

Föreläsning 8: Linjär regression

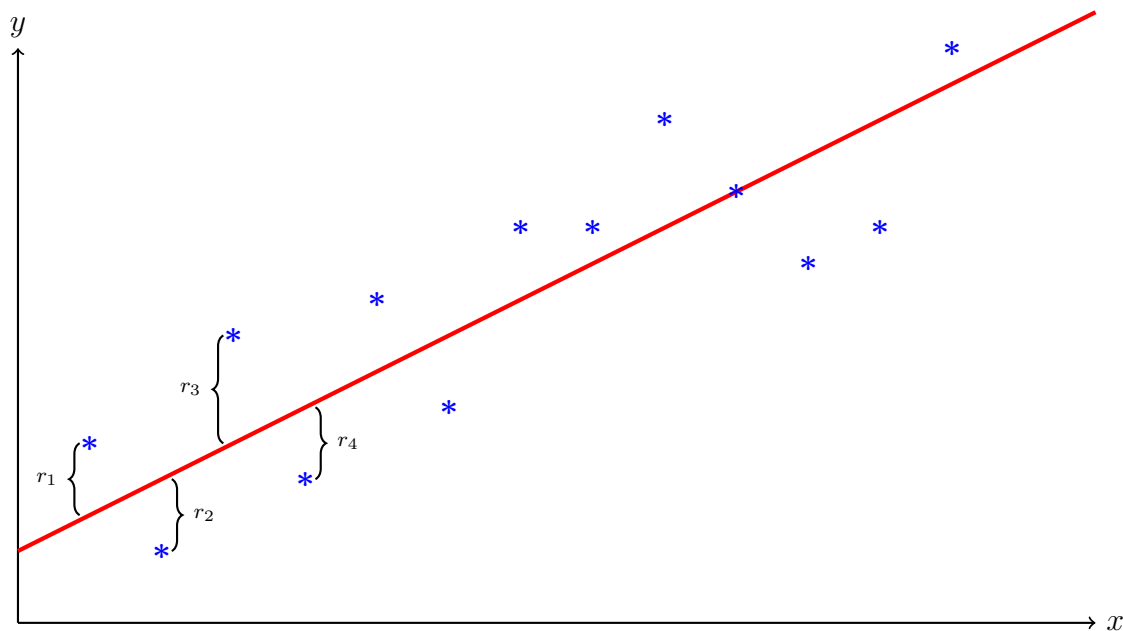
Johan Thim (johan.thim@liu.se)

12 mars 2020

"Your suffering will be legendary, even in hell."

–Pinhead

Vi återgår nu till ett exempel vi stött på redan vid ett flertal tillfällen (föregående föreläsning och föreläsning 2 för att inte tala om tidigare kurser som linjär algebra), nämligen att anpassa en rät linje $y = \beta_0 + \beta_1 x$ efter mätdata (x_i, y_i) , $i = 1, 2, \dots, n$. Grafiskt illustrerat enligt nedan.



Målsättningen är att – givet en mätserie – hitta den linje som approximerar denna serie på lämpligt sätt. Det resulterar i ett par naturliga funderingar.

- (i) Hur hittar man en approximativ linje systematiskt?
- (ii) Om man upprepar försöket, får man samma linje?
- (iii) I vilken mening är linjen optimal? På vilket sätt mäter vi avvikelserna mellan linjen och mätserien?

Vi ska försöka svara på dessa frågor och för att göra det behöver vi ställa upp en modell.



Enkel linjär regression

Definition. Vi kommer betrakta följande modell: givet (x_j, y_j) , $j = 1, 2, \dots, n$, där vi betraktar x_j som fixerade och y_j som observationer av stokastiska variabler

$$Y_j = \beta_0 + \beta_1 x_j + \epsilon_j, \quad j = 1, 2, \dots, n,$$

där $\epsilon_j \sim N(0, \sigma^2)$ antas oberoende (och likafördelade). Den räta linjen $y = \beta_0 + \beta_1 x$ kallas **regressionslinjen**.

Vi använder beteckningen $\mu_j = \beta_0 + \beta_1 x_j = E(Y_j)$.

Är det givet att denna modell är sann? Nej, det är inte självklart utan hänger på vilka förutsättningar datan kommer från. Däremot tenderar modellen att fungera bra i de flesta fall om vi har en rimlig mängd observationer.

Med notationen från figuren ser vi att

$$r_j = y_j - \mu_j$$

om vi tar hänsyn till tecknet (positivt tecken om y_j ligger ovanför regressionslinjen). Ett mått på hur väl linjen approximerar mätserien ges av kvadratsumman

$$\sum_{j=1}^n r_j^2 = \sum_{j=1}^n (y_j - \mu_j)^2.$$

Vi skulle kunna minimera denna summa med avseende på (β_0, β_1) , vilket ger MK-skattningar för β_0 och β_1 . Detta var det som hände i exemplet från föreläsning 2. Istället för att upprepa argumentet går vi över till den generella modellen för linjär regression.



Linjär regression

Definition. Givet $(x_{j1}, x_{j2}, \dots, x_{jk}, y_j)$, $j = 1, 2, \dots, n$, för något positivt heltal k , där vi betraktar x_{ji} som fixerade och y_j som observationer av stokastiska variabler

$$Y_j = \beta_0 + \beta_1 x_{j1} + \beta_2 x_{j2} + \dots + \beta_k x_{jk} + \epsilon_j, \quad j = 1, 2, \dots, n,$$

där $\epsilon_j \sim N(0, \sigma^2)$ antas oberoende (och likafördelade) och $\beta_0, \beta_1, \dots, \beta_k$ är okända parametrar. Den räta linjen

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

kallas **regressionslinjen**.

