

Metody Bayesowskie - Projekt 2

Estera Saidlo, Wiktoria Szczypka, Paweł Warchoł

25 maja 2020

Cel projektu i opis danych

Celem projektu jest zastosowanie bayesowskiego w regresji wielokrotnej na podstawie danych "Boston Housing", które są zbiorowe w R i znajdują się w pakiecie mlbench. Dane dotyczą wartości domów na przedmieściach Bostonu.

Zadania podzielone na 3 części:

- Dane, na podstawie których zostanie stworzony model MNK, który poskrobi do wyznaczenia rozkładu a priori (100 obserwacji).
- Dane użycie do części wstępnej badania (400 obserwacji).
- Dane użycie do prognozy (6 obserwacji).

Prognozy uzyskane za pomocą wnioskowania bayesowskiego zostaną porównane z prognozami Metody Najmniejszych Kwadratów.

Zmienna objaśniana:

MEDV - mediana wartości domów na przedmieściach Bostonu (w tysiącach dolarów).

Zmienne objaśniające:

RM - średnia liczba pokoi w domu.

CRIM - wskaźnik przestępczości na mieścinach.

LSTAT - procent populacji o niższym statusie społecznym.

PTRATIO - stosunek liczby uczniów do nauczycieli w danym obszarze podmiejskim.

INDUS - odsetek niedziałalnych akwów biznesowych w danym obszarze podmiejskim.

#	medv	rm	crim	lstat
#	Min.	5.00	Min.	5.00632
#	1st Qu.	117.02	1st Qu.	5.886
#	Median	21.20	Median	6.208
#	Mean	32.53	Mean	6.285
#	3rd Qu.	126.00	3rd Qu.	6.623
#	Max.	56.00	Max.	18.780
#	ptratio	indus		
#	1st Qu.	137.40	1st Qu.	5.19
#	Median	116.05	Median	5.69
#	Mean	131.45	Mean	5.114
#	3rd Qu.	126.20	3rd Qu.	5.139
#	Max.	122.00	Max.	12.74

Zmienna objaśniana znajduje się w przedziale od 5 do 50 tysięcy dolarów. Wartość maksymalna wydaje się być wartością skrajną, gdyż 3 kwantyl jest na poziomie 25 tysięcy dolarów. Oznacza to, że 75% obserwacji przypada na wartości 25 tysięcy dolarów lub mniej. Średnia liczba pokoi w domu znajduje się w przedziale między 3,5 a 8,5. Średnio w domu znajduje się ok. 6 pokoi. Średni wskaźnik przestępczości na mieszkanka wynosi 3,2, co oznacza, że średnio na 1 mieszkańca przedmieść Bostonu przypadać między 3 a 4 przestępstwa. Wartość maksymalna to prawie 89 przestępstw na jednego mieszkańca, jednak tutaj także można zauważyć, że wartość ta będzie outlierem. Średnio jest 12% ludu o niższym statusie społecznym, maksymalnie jest to prawie 38%, a minimum mniej niż 2%. Stosunek liczby uczniów do liczby nauczycieli jest między 12,6 a 22. Średnio na 1 nauczyciela przypada ok. 18 uczniów. Odsetek niedziałalnych akwów biznesowych znajduje się w przedziale od 0,26% do 27,7%. Średnio jest to 1,14% akwów w miastie.

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.458 on 94 degrees of freedom
## Multiple R-squared:  0.8368, Adjusted R-squared:  0.8281
## F-statistic: 96.41 on 5 and 94 DF, p-value: < 2.2e-16
```

Model MNK stworzony za pomocą Metody Najmniejszych Kwadratów składowy do oszacowania parametrów rozkładu a priori:

Rest z modelu

Code

Model MNK

Został stworzony model MNK na danych przeczyszczonych do właściwej części badania. Na jego podstawie zostanie sprawdzone spełnienie zakładań dotyczących restu oraz dokonana będzie prognoza dla danych z części 3.

Code

```
##
## Call:
## lm(formula = medv ~ ., data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -14.379  -9.952  -1.033   1.702  23.576
##
## Coefficients:
## (Intercept) Estimate Std. Error t value Pr(>|t|)
## (Intercept) 20.16234      4.60995    4.375 1.55e-05 ***
## rm          4.18465      0.40222    10.392 5.92e-16 ***
## crim        -0.07333     0.03449   -2.120  0.0341 *
## lstat       -0.59282     0.06030   -9.818 < 2e-16 ***
## indus        0.08274     0.44633    0.185 0.8560999 ***
## indus       0.01442     0.05306    0.272  0.7860
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Zmienney najbardziej wpływający na oszacowanie wartości domów w Bostonie w modelu nr. 1 są: średnia liczba pokoi w domu oraz odsetek ludu o niższym statusie. Średnia liczba pokoi jest jedną zmienną, której wzrost powoduje wzrost zmiennej objaśnianej. Wzrost pozostałych zmiennych wpływa negatywnie na wzrost wartości domów. Jednak stosunek liczby do nauczycieli nie jest istotny w kontekście kosztowności w jednym mieszkaniu, jednak tutaj także można zauważyć, że informacja dotycząca tej zmiennej w rozkładzie a priori nie będzie pomocna. Nieistotna w tym modelu okazała się zmienna dotycząca odsetka niedziałalnych akwów biznesowych na danym obszarze podmiejskim, natomiast w poprzednim modelu była istotna na poziomie 10%. Wskaźnik przestępczości jest istotny na poziomie 5% w obu modelach.

