

Techniki eksploracji danych

Krzysztof Gajowniczek

Rok akademicki: 2020/2021

- 1 Regresja liniowa
- 2 Regresja grzbietowa
- 3 Literatura

Section 1

Regresja liniowa

Subsection 1

Estymacja analityczna

- Zwykła regresja liniowa metodą najmniejszych kwadratów próbuje dopasować linię prostą do zbioru punktów danych w taki sposób, aby suma kwadratów błędów była minimalna.

$$f(\theta) = \|\mathbf{y} - X\theta\|^2 = \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2$$

gdzie

$$X = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_i^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$

- Nie istnieje dokładne rozwiązanie dla tego problemu, więc staramy się dopasować dane jak najlepiej:

$$\hat{\theta} = \min_{\theta} \|y - X\theta\|^2$$

- Oszacowanie parametrów modelu uzyskuje się za pomocą układu równań normalnych, który wykorzystując notację macierzową ma następującą postać:

$$\hat{\theta} = (X^T X)^{-1} X^T y$$

- Sugeruje to, że $(X^T X)^{-1}$ i $X^T \mathbf{y}$ są wystarczającymi statystykami dla $\hat{\theta}$.
- Wymiarowość X to $n \times p$, dlatego $X^T X$ ma wymiar $p \times p$ oraz $X^T \mathbf{y}$ ma wymiar $p \times 1$.
- Żadne z powyższych nie jest zależne od n , więc wymiar takiej statystyki nie rośnie wraz z danymi.
- Ilość danych jest często większa niż ich wymiar, a także pozwala na ustalenie określonej ilości miejsca na przechowywanie odpowiednich statystyk, niezależnie od tego, ile danych jest podanych.

- Powyższą estymację analityczną można uzyskać także za pomocą macierzy Gramma $K = XX^T$ o wymiarze $n \times n$, będącą iloczynem skalarnym wszystkich możliwych par obserwacji.
- To jest nieparametryczne podejście. Rośnie wraz z danymi.

$$\hat{\theta} = X^T(XX^T)^{-1}\mathbf{y}$$

$$\hat{\theta} = X^T\alpha$$

- gdzie:

$$\alpha = (XX^T)^{-1}\mathbf{y}$$

- Przewidywanie wartości nowych obserwacji uzyskuję się za pomocą:

$$\hat{y}_i = \hat{\boldsymbol{\theta}}^T \mathbf{x}_i = \mathbf{x}_i^T \hat{\boldsymbol{\theta}}$$

- lub (wykorzystując estymaty uzyskane za pomocą macierzy Gramma):

$$\hat{y}_i = \mathbf{x}_i^T \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}$$

- lub:

$$\hat{y}_i = \mathbf{x}_i^T \mathbf{X}^T \boldsymbol{\alpha} = \sum_{j=1}^n \mathbf{x}_i^T \mathbf{x}_j \alpha_j$$

Subsection 2

Estymacja iteracyjna

- Zdefiniujmy na nowo funkcję błędu. Zamiast sumy kwadratów reszt:

$$f(\theta) = \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2$$

- wykorzystywać będziemy zmodyfikowany błąd średnio-kwadratowy:

$$f(\theta) = \frac{1}{2n} \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2$$

$$f(\theta) = \frac{1}{2n} \sum_{i=1}^n (\theta^T \mathbf{x}_i - y_i)^2$$

- Głównym pomysłem jest przyjęcie pochodnej cząstkowej funkcji kosztu względem θ .
- Ten gradient pomnożony przez współczynnik uczenia się staje się regułą aktualizacji szacowanych wartości parametrów.
- Dla każdego θ_j oblicza się:

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} f(\theta)$$

- co po przekształceniach daje:

$$\theta_j = \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n (\theta^T \mathbf{x}_i - y_i) * x_i^{(j)}$$

Section 2

Regresja grzbietowa

- Kiedy w zbiorze danych istnieje wiele skorelowanych zmiennych, współczynniki θ w modelach regresji liniowej, mogą być słabo określone i mieć wiele wariancji.
- Jednym z rozwiązań tego problemu jest ograniczenie parametrów θ tak, aby nie przekraczały pewnego budżetu C , inaczej mówiąc aby parametry nie miały zbyt dużych wartości.
- Jest to równoważne zastosowaniu regularyzacji $L2$, znanej również jako “rozpad wag”.