Vi låter

$$\mu_j = E(Y_j) = \beta_0 + \beta_1 x_{j1} + \beta_2 x_{j2} + \dots + \beta_k x_{jk}.$$

Dubbelindexeringen är lite jobbig att arbeta med, så låt oss gå över till matrisnotation:

$$\begin{aligned} \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} &= \begin{pmatrix} \beta_0 & \beta_1 x_{11} & \cdots & \beta_k x_{1k} \\ \beta_0 & \beta_1 x_{21} & \cdots & \beta_k x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_0 & \beta_1 x_{n1} & \cdots & \beta_k x_{nk} \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix} \\ &= \begin{pmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix} \end{aligned}$$

Vi låter

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \quad X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}, \quad \text{samt } \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix},$$

vilket leder till sambandet

$$\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

Således gäller att

$$E(\mathbf{Y}) = X\boldsymbol{\beta} \quad \text{och} \quad C_{\mathbf{Y}} = \sigma^2 I_n,$$

där I_n är den n -dimensionella enhetsmatrisen. Vi söker MK-skattningen $\hat{\boldsymbol{\beta}}$ för $\boldsymbol{\beta}$, vilket vi kan erhålla genom att minimera den kvadratiske formen

$$Q(\beta_0, \dots, \beta_k) = \sum_{j=1}^n (y_j - \mu_j)^2 = (\mathbf{y} - X\boldsymbol{\beta})^T (\mathbf{y} - X\boldsymbol{\beta}),$$

där $\mathbf{y} = (y_1 \ y_2 \ \cdots \ y_n)^T$. En variant är att sätta igång och derivera, men lite mer elegant gäller följande sats.



Normalekvationerna

Sats. Om $\det X^T X \neq 0$ så ges MK-skattningen av $\boldsymbol{\beta}$ enligt

$$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y}.$$

Bevis. Vi kan skriva vektorn \mathbf{Y} av uppmätta värden som

$$\mathbf{Y} = \beta_0 \mathbf{x}_0 + \beta_1 \mathbf{x}_1 + \cdots + \beta_k \mathbf{x}_k + \boldsymbol{\epsilon},$$

där \mathbf{x}_i , $i = 1, 2, \dots, k$, är kolonn $i + 1$ i matrisen X . Från linjär algebra vet vi att avståndet mellan observationen \mathbf{y} och linjärkombinationer av vektorerna \mathbf{x}_j , dvs

$$|\mathbf{y} - (\hat{\beta}_0 \mathbf{x}_0 + \hat{\beta}_1 \mathbf{x}_1 + \cdots + \hat{\beta}_k \mathbf{x}_k)|,$$

blir minimalt precis då

$$X\hat{\boldsymbol{\beta}} = \sum_{j=1}^k \hat{\beta}_j \mathbf{x}_j$$

är projektionen av \mathbf{y} på det linjära höljet $\text{span}\{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k\}$. Således måste $\mathbf{y} - X\hat{\boldsymbol{\beta}}$ vara vinkelrät mot \mathbf{x}_j för alla $j = 0, 1, \dots, k$. Detta medför att

$$X^T (\mathbf{y} - X\hat{\boldsymbol{\beta}}) = \mathbf{0} \quad \Leftrightarrow \quad X^T X\hat{\boldsymbol{\beta}} = X^T \mathbf{y} \quad \Leftrightarrow \quad \hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y},$$

vilket är precis vad vi ville visa! □



När det gäller $\hat{\beta}$ och dess komponenter kommer vi använda samma beteckningar för observationer av punktskattningen och den stokastiska variabeln. Skulle vi vara konsekventa borde den stokastiska variabeln betecknas $\hat{\mathbf{B}}$, men av tradition görs inte så. Var observant!



Sats. Med förutsättningarna ovan gäller att

$$\hat{\beta} \sim N_{k+1}(\beta, \sigma^2(X^T X)^{-1}).$$

Bevis. Det faktum att $\hat{\beta}$ är normalfördelad följer direkt från faktumet att elementen i $\hat{\beta}$ är linjärkombinationer av normalfördelade variabler. Återstår att visa väntevärde och kovarians:

$$E(\hat{\beta}) = (X^T X)^{-1} X^T E(\mathbf{Y}) = (X^T X)^{-1} X^T X \beta = \beta$$

och

$$\begin{aligned} C_{\hat{\beta}} &= (X^T X)^{-1} X^T C_{\mathbf{Y}} ((X^T X)^{-1} X^T)^T = (X^T X)^{-1} X^T \sigma^2 I_n ((X^T X)^{-1} X^T)^T \\ &= (X^T X)^{-1} X^T \sigma^2 I_n (X^T)^T ((X^T X)^{-1})^T = \sigma^2 (X^T X)^{-1} X^T X ((X^T X)^T)^{-1} = \sigma^2 (X^T X)^{-1}. \end{aligned}$$

Således blir fördelningen precis som beskriven i satsen. □

1 Variansanalys

Vi låter

$$\hat{\mu}_j = \hat{\beta}_0 + \hat{\beta}_1 x_{j1} + \hat{\beta}_2 x_{j2} + \dots + \hat{\beta}_k x_{jk}, \quad j = 1, 2, \dots, n.$$

Ibland betecknas vektorn $\hat{\mu}$ som $\hat{\mathbf{y}}$ eftersom det i någon mening är en skattningen av \mathbf{y} , men det blir lite olyckligt för det är inte \mathbf{y} vi skattar. Vi ska nu studera hur bra den skattning vi tagit fram är.



Regressionsanalysens kvadratsummor

Definition. Vi definierar tre kvadratsummor som beskriver variationen hos y -värdena:

(i) Den **totala variationen** definieras enligt $SS_{\text{TOT}} = \sum_{j=1}^n (y_j - \bar{y})^2$.

(ii) Variationen som förklaras av x_1, x_2, \dots, x_k , definieras enligt $SS_{\text{R}} = \sum_{j=1}^n (\hat{\mu}_j - \bar{y})^2$.

(iii) Variationen som inte förklaras av regressionsmodellen är $SS_{\text{E}} = \sum_{j=1}^n (y_j - \hat{\mu}_j)^2$.

Ibland används beteckningarna Q_{TOT} för SS_{TOT} , Q_{REGR} för SS_{R} och Q_{RES} för SS_{E} . Hur hänger dessa summor ihop? Som tur är finns ett enkelt svar.



