

1 Regresja liniowa cz. II

1.1 Miary dopasowania modelu do danych

Dekompozycja wariancji

Idea:

$$\underbrace{y - \bar{y}}_{\text{Odchylenie całkowite}} = \underbrace{y - \hat{y}}_{\text{Odchylenie nie wyjaśnione (błąd)}} + \underbrace{\hat{y} - \bar{y}}_{\text{Odchylenie wyjaśnione (regresyjne)}}$$

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{Całkowita suma kwadratów SST}} = \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\text{Suma kwadratów błędów SSE}} + \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\text{Suma kwadratów odchyleń regresyjnych SSR}}$$

SSR jest tzw. zmienność wyjaśniona

SSE jest tzw. zmienność niewyjaśniona (wynika z niedoskonałości regresji)

Przykład 1 RYSUNEK!!!

Definicja 1 Współczynnik determinacji r^2 [w programach czasem R^2]

$$r^2 = \frac{SSR}{SST}$$

Wniosek 1

$$r^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

Wartość współczynnika determinacji wyraża zmienność Y , która została wyjaśniona przez zachodzenie liniowego związku między X i Y . Wartość współczynnika korelacji liniowej to miara liniowego związku.

Lemat 1

$$r^2 = \hat{\rho}_{X,Y}^2$$

Własności 1 1. Współczynnik determinacji r^2 jest miarą siły liniowego związku między zmiennymi.

2. $r^2 \in [0, 1]$

3. Wartość współczynnika determinacji nic nam nie mówi o kierunku zależności między zmiennymi.

4. Większe wartości r^2 mogą sugerować dokładniejsze dopasowanie modelu regresji do zbioru obserwacji.

5. $r^2 = 1$ oznacza, że zmienność Y jest w 100% wyjaśniona przez model regresji

6. $r^2 \leq 0.3$ oznacza, że regresja wyjaśnia mniej niż 30% zmienności Y .

Uwaga 1 Sama wartość współczynnika determinacji nie pozwala ocenić jakości dopasowania linii regresji do zbioru obserwacji. Istnieją przykłady, dla których r^2 jest duże, a model jest nieodpowiedni.

Lemat 2 1. $SST = SS_Y$

2. $SSR = \hat{\beta}_1 SS_{XY}$

3. $SSE = SS_Y - \hat{\beta}_1 SS_{XY}$

Ćwiczenie 1 Wskaż przykład danych gdzie model regresji jest nieodpowiedni, a r^2 jest duże. Wsk. parabola.

Definicja 2

$$MSR = \frac{SSR}{1}$$
$$MSE = \frac{SSE}{n-2}$$

MSR to średnie kwadratowe odchylenie regresyjne (mean square regression).
 MSE to średni błąd kwadratowy (mean square error)

Uwaga 2 MSE mierzy/wyraża rozproszenie danych od linii regresji. Jest to miara bezwzględna - jej wielkość w zależności od sytuacji może mieć inne znaczenie.

Uwaga 3 MSE wcześniej oznaczano $\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{SSE}{n-2}$ - bo to również estymator wariancji reszt.

Lemat 3 Jeśli $\beta_1 = 0$ to

$$F = \frac{MSR}{MSE} \sim F_{1, n-2}$$

Obserwacja 1 Stosunek zmienności wyjaśnionej do zmienności niewyjaśnionej $\frac{MSR}{MSE}$ może być potraktowany jako kolejna miara dopasowania modelu do danych. Przy braku liniowego związku między zmiennymi wielkość MSR powinna być bliska zeru. Fakt ten można wykorzystać do przeprowadzenia testu statystycznego badającego zachodzenie liniowego związku.

H_0 : $\beta_1 = 0$ (nie zachodzi związek między X i Y)

H_1 : $\beta_1 \neq 0$

$$\frac{MSR}{MSE} \sim F_{1, n-2}$$

Jest to test jednostronny (prawostronny).