Subsection 1

Estymacja analityczna - forma pierwotna

- Funkcja celu przyjmuje teraz postać:

$$\hat{\theta} = \min_{\theta} \|X\theta - \mathbf{y}\|^2 + \lambda \|\theta\|^2$$

- gdzie λ jest parametrem regularyzacji.
- Po odpowiednich przekształceniach otrzymuje się:

$$\hat{\theta} = (X^T X + \lambda I)^{-1} X^T \mathbf{y}$$

- co jest formą bardzo podobną do zwykłej regresji liniowej, ale tutaj dodaje się λ do każdego elementu diagonalnego $X^T X$.

- Powyższy wzór można zapisać jako:

$$\begin{aligned}(X^T X + \lambda I) \hat{\theta} &= X^T \mathbf{y} = \\ X^T X \hat{\theta} + \lambda \hat{\theta} &= X^T \mathbf{y} = \\ \lambda \hat{\theta} &= X^T \mathbf{y} - X^T X \hat{\theta} = X^T (\mathbf{y} - X \hat{\theta}) = \\ \hat{\theta} &= \lambda^{-1} X^T (\mathbf{y} - X \hat{\theta}) =\end{aligned}$$

$$\hat{\theta} = X^T \alpha$$

- gdzie:

$$\alpha = \lambda^{-1} (\mathbf{y} - X \hat{\theta})$$

- Tak samo jak w przypadku regresji liniowej powyższą estymację analityczną można uzyskać także za pomocą macierzy Gramma:

$$\lambda \alpha = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\theta}}$$

- posługując się podstawieniem ze slajdu nr 8 otrzymujemy (macierz Gramma):

$$\lambda \alpha = \mathbf{y} - \mathbf{X}\mathbf{X}^T \alpha$$

$$\mathbf{X}\mathbf{X}^T \alpha + \lambda \alpha = \mathbf{y}$$

$$(\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I}) \alpha = \mathbf{y}$$

$$\alpha = (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$$

Subsection 2

Estymacja analityczna - wersja jądrowa

- Jądra są używane do obliczania iloczynu skalarnego dwóch wektorów w pewnej przestrzeni cech, nawet bez odwiedzania jej.
- Teraz zastępujemy wszystkie obserwacje, wektorem ich cech w nowej przestrzeni (mogącej mieć mniej lub więcej wymiarów niż oryginalna przestrzeń):

$$\phi : X \rightarrow H$$

- Możemy postrzegać jądro jako:

$$k(\mathbf{x}_i, \mathbf{x}_n) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_n) \rangle$$

$$K(X, X) = k(\mathbf{x}_i, \mathbf{x}_n), \forall i, n$$

- chociaż nie wiemy, czym jest $\phi(.)$ wiemy tylko, że istnieje.

- Jądro liniowe:

$$k(\mathbf{x}_i, \mathbf{x}_n) = \mathbf{x}_i^T \mathbf{x}_n$$

- Jądro wielomianowe stopnia d i stałej c :

$$k(\mathbf{x}_i, \mathbf{x}_n) = (\mathbf{x}_i^T \mathbf{x}_n + c)^d$$

- Jądro o radialnych funkcjach bazowych:

$$k(\mathbf{x}_i, \mathbf{x}_n) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_n\|^2}{2\sigma^2}\right)$$

- Teraz możemy po prostu wziąć rozwiązanie dla regresji grzbietowej i zamienić każde X na $\Phi(X)$:

$$\hat{\theta} = \Phi(X)^T (\Phi(X)\Phi(X)^T + \lambda I)^{-1} \mathbf{y}$$

- Macierz Gramma (będąca iloczynem skalarnym) może być zastąpiona dowolną macierzą kwadratową, będącą przekształceniem nieliniowym do nowej przestrzeni:

$$\alpha = (XX^T + \lambda)^{-1} \mathbf{y} = (K(X, X) + \lambda I)^{-1} \mathbf{y}$$

$$\hat{y}_i = \phi(\mathbf{x}_i^T) \Phi(X)^T \alpha = \sum_{j=1}^n k(\mathbf{x}_i, \mathbf{x}_j) \alpha_j$$

- Dlaczego mielibyśmy chcieć zachować macierz K rozmiaru $n \times n$ zamiast pierwotnej macierzy X rozmiaru $n \times p$ lub macierzy $X^T X$ rozmiaru $p \times p$?
- Nie zrobilibyśmy tego, gdybyśmy dokonali rzeczywistej regresji liniowej lub grzbietowej.
- Możemy teraz jednak zastąpić macierz K z dowolną funkcją jądra na regresję grzbietową w innej przestrzeni, uzyskując regresję nieliniową (i nieparametryczną)!

Section 3

Literatura

- *Friedman, J., Hastie, T., Tibshirani, R. (2001). The elements of statistical learning (Vol. 1, No. 10). New York: Springer series in statistics.*