Sats. Den totala variationen SS_{TOT} kan delas upp enligt

$$SS_{\text{TOT}} = SS_{\text{R}} + SS_{\text{E}}.$$

Bevis. Vi vill uttrycka SS_{TOT} i termer av SS_{R} och SS_{E} , så

$$\begin{aligned} SS_{\text{TOT}} &= \sum_{j=1}^n (y_j - \bar{y})^2 = \sum_{j=1}^n (y_j - \hat{\mu}_j + \hat{\mu}_j - \bar{y})^2 \\ &= \sum_{j=1}^n (y_j - \hat{\mu}_j)^2 + 2 \sum_{j=1}^n (y_j - \hat{\mu}_j)(\hat{\mu}_j - \bar{y}) + \sum_{j=1}^n (\hat{\mu}_j - \bar{y})^2 \\ &= SS_{\text{R}} + 2 \sum_{j=1}^n (y_j - \hat{\mu}_j)\hat{\mu}_j - 2\bar{y} \sum_{j=1}^n (y_j - \hat{\mu}_j) + SS_{\text{E}}. \end{aligned}$$

Vi visar att de båda summorna i mitten summerar till noll. Eftersom $\hat{\boldsymbol{\mu}} = X\hat{\boldsymbol{\beta}}$ gäller det att

$$\begin{aligned} \sum_{j=1}^n (y_j - \hat{\mu}_j)\hat{\mu}_j &= \hat{\boldsymbol{\mu}}^T(\mathbf{y} - \hat{\boldsymbol{\mu}}) = (X\hat{\boldsymbol{\beta}})^T(\mathbf{y} - X\hat{\boldsymbol{\beta}}) = \hat{\boldsymbol{\beta}}^T X^T(\mathbf{y} - X\hat{\boldsymbol{\beta}}) \\ &= \hat{\boldsymbol{\beta}}^T(X^T\mathbf{y} - X^T X(X^T X)^{-1} X^T\mathbf{y}) = \hat{\boldsymbol{\beta}}^T(X^T\mathbf{y} - X^T\mathbf{y}) = 0. \end{aligned}$$

Med andra ord är $\mathbf{y} - \hat{\boldsymbol{\mu}}$ vinkelrät mot $\hat{\boldsymbol{\mu}}$.

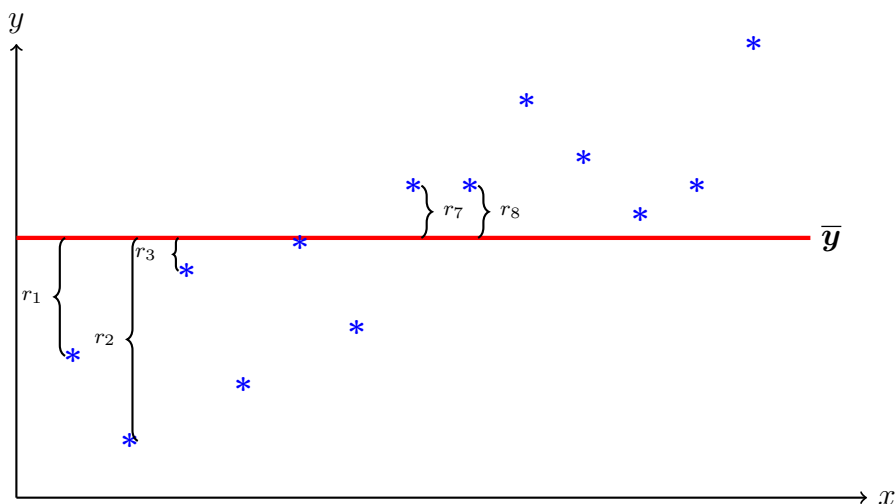
Vidare vet vi att $\hat{\boldsymbol{\beta}}$ minimerar $Q(\boldsymbol{\beta}) = \sum_{j=1}^n (y_j - \mu_j)^2$, så

$$0 = \left. \frac{\partial}{\partial \beta_0} Q(\boldsymbol{\beta}) \right|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} = -2 \sum_{j=1}^n (y_j - \hat{\mu}_j) \Leftrightarrow \sum_{j=1}^n y_j = \sum_{j=1}^n \hat{\mu}_j,$$

vilket var precis det vi behövde. □

1.1 SS_{TOT}

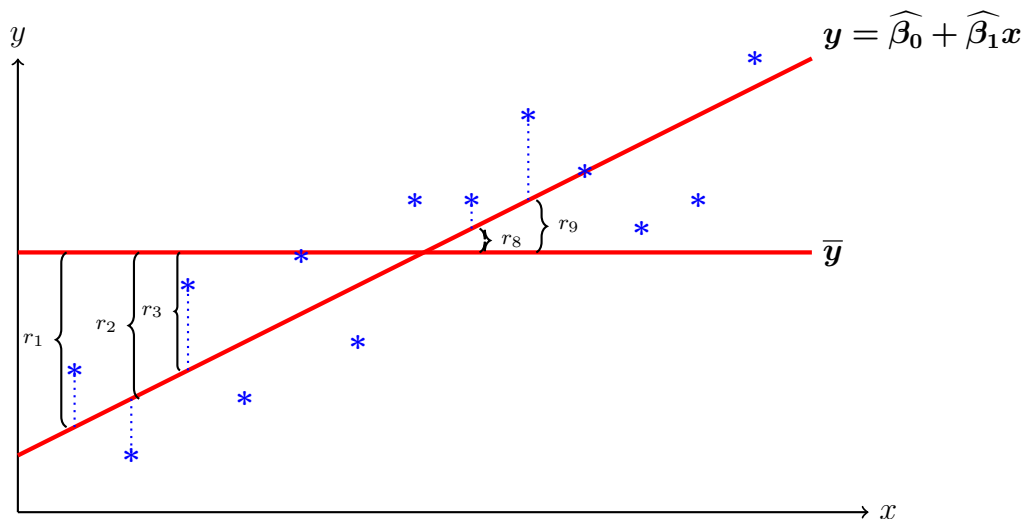
Storheten SS_{TOT} mäter den totala variation av mätvärden jämfört med mätvärdenas medelvärde. I fallet då vi använder modellen $Y = \beta_0 + \epsilon$ ger således SS_{TOT} hela felet i regressionen.



Vi ser att $SS_{TOT} = \sum_{j=1}^n r_j^2 = \sum_{j=1}^n (y_j - \bar{y})^2$.