Parametry rozkładu a priori zostały dobrane na podstawie modelu MNK stworzonego z pierwszej części danych. Wektor β_0 to oszacowane współczynniki modelu. Macierz X_0 to zmienne objaśniające pierwszej części danych, które poprzedza kolumną 1, reprezentującą stałą. Macierz $S^2(X_0)$ to macierz odwrotności macierzy $X_0^T X_0$, na to różnica kwadratów obserwacji i liczby kolumn pierwszej części danych, a σ_0^2 to suma kwadratów reszt z modelu.

Model MNK stworzony na drugiej części danych wskazuje na dużą istotność zmiennych: średnia liczba pokoi w domu, odsetek ludu o niższym statusie oraz stosunek uczniów do nauczycieli. W modelu nr. 1 stosunek uczniów do nauczycieli okazał się być nieistotny. Taka różnica wynika z niejednorodności z niehomogennego rozkładu danych. Można więc wnioskować, że informacja dotycząca tej zmiennej w rozkładzie a priori nie będzie pomocna. Nieistotna w tym modelu okazała się zmienna dotycząca odsetka niedziałalnych akwów biznesowych na danym obszarze podmiejskim, natomiast w poprzednim modelu była istotna na poziomie 10%. Wskaźnik przestępczości jest istotny na poziomie 5% w obu modelach.

- Interpretacja oszacowań parametrów:**
- Przy wzroście średniej liczby pokoi w domu o 1 jednostkę, mediana wartości domów na przedmieściach Bostonu wzrasta o ok. 4,18 tysięcy dolarów przy założeniu, że pozostałe zmienne są na stałym poziomie.
 - Przy wzroście wskaźnika przestępczości na mieszkanka o 1 jednostkę, mediana wartości domów na przedmieściach Bostonu maleje o ok. 0,07 tysięcy dolarów przy założeniu, że pozostałe zmienne są na stałym poziomie.
 - Przy wzroście procentu populacji o niższym statusie społecznym o 1 jednostkę, mediana wartości domów na przedmieściach Bostonu maleje o ok. 0,59 tysięcy dolarów przy założeniu, że pozostałe zmienne są na stałym poziomie.
 - Przy wzroście stosunku liczby uczniów do nauczycieli na danym obszarze podmiejskim o 1 jednostkę, mediana wartości domów na przedmieściach Bostonu maleje o ok. 0,87 tysięcy dolarów przy założeniu, że pozostałe zmienne są na stałym poziomie.
 - Przy wzroście odsetka niedziałalnych akwów biznesowych na danym obszarze podmiejskim o 1 jednostkę, mediana wartości domów na przedmieściach Bostonu wzrasta o ok. 0,01 tysięcy dolarów przy założeniu, że pozostałe zmienne są na stałym poziomie.

Sprawdzenie założeń dotyczących reszt modelu:

Normalność reszt została zbadana za pomocą testu Shapiro-Wilka, którego H_0 stanowi o normalności reszt.

P-value jest mniejsze od 5%, co oznacza, że H_0 należy odrzucić na rzecz H_1 , a więc reszty nie mają rozkładu normalnego. Jednak ze względu na dużą liczbę obserwacji, można założyć asymptotyczną normalność reszt.

W celu zbadania autokorelacji zostanie przeprowadzony test Durbin-Watsona, którego H_0 stanowi o braku autokorelacji reszt.

Wartość p-value wskazuje na odrzucenie H_0 na rzecz H_1 , co oznacza, że występuje autokorelacja reszt.

Heteroskedastyczność została zbadana na podstawie testu Breucha-Pagana o następujących hipotezach:

H_0 : Homoskedastyczność.
 H_1 : Heteroskedastyczność.

Wartość p-value wskazuje na to, że nie ma podstaw by odrzucić H_0 , co oznacza homoskedastyczność modelu.

Badane dane są przekrojowe, więc nie ma sensu sprawdzić ich na autokorelacji.

Rozkład a posteriori

Na podstawie rozkładu a priori wyznaczone zostaną parametry rozkładu a posteriori według poniższych wzorów.

$$\Sigma_1 = (X^T X + \Sigma_0^{-1})^{-1}$$
$$\beta_1 = \Sigma_1 (X^T Y + \Sigma_0^{-1} \beta_0) = \Sigma_1 ((X^T X) \beta + \Sigma_0^{-1} \beta_0)$$
$$\sigma_1 = \sigma_0 + n$$
$$\delta_1 = \delta_0 + y^T y - \beta_1^T \Sigma_1^{-1} \beta_1 + \beta_0^T \Sigma_0^{-1} \beta_0$$

Funkcja calc_pars zwraca obliczone parametry rozkładu.

Wartość parametru σ_1 :

Wartość parametru δ_1 :

Macierz Σ_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :

Macierz β_1 :