Obserwacja 2

$$F = t^2$$

gdzie $t = \frac{\hat{\beta}_1}{SE(\hat{\beta}_1)}$.

Wniosek 2 Wnioskowanie za pomocą testu F i testu t jest równoważne. F jest testem jednostronnym, a t dwustronnym. W obu testach otrzymamy to samo p -value, gdy dla t ustalimy symetryczne przedziały - zbiory krytyczne.

Prognoza

Jednym z zastosowań modelu regresji jest prognozowanie dalszego poziomu zjawiska.

Prognoza punktowa:

Niech x_0 jest ustaloną wartością X . Odpowiadającą jej wartość (teoretyczną) zmiennej Y obliczymy ze wzoru

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$$

Prawdziwa wartość wynosi

$$y_0 = \beta_0 + \beta_1 x_0 + \varepsilon_0$$

Błąd prognozy

$$\hat{y}_0 - y_0 = \hat{\beta}_0 - \beta_0 + (\hat{\beta}_1 - \beta_1) x_0 - \varepsilon_0$$

$$E(\hat{y}_0 - y_0) = 0$$

$$Var(\hat{y}_0 - y_0) = \sigma^2 \left[1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]$$

Wniosek 3 Niepewność mierzona odchyleniem standardowym błędu prognozy rośnie wraz ze zwiększaniem się odległości x_0 od \bar{x} , oraz maleje wraz ze zwiększaniem się liczebności próbki.

Obserwacja 3 Jeśli $x_0 = \bar{x}$ to $\hat{y}_0 = \bar{y}$ oraz

$$Var(\hat{y}_0 - y_0) = Var(\bar{y} - y_0) = \sigma^2 \left[1 + \frac{1}{n} \right] \xrightarrow{n \rightarrow \infty} \sigma^2$$

Uwaga 4 W praktyce nie znamy σ^2 zastępujemy go estymatorem $\hat{\sigma}^2$. $\sqrt{\hat{\sigma}^2 \left[1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]}$ nazywamy standardowym błędem predykcji.

Prognoza przedziałowa:

Twierdzenie 1 $(1 - \alpha)$ 100% przedział predykcji zmiennej Y dla ustalonej wartości x_0 zmiennej objaśniającej jest określony wzorem

$$\hat{y}_0 \pm t_{n-2, \alpha/2} \sqrt{\hat{\sigma}^2 \left[1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_X} \right]}$$

gdzie $P(t > t_{n-2, \alpha/2}) = \alpha/2$, $t \sim t_{n-2}$.

Prognoza przeciętnego poziomu zmiennej Y

Motywacja:

y_0 - prognoza przychodów ze sprzedaży

$E(y_0)$ - prognoza przeciętnych przychodów ze sprzedaży

Mając ustaloną wartość x_0 jesteśmy zainteresowani prognozą $E(y_0)$ a nie y_0 .

Zauważmy, że $E(y_0) = \beta_0 + \beta_1 x_0$, wielkość tę przybliżamy wartością $\hat{E}(y_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0$. **Wartość prognozy dla $E(y_0)$ i y_0 są takie same** ($\widehat{E}(y_0) = \hat{y}_0$)
! Zobaczmy czy zmieni się standardowy błąd predykcji:

$$E(y_0) - \hat{E}(y_0) = (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1) x_0$$

$$\text{Var}(E(y_0) - \hat{E}(y_0)) = \sigma^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]$$

Twierdzenie 2 $(1 - \alpha)$ 100% przedział predykcji średniego poziomu zmiennej Y dla ustalonej wartości x_0 zmiennej objaśniającej jest określony wzorem

$$\hat{E}(y_0) \pm t_{n-2, \alpha/2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_X} \right]}$$

gdzie $P(t > t_{n-2, \alpha/2}) = \alpha/2$, $t \sim t_{n-2}$.

Uwaga 5 Zgodnie z oczekiwaniami przedział ufności dla przeciętnego poziomu jest węższy niż dla samego poziomu.