1.2 SS_R och SS_E

Om vi istället använder modellen $y = \beta_0 + \beta_1 x + \epsilon$ blir summorna SS_E och SS_R relevanta (SS_{TOT} ser fortfarande ut som ovan). Låt oss först illustrera SS_R .

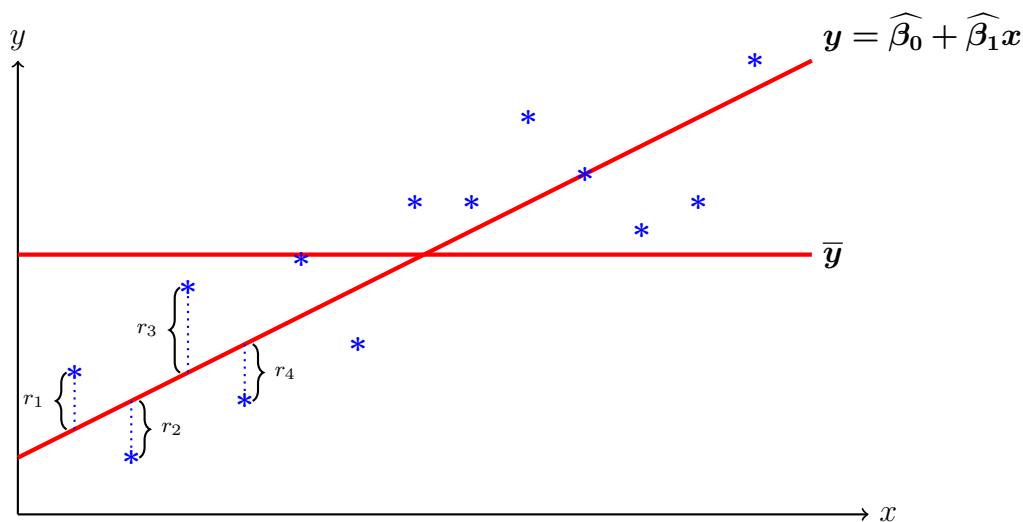


Vi ser att

$$SS_R = \sum_{j=1}^n r_j^2 = \sum_{j=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_j - \bar{y})^2$$

och SS_R mäter alltså hur mycket den skattade regressionslinjen (i mätpunkterna x_j) skiljer sig från medelvärdet av y -värdena.

När det gäller SS_E är det istället skillnaden mellan den skattade regressionslinjen och mätvärdena vi betraktar.



Kvadratsumman av dessa avvikelser blir

$$SS_E = \sum_{j=1}^n r_j^2 = \sum_{j=1}^n (y_j - (\hat{\beta}_0 + \hat{\beta}_1 x_j))^2.$$

2 Projektionsmatriser

Eftersom vi arbetar med matriser och MK-lösningen i princip är projektionen på ett underrum av \mathbf{R}^n så är följande resultat inte allt för förvånande.



Hatt-matrisen

Definition. Vi definierar $H = X(X^T X)^{-1} X^T$. Vi definierar även J som en matris vars samtliga element är 1. Vi skriver J_{nm} om vi vill markera dimensionen.

Varför kallar vi H för hatt-matrisen? Ganska enkelt:

$$H\mathbf{y} = X(X^T X)^{-1} X^T \mathbf{y} = X\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\mu}} = \hat{\mathbf{y}}.$$

Avbildningen sätter alltså hatten på!



Sats. Matriserna H , $I - H$ och $H - \frac{1}{n}J$ är projektionsmatriser P som uppfyller $P^2 = P^T = P$.

Bevis. Direkt från definitionen av H blir

$$H^2 = (X(X^T X)^{-1} X^T)(X(X^T X)^{-1} X^T) = X(X^T X)^{-1} X^T = H$$

och

$$H^T = (X(X^T X)^{-1} X^T)^T = X(X^T X)^{-T} X^T = X((X^T X)^T)^{-1} X^T = H.$$

Därför följer det att $(I - H)^2 = I^2 - 2H + H^2 = I - H$ och att $(I - H)^T = I^T - H^T = I - H$. För den sista operatören skriver vi

$$\left(H - \frac{1}{n}J\right)^2 = H^2 - \frac{1}{n}HJ - \frac{1}{n}JH + \frac{1}{n^2}J^2 = H - \frac{1}{n}HJ - \frac{1}{n}JH + \frac{1}{n}J$$

ty J^2 är matrisen med samtliga element lika med n . Vidare gäller att J kommuterar med alla kvadratiska matriser, så $JH = HJ$. Här kan vi se att $H\mathbf{1} = \mathbf{1}$ (där $\mathbf{1}$ är en vektor med samtliga element lika med 1) eftersom lösningen till regressionsproblemet om \mathbf{y} är konstant är just den konstanten (lösningen är exakt). Alltså blir

$$\left(H - \frac{1}{n}J\right)^2 = H - \frac{1}{n}J.$$

Givetvis gäller även att $\left(H - \frac{1}{n}J\right)^T = H - \frac{1}{n}J$. □

En följsats av detta är att vi kan skriva kvadratsummorna som kvadratiska former.



Sats. $SS_{\text{TOT}} = \mathbf{y}^T \left(I - \frac{1}{n}J\right) \mathbf{y}$, $SS_{\text{R}} = \mathbf{y}^T \left(H - \frac{1}{n}J\right) \mathbf{y}$, samt $SS_{\text{E}} = \mathbf{y}^T (I - H) \mathbf{y}$.

3 Förklaringsgrad

När vi utför regressionsanalys vill vi ofta ha ett mått som preciserar hur bra modellen passar de uppmätta värdena. Ett sådant mått är förklaringsgraden.



Definition. Förklaringsgraden R^2 definieras som

$$R^2 = \frac{SS_R}{SS_{TOT}} = \frac{SS_{TOT} - SS_E}{SS_{TOT}} = 1 - \frac{SS_E}{SS_{TOT}}.$$

Ibland använder man lite andra varianter av R^2 och en mycket vanlig är den så kallade **justerade förklaringsgraden**

$$R_{adj}^2 = 1 - \frac{SS_E/(n - k - 1)}{SS_{TOT}/(n - 1)}.$$

Om inte n är förhållandevis stor i jämförelse med k (vilket är nödvändigt för att regressionen ska vara vettig) så blir den justerade graden sämre.

