy \mathbf{X} , to wektor predykcji można wyznaczyć jako $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$.

Zadanie dla chętnych: Jak wybór wartości λ wpływa na wynik (sumaryczny błąd) na zbiorze treningowym i testowym? Rozwiąż problem dla wielu wartości λ wyznaczonych jako kolejne potęgi dwójki, np. $\lambda = 2^{-15}, 2^{-14}, 2^{-13}, \dots$, aż do $\lambda = 2^{20}$. Dla każdej takiej wartości wyznacz błąd na zbiorze treningowym i testowym, a następnie umieść oba błędy na wykresie w zależności od λ (na osi X zastosuj skalę logarytmiczną). Zinterpretuj wyniki.

5 Pytania

Odpowiedz na następujące pytania:

- Jaka będzie wartość współczynnika regresji, gdy mamy tylko jedną cechę wejściową x (i brak wyrazu wolnego)? W szczególności, jak wygląda taki współczynnik, gdy cecha wejściowa jest stale równa jeden?
- Jak będzie wyglądał ten współczynnik w przypadku regresji grzbietowej? Zinterpretuj zmiany.
- Spróbuj wyznaczyć wzór na prostą regresję liniową jednej zmiennej z wyrazem wolnym:

$$w_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2}, \quad w_0 = \bar{y} - w_1 \bar{x}.$$

- Rozważ regresję poprzez minimalizację sumy wartości bezwzględnych. Spróbuj udowodnić, że jeśli mamy tylko wyraz wolny, to przewidywana wartość \hat{y} będzie równanie medianie z danych.