4 Regressionsanalysens huvudsats

Innan vi ger oss in på huvudsatsen visar vi en hjälpsats från linjär algebra.



Sats. Låt $\mathbf{x} = (x_1 \ x_2 \ \cdots \ x_n)$ vara sådan att

$$\mathbf{x}^T \mathbf{x} = \sum_{j=1}^n x_j^2 = \mathbf{x}^T A \mathbf{x} + \mathbf{x}^T B \mathbf{x}$$

för $A, B \in \mathbf{R}^{n \times n}$ positivt semi-definita och symmetriska med $\text{rank}(A) = r$ och $\text{rank}(B) = n - r$ för något heltal r så att $0 < r < n$. Då finns en ON-matris C så att med $\mathbf{x} = C\mathbf{y}$ gäller att

$$\mathbf{x}^T A \mathbf{x} = \sum_{j=1}^r y_j^2 \quad \text{och} \quad \mathbf{x}^T B \mathbf{x} = \sum_{j=r+1}^n y_j^2.$$

Bevis. Eftersom A är positivt semi-definit har A endast icke-negativa egenvärden som vi ordnar enligt $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$. Vi vet även att $\text{rank}(A) = r$, så $\lambda_{r+1} = \cdots = \lambda_n = 0$. Vidare är A diagonaliserbar eftersom A är symmetrisk, så det finns en ON-matris C så att

$$C^T A C = D = \begin{pmatrix} \lambda_1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 & 0 & 0 & 0 \\ 0 & \vdots & \ddots & \vdots & 0 & 0 & 0 \\ 0 & 0 & 0 & \lambda_r & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

Med $\mathbf{x} = C\mathbf{y}$ gäller då att

$$\mathbf{x}^T \mathbf{x} = (C\mathbf{y})^T (C\mathbf{y}) = \mathbf{y}^T C^T C \mathbf{y} = \mathbf{y}^T \mathbf{y}$$

och

$$\mathbf{x}^T A \mathbf{x} = (C\mathbf{y})^T A C \mathbf{y} = \mathbf{y}^T C^T A C \mathbf{y} = \mathbf{y}^T D \mathbf{y} = \sum_{j=1}^r \lambda_j y_j^2.$$

Vi har även

$$\mathbf{x}^T B \mathbf{x} = (C\mathbf{y})^T B C \mathbf{y} = \mathbf{y}^T C^T B C \mathbf{y}$$

och då blir

$$\sum_{j=1}^n y_j^2 = \mathbf{y}^T \mathbf{y} = \mathbf{x}^T \mathbf{x} = \mathbf{x}^T A \mathbf{x} + \mathbf{x}^T B \mathbf{x} = \sum_{j=1}^r \lambda_j y_j^2 + \mathbf{y}^T C^T B C \mathbf{y}$$

så

$$\mathbf{y}^T C^T B C \mathbf{y} = \sum_{j=1}^r (1 - \lambda_j) y_j^2 + \sum_{j=r+1}^n y_j^2.$$

Det faktum att $\text{rank}(B) = n - r$ visar att $\lambda_1 = \lambda_2 = \dots = \lambda_r = 1$ och vi erhåller identiteten

$$\mathbf{y}^T C^T B C \mathbf{y} = \sum_{j=r+1}^n y_j^2.$$

Alla beteckningar är nu ur vägen och vi är framme vid huvudresultatet.



Regressionsanalysens huvudsats

Sats. Med förutsättningarna ovan gäller följande för SS_E och SS_R sedda som stokastiska variabler.

(i) $\frac{SS_E}{\sigma^2} = \frac{1}{\sigma^2} \sum_{j=1}^n (Y_j - \hat{\mu}_j)^2 \sim \chi^2(n - k - 1).$

(ii) Givet att $\beta_1 = \beta_2 = \dots = \beta_k = 0$ är $\frac{SS_R}{\sigma^2} = \frac{1}{\sigma^2} \sum_{j=1}^n (\hat{\mu}_j - \bar{Y})^2 \sim \chi^2(k).$

(iii) Både SS_R och $\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{Y}$ är oberoende av SS_E .

Bevis. Eftersom H är en projektionsmatris har H endast egenvärdena $\lambda = 0$ och $\lambda = 1$. Således gäller det finurliga att matrisens rang är lika med summan av egenvärdena, vilka kan beräknas genom att ta spåret¹ av matrisen:

$$\text{rank}(H) = \text{tr}(H) = \text{tr}(X(X^T X)^{-1} X^T) = \text{tr}(X^T X (X^T X)^{-1}) = \text{tr}(I_{k+1}) = k + 1,$$

Det följer sedan att

$$\text{rank}(I - H) = n - k - 1.$$

¹se sista (bonus)avsnittet för lite detaljer kring spår av matriser.

Eftersom $HX = X$ så gäller att

$$\mathbf{Y}^T(I - H)\mathbf{Y} = \boldsymbol{\epsilon}^T(I - H)\boldsymbol{\epsilon}.$$

Detta är användbart eftersom $E(\boldsymbol{\epsilon}) = \mathbf{0}$. Enligt föregående hjälpsats gäller att sambandet

$$\boldsymbol{\epsilon}^T\boldsymbol{\epsilon} = \boldsymbol{\epsilon}^T(I - H)\boldsymbol{\epsilon} + \boldsymbol{\epsilon}^T H\boldsymbol{\epsilon}$$

medför att det finns en ON-matris C så att $\mathbf{Z} = C\boldsymbol{\epsilon}$ reducerar likheten till

$$\boldsymbol{\epsilon}^T\boldsymbol{\epsilon} = \sum_{j=1}^{n-k-1} Z_j^2 + \sum_{j=n-k}^n Z_j^2$$

där

$$\boldsymbol{\epsilon}^T(I - H)\boldsymbol{\epsilon} = \mathbf{Y}^T(I - H)\mathbf{Y} = \text{SS}_E = \sum_{j=1}^{n-k-1} Z_j^2.$$

Eftersom komponenterna i $\boldsymbol{\epsilon}$ är oberoende och $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 I)$ följer det att

$$E(\mathbf{Z}) = C\mathbf{0} = \mathbf{0} \quad \text{och} \quad C_{\mathbf{Z}} = C_{C\boldsymbol{\epsilon}} = \sigma^2 C^T I C = \sigma^2 I$$

så $\mathbf{Z} \sim N(\mathbf{0}, \sigma^2 I)$. Komponenterna i \mathbf{Z} är alltså oberoende och $Z_j \sim N(0, \sigma^2)$. Detta medför att

$$\frac{\text{SS}_E}{\sigma^2} = \frac{1}{\sigma^2} \sum_{j=1}^{n-k-1} Z_j^2 \sim \chi^2(n - k - 1).$$

Vi erhåller även att SS_E och $H\boldsymbol{\epsilon}$ är oberoende (de delar inga variabler Z_j). Detta medför att SS_E och $\hat{\boldsymbol{\beta}}$ är oberoende. Det faktum att SS_E och SS_R är oberoende följer av att kovariansen

$$\begin{aligned} C \left(\left(H - \frac{1}{n} J \right) \mathbf{Y}, (I - H)\mathbf{Y} \right) &= \left(H - \frac{1}{n} J \right) C(\mathbf{Y}, \mathbf{Y})(I - H)^T = \sigma^2 \left(H - \frac{1}{n} J \right) (I - H) \\ &= \sigma^2 \left(H - H^2 - \frac{1}{n} J + \frac{1}{n} JH \right) = \sigma^2 \left(-\frac{1}{n} J + \frac{1}{n} HJ \right) \\ &= \sigma^2 \left(-\frac{1}{n} J + \frac{1}{n} J \right) = 0, \end{aligned}$$

så dessa vektorer är okorrelerade och normalfördelade, så oberoende. Det följer direkt att SS_E och SS_R är oberoende (som funktioner av oberoende variabler).

Om vi fokuserar på fördelningen för SS_R så kan vi på samma sätt som ovan utnyttja att $(H - \frac{1}{n} J)$ är en projektionsmatris, så

$$\text{rank} \left(H - \frac{1}{n} J \right) = \text{tr} \left(H - \frac{1}{n} J \right) = \text{tr}(H) - \text{tr} \left(\frac{1}{n} J \right) = k + 1 - 1 = k.$$

Matrisen J är synnerligen rangdefekt med $\text{rank}(J) = 1$ (eftersom alla kolonner är lika spänner vi bara upp ett en-dimensionellt rum).

Nu stöter vi på lite problem. Vi ser att

$$\left(H - \frac{1}{n} J \right) \mathbf{Y} = \left(I - \frac{1}{n} J \right) X\boldsymbol{\beta} + \left(H - \frac{1}{n} J \right) \boldsymbol{\epsilon}$$

där den första termen *inte* försvinner såvida inte $\beta_1 = \beta_2 = \dots = \beta_k = 0$. Men under detta antagande har vi åter igen en situation där vi kan ”byta ut” \mathbf{Y} mot $\boldsymbol{\epsilon}$. Föregående hjälpsats visar – analogt med föregående argument – att

$$\mathbf{Y}^T \left(H - \frac{1}{n} J \right) \mathbf{Y} = \sum_{j=1}^k Z_j^2,$$

där Z_j är oberoende och $Z_j \sim N(0, \sigma^2)$, vilket medför att $\frac{1}{\sigma^2} \text{SS}_R \sim \chi^2(k)$, under förutsättningen att $\beta_1 = \beta_2 = \dots = \beta_k = 0$. \square

Kommentar: utan villkoret $\beta_1 = \beta_2 = \dots = \beta_k = 0$ kan man fortfarande genomföra argumentet, men resultatet blir en icke-centrerad $\chi^2(k)$ -fördelning (något som inte ingår i kursen).

5 Hypotestester och kofidensintervall

Vi har nu samlat på oss en ordentlig verktygslåda, så det kanske är dags att se hur vi använder de olika delarna.

5.1 Skattning av σ^2

Variansen σ^2 skattar vi med

$$s^2 = \frac{\text{SS}_E}{n - k - 1}.$$

Denna skattning är väntevärdesriktig:

$$E(S^2) = \frac{\sigma^2}{n - k - 1} E \left(\frac{\text{SS}_E}{\sigma^2} \right) = \sigma^2 \frac{n - k - 1}{n - k - 1} = \sigma^2$$

ty $\frac{1}{\sigma^2} \text{SS}_E \sim \chi^2(n - k - 1)$.

5.2 Finns det något vettigt i modellen?

En rimlig fråga efter utförd regression är om modellen faktiskt säger något. Vi brukar testa denna hypotes enligt följande. Vi testar nollhypotesen

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$$

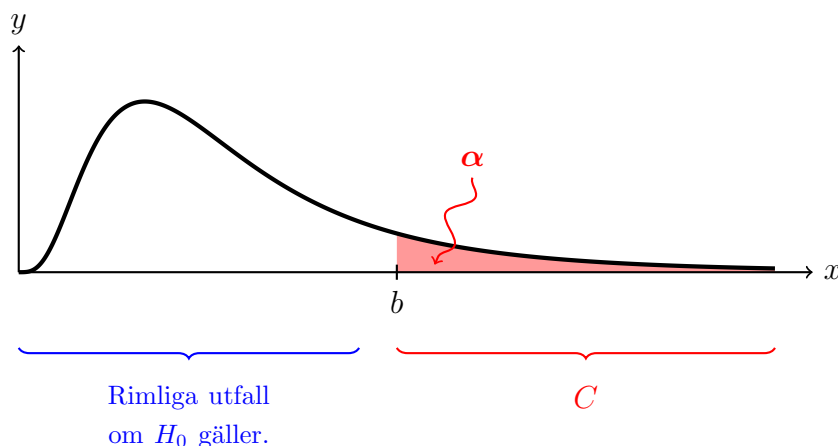
mot

$$H_1 : \text{minst något } \beta_i \neq 0.$$

En naturlig testvariabel ges av

$$v = \frac{\text{SS}_R/k}{\text{SS}_E/(n - k - 1)}.$$

Huvudsatsen medför att $V \sim F(k, n - k - 1)$ (om H_0 är sann) och vi förkastar H_0 om v är stor. Hur stor avgörs av signifikansnivån och lite slagning i F -tabeller.



Vi hittar gränsen b ur tabell (eller med `finv(1-alpha, k, n-k-1)` i MATLAB) så att

$$P(V > b) = \alpha.$$

5.3 Enskilda koefficienter

Eftersom $\hat{\beta} \sim N(\beta, \sigma^2(X^T X)^{-1})$ introducerar vi beteckningen

$$(X^T X)^{-1} = \begin{pmatrix} h_{00} & h_{01} & \cdots & h_{0k} \\ h_{10} & h_{11} & \cdots & h_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ h_{k0} & h_{k1} & \cdots & h_{kk} \end{pmatrix}.$$

Då är $\hat{\beta}_i \sim N(\beta_i, \sigma^2 h_{ii})$. Om vi känner σ^2 kan vi använda att

$$Z = \frac{\hat{\beta}_i - \beta_i}{\sigma \sqrt{h_{ii}}} \sim N(0, 1)$$

för att testa hypotesen $H_0 : \beta_i = 0$ eller för att ställa upp konfidensintervall I_{β_i} .

Nu är det extremt sällan vi känner σ^2 exakt, men vi vet att $\hat{\beta}$ är oberoende av SS_E enligt huvudsatsen, så $\hat{\beta}_i$ är oberoende av SS_E . Vidare är $s^2 = SS_E/(n - k - 1)$ en skattning av σ^2 vi känner väl, så

$$\frac{(\hat{\beta}_i - \beta_i)/(\sigma \sqrt{h_{ii}})}{\sqrt{((n - k - 1)S^2/\sigma^2)/(n - k - 1)}} = \frac{\hat{\beta}_i - \beta_i}{S \sqrt{h_{ii}}} \sim t(n - k - 1)$$

enligt Gossets sats.

5.4 Hypotestest: $H_0 : \beta_i = 0$

Vi kan även utföra hypotestester för en ensklid koefficient β_i för att se om den förklaringsvariabeln tillför något signifikant (givetvis går det att ställa upp konfidensintervall också för mer kvantitativt innehåll).

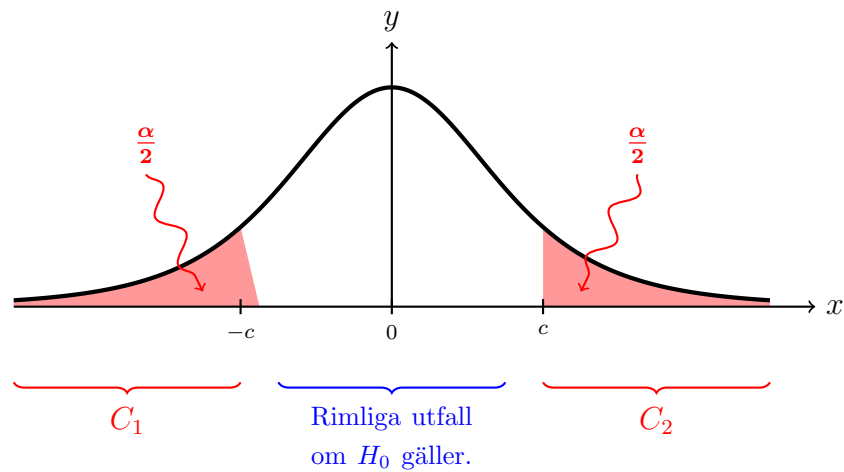
Låt $H_0 : \beta_i = 0$ och $H_1 : \beta_i \neq 0$. Vi vet enligt ovan att

$$T = \frac{\hat{\beta}_i - \beta_i}{S \sqrt{h_{ii}}} \sim t(n - k - 1),$$

så det kritiska området finner vi enligt

$$C = \{t \in \mathbf{R} : |t| > c\}$$

för lämpligt tal $c > 0$ beroende på signifikansnivån och antalet frihetsgrader.



Gränsen hittar vi i tabell genom att leta reda på ett tal $c = F_T^{-1}(1 - \alpha/2)$ (sitter du med MATLAB kan du använda `c = -tinv(alpha/2)`).

6 Exempel: enkel linjär regression

Vi diskuterade enkel linjär regression tidigare i samband med MK-skattningar (föreläsning 2). Låt oss visa att vi får samma resultat med den metod vi nu tagit fram (och fördelningar för ingående storheter på ett enkelt sätt).

Vi låter x_1, x_2, \dots, x_n vara fixerade tal och y_1, y_2, \dots, y_n vara observationer från

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i,$$

där $\epsilon_i \sim N(0, \sigma^2)$ är oberoende. Alltså precis den modell vi använt tidigare. Syntesmatrisen X ges av

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$$

så

$$X^T X = \begin{pmatrix} n & n\bar{x} \\ n\bar{x} & \sum_{i=1}^n x_i^2 \end{pmatrix} \Rightarrow (X^T X)^{-1} = \frac{1}{n \sum_{i=1}^n x_i^2 - (n\bar{x})^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{pmatrix}$$

om $\det(X^T X) \neq 0$. Således blir

$$\begin{aligned} \hat{\beta} &= \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = (X^T X)^{-1} X^T \mathbf{y} \\ &= \frac{1}{n \sum_{i=1}^n x_i^2 - (n\bar{x})^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{pmatrix} \begin{pmatrix} n\bar{y} \\ \sum_{i=1}^n x_i y_i \end{pmatrix} \\ &= \frac{1}{n \sum_{i=1}^n x_i^2 - (n\bar{x})^2} \begin{pmatrix} n\bar{y} \sum_{i=1}^n x_i^2 - n\bar{x} \sum_{i=1}^n x_i y_i \\ -n^2 \bar{x} \bar{y} + n \sum_{i=1}^n x_i y_i \end{pmatrix}, \end{aligned}$$

vilket ger att

$$\hat{\beta}_1 = \frac{\sum_{j=1}^n x_j y_j - n\bar{x}\bar{y}}{\sum_{j=1}^n x_j^2 - n\bar{x}^2} = \frac{\sum_{j=1}^n y_j(x_j - \bar{x})}{\sum_{j=1}^n (x_j - \bar{x})^2} = \frac{\sum_{j=1}^n (y_j - \bar{y})(x_j - \bar{x})}{\sum_{j=1}^n (x_j - \bar{x})^2} \quad \text{och} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

Det följer nu att

$$\hat{\beta}_0 \sim N\left(\beta_0, \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}\right) \quad \text{och} \quad \hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right).$$

Någonstans nu inser vi vilket kraftigt syntaktiskt verktyg matriser och linjär algebra är.. Eftersom σ är okänd skattar vi σ^2 med

$$s^2 = \frac{\text{SS}_E}{n-2} = \frac{\sum_{j=1}^n (y_j - \hat{\beta}_0 - \hat{\beta}_1 x_j)^2}{n-2}$$

där $S^2 \sim \chi^2(n-2)$. Vi kan skriva om SS_E genom att utnyttja att

$$\begin{aligned} \text{SS}_R &= \sum_{j=1}^n (\hat{\mu}_j - \bar{y})^2 = \hat{\beta}_1^2 \sum_{j=1}^n (x_j - \bar{x})^2 = \left(\frac{\sum_{j=1}^n (y_j - \bar{y})(x_j - \bar{x})}{\sum_{j=1}^n (x_j - \bar{x})^2} \right)^2 \sum_{j=1}^n (x_j - \bar{x})^2 \\ &= r^2 \sum_{j=1}^n (y_j - \bar{y})^2 \end{aligned}$$

så

$$\text{SS}_E = \text{SS}_{\text{TOT}} - \text{SS}_R = (1 - r^2) \sum_{j=1}^n (y_j - \bar{y})^2.$$

Vi kan här se att $r = 1$ ger perfekt matching så alla (x_j, y_j) ligger på en rät linje.

7 Bonus: determinanter och spår

För en kvadratisk matris A med dimension $n \times n$ definierar vi som bekant det **karaktéristiska polynomet** enligt

$$p(\lambda) = \det(A - \lambda I).$$

Bekant från linjär algebra är att $p(\lambda)$ kan representeras enligt

$$\begin{aligned} p(\lambda) &= (-1)^n (\lambda - \lambda_1)(\lambda - \lambda_2) \cdots (\lambda - \lambda_n) \\ &= (-1)^n (\lambda^n - \lambda^{n-1}(\lambda_1 + \lambda_2 + \cdots + \lambda_n) + \cdots + (-1)^n \lambda_1 \lambda_2 \cdots \lambda_n), \end{aligned} \quad (1)$$

där λ_i , $i = 1, 2, \dots, n$, är matrisens egenvärden.

Spåret för A betecknar $\text{tr } A$ och definieras som summan av diagonalelementen:

$$\text{tr } A = \sum_{i=1}^n a_{ii}.$$

Spåret är linjär och uppfyller alltså att

$$\text{tr}(cA) = c \text{tr } A \quad \text{och} \quad \text{tr}(A + B) = \text{tr } A + \text{tr } B.$$

Dessutom gäller att

$$\operatorname{tr}(AB) = \sum_{i=1}^n \sum_{j=1}^m a_{ij} b_{ji} = \sum_{j=1}^m \sum_{i=1}^n b_{ji} a_{ij} = \operatorname{tr}(BA)$$

om A är $n \times m$ och B är $m \times n$. Denna likhet leder direkt till att

$$\operatorname{tr}(PAP^{-1}) = \operatorname{tr}(P^{-1}(PA)) = \operatorname{tr} A,$$

vilket innebär att spåret (likt determinanten) är en invariant.

Vi vet även att $\det A = \lambda_1 \lambda_2 \cdots \lambda_n$, dvs determinanten ges av produkten av alla egenvärden. Finns något liknande för spåret? Svaret är vad man kanske skulle kunna gissa, nämligen att

$$\operatorname{tr} A = \lambda_1 + \lambda_2 + \cdots + \lambda_n.$$

Enklaste beviset är att utnyttja att alla matriser A kan skrivas på Jordans normalform, dvs att det finns en basbytesmatris så att $J = PAP^{-1}$, där J har $\lambda_1, \lambda_2, \dots, \lambda_n$ på diagonalen och är nästan en diagonalmatris. Utnyttjar vi att spåret är invariant är vi klara. Kanske ligger detta argument lite utanför grundkursen i linjär algebra, så låt oss studera det karakteristiska polynomet lite närmare som alternativ. Genom att utveckla efter exempelvis rad 1 kan determinanten skrivas

$$p(\lambda) = (a_{11} - \lambda) \begin{vmatrix} a_{22} - \lambda & a_{23} & \cdots & a_{2n} \\ a_{32} & a_{33} - \lambda & \cdots & a_{3n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n2} & a_{n3} & \cdots & a_{nn} - \lambda \end{vmatrix} + p_1(\lambda),$$

där $\operatorname{grad} p_1(\lambda) \leq n - 2$ eftersom vi tagit bort en λ -term. Om vi utför samma operation på denna determinant ser vi att

$$p(\lambda) = (a_{11} - \lambda)(a_{22} - \lambda) \begin{vmatrix} a_{33} - \lambda & a_{34} & \cdots & a_{3n} \\ a_{43} & a_{44} - \lambda & \cdots & a_{4n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n3} & a_{n4} & \cdots & a_{nn} - \lambda \end{vmatrix} + p_2(\lambda),$$

där $\operatorname{grad} p_2(\lambda) \leq n - 2$ med något (nytt) polynom $p_2(\lambda)$.

Upprepa proceduren n gånger och vi erhåller att

$$\begin{aligned} p(\lambda) &= (a_{11} - \lambda)(a_{22} - \lambda) \cdots (a_{nn} - \lambda) + p_n(\lambda) \\ &= (-1)^n \lambda^n + (-1)^{n-1} \lambda^{n-1} (a_{11} + a_{22} + \cdots + a_{nn}) + p_{n+1}(\lambda) \\ &= (-1)^n (\lambda^n - \lambda^{n-1} \operatorname{tr} A) + p_{n+1}(\lambda), \end{aligned} \tag{2}$$

där $p_{n+1}(\lambda)$ är något polynom med grad högst $n - 2$. Om vi jämför (1) med (2) ser vi att

$$\operatorname{tr} A = \lambda_1 + \lambda_2 + \cdots + \lambda_n.$$